

Matematyka stosowana

Symulacje stochastyczne i metody Monte Carlo

Wojciech NIEMIRO

`wniem@mimuw.edu.pl`

<http://adres.strony1.www>

Uniwersytet Warszawski, 2013



Streszczenie. Jest to wykład na temat symulacji zjawisk losowych. Obejmuje też wstęp do metod Monte Carlo (MC), czyli algorytmów zrandomizowanych. Symulacje komputerowe są nie tylko prostym i skutecznym narzędziem badania procesów losowych, ale też znajdują zastosowania do obliczeń wielkości deterministycznych. Metody Monte Carlo a w szczególności algorytmy MCMC, oparte na łańcuchach Markowa należą do standardowych narzędzi obliczeniowych, między innymi w statystyce Bayesowskiej i fizyce statystycznej. Przedmiot jest przeznaczony dla wszystkich studentów lubiących rachunek prawdopodobieństwa, matematyków i informatyków. Powinien zachęcić słuchaczy do „zobaczenia losowości” – przy wykorzystaniu pięknego i dostępnego za darmo pakietu R. Wykład jest utrzymany na elementarnym poziomie.

W pierwszej części wykładu omówione zostaną sposoby generowania zmiennych losowych o zadanym rozkładzie prawdopodobieństwa i prostych procesów stochastycznych. Druga część poświęcona będzie ogólnym zasadom konstrukcji algorytmów Monte Carlo, szacowania ich dokładności i redukcji błędów. Sporo miejsca zajmują Markowskie algorytmy Monte Carlo, MCMC.

Wersja internetowa wykładu:

<http://mst.mimuw.edu.pl/lecture.php?lecture=sst>

(może zawierać dodatkowe materiały)



Niniejsze materiały są dostępne na [licencji Creative Commons 3.0 Polska](#):
Uznanie autorstwa — Użycie niekomercyjne — Bez utworów zależnych.

Copyright © W.Niemirow, Uniwersytet Warszawski, Wydział Matematyki, Informatyki i Mechaniki, 2013. Niniejszy plik PDF został utworzony 21 lutego 2013.



Projekt współfinansowany przez Unię Europejską w ramach Europejskiego Funduszu Społecznego.



Skład w systemie L^AT_EX, z wykorzystaniem m.in. pakietów beamer oraz listings. Szablony podręcznika i prezentacji: Piotr Krzyżanowski; koncept: Robert Dąbrowski.

Spis treści

1. Wprowadzenie	5
1.1. Od autora	5
1.2. Przykłady	5
2. Podstawy R i ćwiczenia komputerowe	11
2.1. Początki	11
2.2. Ćwiczenia	14
3. Generowanie zmiennych losowych I. Ogólne metody	15
3.1. Przykłady	15
3.2. Metoda przekształceń	18
3.2.1. Odwrócenie dystrybuanty	18
3.3. Metoda eliminacji	20
3.3.1. Ogólny algorytm	20
3.3.2. Eliminacja w R	22
3.4. Metoda kompozycji	23
4. Generowanie zmiennych losowych II. Specjalne metody	24
4.1. Rozkłady dyskretne	24
4.2. Schematy kombinatoryczne	25
4.2.1. Pobieranie próbki bez zwracania	25
4.2.2. Permutacje losowe	26
4.3. Specjalne metody eliminacji	27
4.3.1. Iloraz zmiennych równomiernych	27
4.3.2. Gęstości przedstawione szeregami	29
5. Generowanie zmiennych losowych III. Rozkłady wielowymiarowe	31
5.1. Ogólne metody	31
5.1.1. Metoda rozkładów warunkowych	31
5.1.2. Metoda przekształceń	32
5.2. Kilka ważnych przykładów	33
5.2.1. Rozkłady sferyczne i eliptyczne	33
5.2.2. Rozkłady Dirichleta	36
6. Symulowanie procesów stochastycznych I.	40
6.1. Stacjonarne procesy Gaussowskie	40
6.2. Procesy Poissona	42
6.2.1. Jednorodny proces Poissona na półprostej	42
6.2.2. Niejednorodne procesy Poissona w przestrzeni	45
7. Symulowanie procesów stochastycznych II. Procesy Markowa	49
7.1. Czas dyskretny, przestrzeń dyskretna	49
7.2. Czas dyskretny, przestrzeń ciągła	50
7.3. Czas ciągły, przestrzeń dyskretna	51
8. Algorytmy Monte Carlo I. Obliczanie całek	55
8.1. Losowanie istotne	56
8.2. Efektywność estymatorów MC	56
8.3. Wazona eliminacja	59
9. Algorytmy Monte Carlo II. Redukcja wariancji	65

9.1. Losowanie warstwowe	65
9.2. Zmienne kontrolne	67
9.3. Zmienne antytetyczne	68
10. Markowskie Monte Carlo I. Wprowadzenie	71
10.1. Co to jest MCMC ?	71
10.2. Łańcuchy Markowa	71
10.2.1. Rozkład stacjonarny	72
10.2.2. Twierdzenia graniczne dla łańcuchów Markowa	73
11. Markowskie Monte Carlo II. Podstawowe algorytmy	76
11.1. Odwracalność	76
11.2. Algorytm Metropolis-Hastingsa	76
11.3. Próbnik Gibbsa	78
12. Markowskie Monte Carlo III. Przykłady zastosowań	82
12.1. Statystyka bayesowska	82
12.1.1. Hierarchiczny model klasyfikacji	82
12.1.2. Próbnik Gibbsa w modelu hierarchicznym	84
12.2. Estymatory największej wiarygodności	86
12.2.1. Model auto-logistyczny	86
13. Markowskie Monte Carlo IV. Pola losowe	89
13.1. Definicje	89
13.2. Generowanie markowskich pól losowych	90
13.3. Rekonstrukcja obrazów	91
14. Markowskie Monte Carlo V. Elementarna teoria łańcuchów Markowa	94
14.1. Podstawowe określenia i oznaczenia	94
14.2. Regeneracja	95
14.3. Łańcuchy sprzężone i zbieżność rozkładów	102
14.3.1. Odległość pełnego wahan	102
14.3.2. Sprzęganie	103
15. Markowskie Monte Carlo VI. Oszacowania dokładności	106
15.1. Reprezentacja spektralna macierzy odwracalnej	106
15.1.1. Oszacowanie szybkości zbieżności	108
15.1.2. Oszacowanie normy pełnego wahan	108
15.1.3. Oszacowanie obciążenia estymatora	109
15.2. Oszacowanie błędu średniokwadratowego estymatora	110
15.2.1. Asymptotyczna wariancja	110
15.2.2. Oszacowanie BSK	111

1. Wprowadzenie

1.1. Od autora

Jest kilka ważnych powodów, dla których warto się zająć symulacjami stochastycznymi:

- Symulacje stochastyczne są prostym sposobem badania zjawisk losowych.
- Ściśle związane z symulacjami stochastycznymi są metody obliczeniowe nazywane „Monte Carlo” (MC). Polegają one na wykorzystaniu „sztucznie generowanej” losowości w celu rozwiązania zadań deterministycznych. Metody MC są proste i skuteczne. Dla pewnych problemów MC jest jedynym dostępnym narzędziem obliczeniowym. Dla innych problemów MC jest co prawda mniej efektywne od metod numerycznych, ale za to dużo łatwiejsze!
- W moim przekonaniu symulacje stochastyczne są wspaniałą pomocą przy nauce rachunku prawdopodobieństwa. Pozwalają lepiej „zrozumieć losowość”.
- Symulacje stochastyczne są dostępne dla każdego. W szczególności, „otoczenie” **R**, które stanowi naprawdę potężne narzędzie, jest rozpowszechniane **za darmo!**

Jest wreszcie powód najważniejszy:

- Symulacje stochastyczne są **świetną zabawą!**

Literatura na temat symulacji stochastycznych jest bardzo obszerna. Godna polecenia jest książka Zielińskiego i Wieczorkowskiego [23], poświęcona w całości generatorom zmiennych losowych. Przedstawia ona bardziej szczegółowo zagadnienia, odpowiadające Rozdziałom 2–4 niniejszego skryptu i zawiera materiał, który zdecydowałem się pominąć: wytwarzanie „liczb losowych” o rozkładzie jednostajnym i testowanie generatorów. Podobne zagadnienia są przedstawione trochę w innym stylu w monografii Ripleya [18], która również zawiera wstęp do metod Monte Carlo. Zaawansowane wykłady można znaleźć w nowoczesnych monografiach Asmussena i Glynn’a [2], Liu [15], Roberta i Caselli [19]. Pierwsza z nich jest zorientowana bardziej na wyniki teoretyczne, zaś druga bardziej na zastosowania. Świetnym wstępem do metod MCMC są prace Geyera [7] i [8]. Teoria łańcuchów Markowa z uwzględnieniem zagadnień istotnych dla MCMC jest przystępnie przedstawiona w książce Brémaud [4]. Podstawy teoretycznej analizy zrandomizowanych algorytmów (tematy poruszane w Rozdziale 15 skryptu) są znakomicie przedstawione w pracach Jerruma i Sinclaire’a [11] oraz Jerruma [12].

1.2. Przykłady

Zacznę od kilku przykładów zadań obliczeniowych, które można rozwiązywać symulując losowość. Wybrałem przykłady najprostsze, do zrozumienia których wystarcza zdrowy rozsądek i nie potrzeba wielkiej wiedzy.

Przykład 1.1 (Igła Buffona, 1777). Podłoga jest nieskończoną płaszczyzną, podzieloną równoległymi prostymi na „deski” szerokości d . Rzucamy „losowo” igłę o długości l . Elementarne rozważania prowadzą do wniosku, że

$$p = \mathbb{P}(\text{igła przetnie którąś z prostych}) = \frac{2l}{\pi d}.$$

Buffon zauważył, że tę prostą obserwację można wykorzystać do... obliczania liczby π metodą „statystyczną”. Powtórzmy nasze doświadczenie niezależnie n razy. Oczywiście, mamy do czynienia ze schematem Bernoulliego, w którym „sukcesem” jest przecięcie prostej. Niech

$$\begin{aligned}\hat{p}_n &= \frac{\text{liczba sukcesów}}{\text{liczba doświadczeń}} \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{I}(\text{w } i\text{-tym doświadczeniu igła przecięła prosta})\end{aligned}$$

Wielkość \hat{p}_n jest empirycznym odpowiednikiem prawdopodobieństwa p i powinna przybliżać to prawdopodobieństwo, przynajmniej dla dużych n . W takim razie, możemy za przybliżenie liczby π przyjąć

$$\pi_n = \frac{2l}{\hat{p}_n d}.$$

Statystyk powiedziałby, że *zmienna losowa* π_n jest *estymatorem* liczby π . To wszystko jest bardzo proste, ale parę kwestii wymaga uściślenia. Jak duża ma być liczba powtórzeń, żeby przybliżenie było odpowiednio dokładne? Ale przecież mamy do czynienia ze „ślepyim losem”! Czyż pech nie może sprawić, że mimo dużej liczby doświadczeń przybliżenie jest kiepskie? Czy odpowiednio dobierając l i d możemy poprawić dokładność? A może da się zaprojektować lepsze doświadczenie?

Przykład 1.2 (Sieć zawodnych połączeń). Niech \mathcal{V}, \mathcal{E} będzie grafem skierowanym spójnym. Krawędzie reprezentują „połączenia” pomiędzy wierzchołkami. Krawędź $e \in \mathcal{E}$, niezależnie od pozostałych, ulega awarii ze znanym prawdopodobieństwem p_e . W rezultacie powstaje losowy podzbiór $C \subseteq \mathcal{E}$ sprawnych połączeń o rozkładzie prawdopodobieństwa

$$P(C) = \prod_{e \notin C} p_e \prod_{e \in C} (1 - p_e).$$

Jest jasne, jak można symulować to zjawisko: dla każdej krawędzi e „losujemy” zmienną U_e o rozkładzie równomiernym na $[0, 1]$ i przyjmujemy

$$\begin{cases} e \notin C & \text{jeśli } U_e \leq p_e; \\ e \in C & \text{jeśli } U_e > p_e. \end{cases}$$

Powiedzmy, że interesuje nas możliwość znalezienia ścieżki (ciągu krawędzi) wiodącej z ustalonego wierzchołka v_0 do innego wierzchołka v_1 . Niech

$$\theta = \mathbb{P}(\text{w zbiorze } C \text{ istnieje ścieżka z } v_0 \text{ do } v_1).$$

Generujemy niezależnie n kopii C_1, \dots, C_n zbioru C . Nieznane prawdopodobieństwo θ przybliżamy przez odpowiednik próbkowy:

$$\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(\text{w zbiorze } C_i \text{ istnieje ścieżka z } v_0 \text{ do } v_1).$$

Przykład 1.3 (Skomplikowana całka). Załóżmy, że X_1, \dots, X_m są niezależnymi zmiennymi losowymi o jednakowym rozkładzie $N(0, 1)$. Niech

$$R = \max_{k=1}^m \sum_{i=1}^k X_i - \min_{k=1}^m \sum_{i=1}^k X_i.$$

Chcemy obliczyć dystrybuantę zmiennej losowej R ,

$$H(x) = \mathbb{P}(R \leq x).$$

Zauważmy, że z definicji,

$$H(x) = \int \cdots \int_{\Omega_m} (2\pi)^{-m/2} \exp \left[-\frac{1}{2} \sum_{i=1}^m x_i^2 \right] dx_1 \cdots dx_m,$$

gdzie $\Omega_m = \left\{ (x_1, \dots, x_m) : \max_{k=1}^m \sum_{i=1}^k x_i - \min_{k=1}^m \sum_{i=1}^k x_i \right\}$. W zasadzie, jest to więc zadanie obliczenia całki. Jednak skomplikowany kształt wielowymiarowego zbioru Ω_m powoduje, że zastosowanie standardowych metod numerycznych jest utrudnione. Z kolei symulowanie zmiennej losowej R jest bardzo łatwe, wprost z definicji. Można wygenerować wiele niezależnych zmiennych R_1, \dots, R_n o rozkładzie takim samym jak R . Wystarczy teraz policzyć ile spośród tych zmiennych jest $\leq x$:

$$\hat{H}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(R_i \leq x).$$

Schemat symulacyjnego obliczania prawdopodobieństwa jest taki sam jak w dwu poprzednich przykładach. Podobnie zresztą jak w tamtych przykładach, podstawowy algorytm można znacznie ulepszać, ale dyskusję na ten temat odłożymy na później.

Przykład 1.4 (Funkcja harmoniczna). Następujące zadanie jest dyskretnym odpowiednikiem sławnego zagadnienia Dirichleta. Niech \mathcal{D} będzie podzbiorem kraty całkowitoliczbowej \mathbb{Z}^2 . Oznaczmy przez $\partial\mathcal{D}$ brzeg tego zbioru, zdefiniowany następująco:

$$(x, y) \in \partial\mathcal{D} \equiv (x, y) \notin \mathcal{D} \text{ i } [(x+1, y) \in \mathcal{D} \text{ lub } (x-1, y) \in \mathcal{D} \\ \text{lub } (x, y+1) \in \mathcal{D} \text{ lub } (x, y-1) \in \mathcal{D}].$$

Powiemy, że $u : \mathcal{D} \cup \partial\mathcal{D} \rightarrow \mathbb{R}$ jest funkcją harmoniczną, jeśli dla każdego punktu $(x, y) \in \mathcal{D}$,

$$u(x, y) = \frac{1}{4} \sum_{(x', y') \in \partial(x, y)} u(x', y'),$$

gdzie sumowanie rozciąga się na 4 punkty (x', y') sąsiadujące z (x, y) , to znaczy $\partial(x, y) = \{(x+1, y), (x-1, y), (x, y+1), (x, y-1)\}$.

Mamy daną funkcję na brzegu: $\bar{u} : \partial\mathcal{D} \rightarrow \mathbb{R}$. Zadanie polega na skonstruowaniu jej rozszerzenia harmonicznego, to znaczy takiej funkcji harmonicznej $u : \mathcal{D} \cup \partial\mathcal{D} \rightarrow \mathbb{R}$, że $u(x, y) = \bar{u}(x, y)$ dla $(x, y) \in \partial\mathcal{D}$.

Wyobraźmy sobie błądzenie losowe po kracie, startujące w punkcie $(x, y) \in \mathcal{D}$. Formalnie jest to ciąg losowych punktów określonych następująco:

$$(X_0, Y_0) = (x, y); \\ (X_{k+1}, Y_{k+1}) = (X_k, Y_k) + (\xi_{k+1}, \eta_{k+1}),$$

gdzie (ξ_k, η_k) są niezależnymi wektorami losowymi o jednakowym rozkładzie prawdopodobieństwa:

$$\begin{aligned}\mathbb{P}((\xi_k, \eta_k) = (0, 1)) &= \mathbb{P}((\xi_k, \eta_k) = (0, -1)) \\ &= \mathbb{P}((\xi_k, \eta_k) = (1, 0)) = \mathbb{P}((\xi_k, \eta_k) = (-1, 0)) = \frac{1}{4}.\end{aligned}$$

Błądźmy tak długo, aż natrafimy na brzeg obszaru. Formalnie, określamy moment zatrzymania $T = \min\{k : (X_k, Y_k) \in \partial\mathcal{D}\}$. Łatwo zauważyć, że funkcja

$$u(x, y) := \mathbb{E}[\bar{u}(X_T, Y_T) | (X_0, Y_0) = (x, y)]$$

jest rozwiązaniem zagadnienia! Istotnie, ze wzoru na prawdopodobieństwo całkowite wynika, że

$$u(x, y) = \frac{1}{4} \sum_{(x', y') \in \partial(x, y)} \mathbb{E}[\bar{u}(X_T, Y_T) | (X_1, Y_1) = (x', y')].$$

Wystarczy teraz spostrzeżenie, że rozkład zmiennej losowej $u(X_T, Y_T)$ pod warunkiem $(X_1, Y_1) = (x', y')$ jest taki sam jak pod warunkiem $(X_0, Y_0) = (x', y')$, bo błądzenie „rozpoczyna się na nowo”.

Algorytm Monte Carlo obliczania $u(x, y)$ oparty na powyższym spostrzeżeniu został wynaleziony przez von Neumanna i wygląda następująco:

- Powtórz wielokrotnie, powiedzmy n razy, niezależnie doświadczenie:
„błądź startując startując z (x, y) aż do brzegu; oblicz $\bar{u}(X_T, Y_T)$ ”
- Uśrednij wyniki n doświadczeń.

Dla bardziej formalnego zapisu algorytmu będę się posługiwał pseudo-kodem, który wydaje się zrozumiały bez dodatkowych objaśnień:

Listing.

```
{ 'Gen' oznacza 'Generuj' }
U := 0;
for j = 1 to n
  begin
    (X, Y) := (x, y);
    while (X, Y) ∈ D
      begin
        Gen (ξ, η) ~ U{(0, 1), (0, -1), (1, 0), (-1, 0)}; [ rozkład jednostajny na zbiorze 4-punktowym ]
        (X, Y) := (X, Y) + (ξ, η)
      end
      U := U +  $\bar{u}(X, Y)$ 
    end
  end
U := U/n
```

Przykład 1.5 (Problem „plecakowy”). Załóżmy, że $a = (a_1, \dots, a_m)^\top$ jest wektorem o współrzędnych naturalnych ($a_i \in \{1, 2, \dots\}$) i $b \in \{1, 2, \dots\}$. Rozważamy wektory $x = (x_1, \dots, x_m)^\top$ o współrzędnych zero-jedynkowych ($x_i \in \{0, 1\}$). Interesuje nas *liczba rozwiązań nierówności*

$$x^\top a = \sum_{i=1}^m x_i a_i \leq b,$$

a więc liczność $|\mathcal{X}(b)|$ zbioru $\mathcal{X}(b) = \{x \in \{0, 1\}^m : x^\top a \leq b\}$. Czytelnik domyśla się z pewnością, skąd nazwa problemu. Dokładne obliczenie liczby rozwiązań jest trudne. Można spróbować zastosować prostą metodę Monte Carlo. Jeśli zmienna losowa X ma rozkład jednostajny na przestrzeni $\{0, 1\}^m$, to $\mathbb{P}(X \in \mathcal{X}(b)) = |\mathcal{X}(b)|/2^m$. Wystarczy zatem oszacować to prawdopodobieństwo tak jak w trzech pierwszych przykładach w tym rozdziale. Generowanie zmiennej losowej o rozkładzie jednostajnym na przestrzeni $\{0, 1\}^m$ sprowadza się do przeprowadzenia m rzutów monetą i jest dziecinnie proste. Na czym więc polega problem? Otóż szacowane prawdopodobieństwo może być astronomicznie małe. Dla, powiedzmy $a = (1, \dots, 1)^\top$ i $b = m/3$, to prawdopodobieństwo jest $\leq e^{-m/18}$ (proszę się zastanowić jak uzasadnić tę nierówność). Przeprowadzanie ciągu doświadczeń Bernoulliego z takim prawdopodobieństwem sukcesu jest bardzo nieefektywne – na pierwszy sukces oczekiwać będziemy średnio $\geq e^{m/18}$, co dla dużych m jest po prostu katastrofalne.

Metoda, którą naszkicuję należy do rodziny algorytmów MCMC (Monte Carlo opartych na łańcuchach Markowa). Algorytm jest raczej skomplikowany, ale o ile mi wiadomo jest najefektywniejszym ze znanych sposobów rozwiązania zadania. Bez straty ogólności możemy przyjąć, że $a_1 \leq \dots \leq a_m$. Niech $b_0 = 0$ oraz $b_j = \min\{b, \sum_{i=1}^j a_i\}$. Rozważmy *ciąg* zadań plecakowych ze zmniejszającą się prawą stroną nierówności, równą kolejno $b = b_m, b_{m-1}, \dots, b_1, b_0 = 0$. Niech więc $\mathcal{X}(b_j) = \{x \in \{0, 1\}^m : x^\top a \leq b_j\}$. Zachodzi następujący „wzór teleskopowy”:

$$|\mathcal{X}(b)| = |\mathcal{X}(b_m)| = \frac{|\mathcal{X}(b_m)|}{|\mathcal{X}(b_{m-1})|} \cdot \frac{|\mathcal{X}(b_{m-1})|}{|\mathcal{X}(b_{m-2})|} \dots \frac{|\mathcal{X}(b_1)|}{|\mathcal{X}(b_0)|} \cdot |\mathcal{X}(b_0)|.$$

Oczywiście, $|\mathcal{X}(b_0)| = 1$, a więc możemy uznać, że zadanie sprowadza się do obliczenia ilorazów $|\mathcal{X}(b_{j-1})|/|\mathcal{X}(b_j)|$ (pomijamy w tym miejscu subtelności związane z dokładnością obliczeń). Gdybyśmy umieli efektywnie generować zmienne losowe o rozkładzie jednostajnym na przestrzeni $X(b_j)$, to moglibyśmy postępować w dobrze już znany sposób: liczyć „sukcesy” polegające na wpadnięciu w zbiór $X(b_{j-1}) \subseteq X(b_j)$. Rzecz jasna, możemy losować z rozkładu jednostajnego na kostce $\{0, 1\}^m$ i eliminować punkty poza zbiorem $X(b_j)$, ale w ten sposób wpadlibyśmy w pułapkę, od której właśnie uciekamy: co zrobić jeśli przez sito eliminacji przechodzi jedno na $\geq e^{m/18}$ losowań?

Opiszemy pewne wyjście, polegające na zastosowaniu błędzenia losowego po zbiorze $X(b_j)$. Dla ustalenia uwagi przyjmijmy $j = m$, czyli $b_j = b$. Losowy ciąg punktów $X(0), X(1), \dots, X(n), \dots$ generujemy rekurencyjnie w taki sposób:

Listing.

```

X(0) := 0;
for n = 0 to ∞
  begin
    X := X(n); [ gdzie X = (X1, ..., Xm)
    Gen I ~ U{1, ..., m}; [ losujemy indeks do zmiany ]
    Y := X; YI = 1 - XI;
    if Y⊤a ≤ b then X(n+1) := Y
      else X(n+1) := X
    end
  end
end

```

Zaczynamy błędzenie w punkcie $0 = (0, \dots, 0)^\top$. Jeśli w chwili n jesteśmy w punkcie X , to próbujemy przejść do nowego punktu Y , utworzonego przez zmianę jednej, losowo wybranej współrzędnej punktu X (0 zamieniamy na 1 lub z 1 na 0; pozostałe współrzędne zostawiamy bez zmian). Jeśli „proponowany punkt Y nie wypadł” z rozważanej przestrzeni $\mathcal{X}(b)$, to przechodzimy do punktu Y . W przeciwnym wypadku stoimy w punkcie X .

Rzecz jasna, generowane w ten sposób zmienne losowe $X(n)$ nie mają dokładnie rozkładu jednostajnego ani tym bardziej nie są niezależne. Jednak dla dużych n zmienna $X(n)$ *ma w przybliżeniu* rozkład jednostajny:

$$\mathbb{P}(X(n) = x) \rightarrow \frac{1}{|\mathcal{X}(b)|} \quad (n \rightarrow \infty)$$

dla każdego $x \in \mathcal{X}(b)$. Co więcej, z prawdopodobieństwem 1 jest prawdą, że

$$\frac{1}{n} \sum_{i=1}^n \mathbb{I}(X(i) = x) \rightarrow \frac{1}{|\mathcal{X}(b)|} \quad (n \rightarrow \infty).$$

Innymi słowy, ciąg $X(n)$ spełnia prawo wielkich liczb – i może być użyty do szacowania liczności podzbiorów przestrzeni $\mathcal{X}(b)$ *zamiast* trudnego do symulowania ciągu niezależnych zmiennych o rozkładzie jednostajnym.

2. Podstawy R i ćwiczenia komputerowe

Oczekuję, że Czytelnik uruchomił pakiet (lub „otoczenie”, *environment*) **R**. Nie zakładam żadnej wstępnej znajomości tego języka programowania. Oczywiście, podstaw programowania w R można się nauczyć z odpowiednich źródeł. Ale prawdę mówiąc, prostych rzeczy można się szybko domyślić i zacząć zabawę natychmiast.

2.1. Początki

Przykład 2.1 (Rozkład prawdopodobieństwa i próbka losowa). Wylosujmy „próbkę” średniego rozmiaru, powiedzmy $n = 100$ z rozkładu normalnego $N(0, 1)$:

```
> n <- 100
> X <- rnorm(n)
> X
```

(aby przyjrzeć się funkcji `rnorm`, napiszmy `?rnorm`). Wektor `X` zawiera „realizacje” niezależnych zmiennych losowych X_1, \dots, X_{100} . Możemy obliczyć średnią $\bar{X} = n^{-1} \sum_{i=1}^n X_i$ i wariancję próbkową $\bar{X} = (n - 1)^{-1} \sum_{i=1}^n (X_i - \bar{X})^2$. Jakich wyników oczekujemy? Zobaczmy:

```
> mean(X)
> var(X)
```

Porównajmy kwantyl empiryczny rzędu, powiedzmy $p = 1/3$, z teoretycznym:

```
> quantile(X,p)
> qnorm(p)
```

Możemy powtórzyć nasze „doświadczenie losowe” i zobaczyć jakie są losowe fluktuacje wyników. Warto wykonać nasz mini-programik jeszcze raz, albo parę razy. Teraz spróbujmy „zobaczyć” próbkę. Najlepiej tak:

```
> hist(X,prob=TRUE) ( proszę się dowiedzieć co znaczy parametr 'prob' ! )
> rug(X)
```

Możemy łatwo narysować wykres gęstości. Funkcja obliczająca gęstość φ rozkładu $N(0, 1)$ nazywa się `dnorm`, a `curve` jest funkcją rysującą wykresy.

```
> curve(dnorm(x),col="blue",add=TRUE)
```

(nawiasem mówiąc, zachęcam *początkujących* pRobabilistów do oznaczania wygenerowanych wektorów losowych dużymi literami, np. `X`, a deterministycznych zmiennych i wektorów – małymi, np. `x`, podobnie jak na rachunku prawdopodobieństwa).

Podobnie, możemy porównać dystrybuantę empiryczną z prawdziwą dystrybuantą Φ . Funkcja obliczająca dystrybuantę empiryczną nazywa się `ecdf`, zaś dystrybuantę rozkładu $N(0, 1)$ – `pnorm`.

```
> plot(ecdf(X))
```

```
> curve(pnorm(x), from=xmin, to=xmax, col="blue", add=TRUE)
```

Test Kołmogorowa-Smirnowa oblicza maksymalną odległość $D = \sum_x |\hat{F}_n(x) - F(x)|$, gdzie \hat{F}_n jest dystrybucją empiryczną i $F = \Phi$.

```
> ks.test(X, pnorm, exact=TRUE)
```

Przypomnijmy sobie z wykładu ze statystyki, co znaczy podana przez test p -wartość. Jakiego wyniku „spodziewaliśmy się”?

Jest teraz dobra okazja, aby pokazać jak się symulacyjnie bada rozkład zmiennej losowej. Przypuśćmy, że interesuje nas rozkład prawdopodobieństwa p -wartości w przeprowadzonym powyżej doświadczeniu (polegającym na wylosowaniu próbki X i przeprowadzeniu testu KS). Tak naprawdę powinniśmy znać odpowiedź bez żadnych doświadczeń, ale możemy udawać niewiedzę i symulować.

Przykład 2.2 (Powtarzanie doświadczenia symulacyjnego). Symulacje polegają na powtórzeniu całego doświadczenia wiele razy, powiedzmy $m = 10000$ razy, zanotowaniu wyników i uznaniu powstałego rozkładu empirycznego za przybliżenie badanego rozkładu prawdopodobieństwa. Bardzo ważne jest zrozumienie różnej roli jaką pełni tu n (rozmiar próbki w pojedynczym doświadczeniu, a więc „parametr badanego zjawiska”) i m (liczba powtórzeń, która powinna być możliwie największa aby zwiększyć dokładność badania symulacyjnego). Następujący programik podkreśla logiczną konstrukcję powtarzanego doświadczenia:

```
> m <- 10000
> n <- 100
>
> # Przygotowujemy wektory w którym zapiszemy wyniki:
> D <- c()
> P <- c()
>
> for (i in 1:m)
> {
> X <- rnorm(n)
> Test <- ks.test(X, pnorm, exact=TRUE)
> D[i] <- Test$statistic
> P[i] <- Test$p.value
> } # koniec pętli for
>
> # Analizujemy wyniki:
> hist(D, prob=TRUE)
> hist(P, prob=TRUE)
```

Co prawda powyższy programik spełnia swoją rolę, jednak struktura pakietu R jest dostosowana do innego stylu pisania programów. Zamiast pętli `for` zalecane jest używanie funkcji które powtarzają pewne operacje i posługiwanie się, kiedykolwiek to możliwe, całymi wektorami a nie pojedynczymi komponentami. Poniższy fragment kodu zastępuje pętlę `for`

```
> DiP <- replicate(m, ks.test(rnorm(n), pnorm, exact=TRUE)[1:2])
>
> DiP <- t(DiP) # transpozycja macierzy ułatwia oglądanie wyników
> D <- as.numeric((DiP)[,1]) # pierwsza kolumna macierzy
> P <- as.numeric((DiP)[,2]) # druga kolumna macierzy
```

```
> # obiekty typu „list” przerobimy na wektory liczbowe:
> D <- as.numeric(D); P <- as.numeric(D)
```

Symulacje dają okazję „namacalnego” przedstawienia twierdzeń probabilistycznych (i nie tylko twierdzeń, ale także stwierdzeń, przypuszczeń prawdziwych lub fałszywych).

Przykład 2.3 (Mocne Prawo Wielkich Liczb). Wylosujmy próbkę z „jakiegoś” rozkładu prawdopodobieństwa. Weźmy na przykład niezależne zmienne o rozkładzie wykładniczym, $X_1, \dots, X_n, \dots \text{Ex}(2)$. Obliczmy średnie empiryczne

$$M_n = \frac{S_n}{n} = \frac{1}{n} \sum_{i=1}^n X_i.$$

Zróbmy wykres ciągu średnich $S_1/1, S_2/2, \dots, S_n/n, \dots$

```
> nmax <- 1000 # komputerowy odpowiednik „ $n \rightarrow \infty$ ”
> n <- (1:nmax)
> lambda <- 2 > X <- rexp(nmax,rate=lambda)
> S <- cumsum(X) # ciąg narastających sum
> M <- S/n # działania w R (np. dzielenie) są wykonywane „po współrzędnych”
> plot(n,M,type="l") # „zakłęcie” type="l" powoduje narysowanie łamanej
```

Teraz spóbuemy podobne doświadczenie zrobić dla zmiennych $X_i = 1/U_i - 1$, gdzie $U_i \sim U(0, 1)$ (X_i są próbką z tak zwanego rozkładu Pareto).

```
> # potrzebujemy dużej próbki, żeby się zorientować, co się dzieje...
> nmax <- 100000
> n <- (1:nmax)
> X <- 1/runif(nmax)-1
> M <- cumsum(X)/n
> plot(log(n),M,type="l") # i zrobimy wykres w skali logarytmicznej
```

Przykład 2.4 (Centralne Twierdzenie Graniczne). CTG jest jednym ze „słabych” twierdzeń granicznych rachunku prawdopodobieństwa, to znaczy dotyczy zbieżności rozkładów. Symulacyjne „sprawdzenie” lub ilustracja takich twierdzeń wymaga powtarzania doświadczenia wiele razy, podobnie jak w Przykładzie 2.2. Pojedyncze doświadczenie polega na obliczeniu sumy

$$S_n = \sum_{i=1}^n X_i,$$

gdzie X_1, \dots, X_n jest próbką z „jakiegoś” rozkładu prawdopodobieństwa i n jest „duże”. Weźmy na przykład niezależne zmienne o rozkładzie wykładniczym, jak w Przykładzie 2.3.

```
> m <- 10000
> n <- 100
> lambda <- 2
> S <- replicate(m, sum(rexp(n,rate=lambda)))
> hist(S,prob=TRUE)
> curve(dnorm(x,mean=n/lambda,sd=sqrt(n)/lambda),col="blue",add=TRUE)
> # wydaje się na podstawie obrazka, że dopasowanie jest znakomite
> ks.test(S,pnorm,mean=n/lambda,sd=sqrt(n)/lambda)
> # ale test Kołmogorowa-Smirnowa „widzi” pewne odchylenie od rozkładu normalnego
>
```

Chciałbym podkreślić raz jeszcze różnicę pomiędzy metodologią sprawdzania „mocnego” twierdzenia granicznego w Przykładzie 2.3 i „słabego” w Przykładzie 2.4. W pierwszym przypadku chcieliśmy zobrazować zbieżność ciągu *zmiennych losowych*, a w drugim – zbieżność ciągu *rozkładów prawdopodobieństwa*.

2.2. Ćwiczenia

Ćwiczenie 2.1. Wyjaśnić dlaczego w Przykładzie 2.4 nie musieliśmy (choć mogliśmy) rozpatrywać *unormowanych* zmiennych losowych

$$\frac{S_n - n\mu}{\sqrt{n}\sigma},$$

gdzie $\mu = \mathbb{E}X_1$ i $\sigma^2 = \text{Var} X_1$.

Ćwiczenie 2.2. W Przykładzie 2.4 znany jest *dokładny* rozkład prawdopodobieństwa sumy S_n . Co to za rozkład? Sprawdzić symulacyjnie zgodność z tym rozkładem.

Ćwiczenie 2.3. Przeprowadzić podobne doświadczenie jak w Przykładzie 2.4 (CTG), dla $n = 12$ i $X_i \sim U(0, 1)$. Skomentować wynik.

Ćwiczenie 2.4. Przeprowadzić podobne doświadczenie jak w Ćwiczeniu 2.3, ale zastępując sumę przez medianę dla $n = 13$ i $X_i \sim U(0, 1)$. Wypróbować przybliżenie normalne (jeśli teoretyczna wartość wariancji mediany nie jest znana, można zastąpić ją przez przybliżenie empiryczne). Skomentować wynik.

Ćwiczenie 2.5. W Ćwiczeniu 2.3, znany jest *dokładny* rozkład prawdopodobieństwa mediany. Co to za rozkład? Sprawdzić symulacyjnie zgodność z tym rozkładem. Można wziąć mniejsze, nieparzyste n . Dlaczego nieparzyste?

3. Generowanie zmiennych losowych I. Ogólne metody

3.1. Przykłady

Moje wykłady ograniczają się do zagadnień leżących w kompetencji rachunku prawdopodobieństwa i statystyki. U podstaw symulacji stochastycznych leży generowanie „liczb pseudo-losowych”, naśladujących zachowanie zmiennych losowych o rozkładzie jednostajnym. Jak to się robi, co to jest „pseudo-losowość”, czym się różni od „prawdziwej losowości”? To są fascynujące pytania, którymi zajmuje się: teoria liczb, teoria (chaotycznych) układów dynamicznych oraz filozofia. Dyskusja na ten temat przekracza ramy tych wykładów. Z punktu widzenia użytkownika, „liczby losowe” są bardzo łatwo dostępne, bo ich generatory są wbudowane w systemy komputerowe. Przyjmę pragmatyczny punkt widzenia i zacznę od następującego założenia.

Założenie 3.1. Mamy do dyspozycji potencjalnie nieskończony ciąg niezależnych zmiennych losowych U_1, \dots, U_n, \dots o jednakowym rozkładzie $U(0, 1)$.

W języku algorytmicznym: przyjmujemy, że każdorazowe wykonanie instrukcji zapisanej w pseudokodzie

Listing.

```
Gen  $U \sim U(0, 1)$ 
```

wygeneruje kolejną (nową) zmienną U_n . Innymi słowy, zostanie wykonane nowe, niezależne doświadczenie polegające na wylosowaniu przypadkowo wybranej liczby z przedziału $]0, 1[$.

Przykład 3.1. Wykonanie pseudokodu

Listing.

```
for  $i = 1$  to 10
  begin
    Gen  $U \sim U(0, 1)$ 
    write  $U$ 
  end
```

da, powiedzmy, taki efekt:

```
0.32240106    0.38971803    0.35222521    0.22550039    0.04162166
0.13976025    0.16943910    0.69482111    0.28812341    0.58138865
```

Nawiasem mówiąc, rzeczywisty kod w R, który wyprodukował nasze 10 liczb losowych był taki:

```
U <- runif(10); U
```

Język R jest zwięzły i piękny, ale nasz pseudokod ma pewne walory dydaktyczne.

Zajmiemy się teraz pytaniem, jak „wyprodukować” zmienne losowe o różnych rozkładach, wykorzystując zmienne U_1, U_2, \dots . W tym podrozdziale pokażę kilka przykładów, a w następnym przedstawię rzecz nieco bardziej systematycznie.

Przykład 3.2 (Rozkład Wykładniczy). To jest wyjątkowo łatwy do generowania rozkład – wystarczy taki algorytm:

Listing.

```
Gen  $U$ ;  $X := -\frac{1}{\lambda} \ln U$ 
```

Na wyjściu, $X \sim \text{Ex}(\lambda)$. Żeby się o tym przekonać, wystarczy obliczyć dystrybuantę tej zmiennej losowej: $\mathbb{P}(X \leq x) = \mathbb{P}(-\frac{1}{\lambda} \log U \leq x) = \mathbb{P}(U \geq e^{-\lambda x}) = 1 - e^{-\lambda x}$. Jest to najprostszy przykład ogólnej metody „odwracania dystrybuanty”, której poświęcę następny podrozdział.

Przykład 3.3 (Generacja rozkładu normalnego). Zmienna losowa

$$X = \sum_{i=1}^{12} U_i - 6$$

ma w przybliżeniu standardowy rozkład normalny $N(0, 1)$. Wynika to z Centralnego Twierdzenia Granicznego (jeśli uznamy, że liczba 12 jest dostatecznie bliska ∞ ; zauważmy, że $\mathbb{E}X = 0$ i $\text{Var} X = 1$). Oczywiście, w czasach szybkich komputerów ta przybliżona metoda zdecydowanie nie jest polecana. Jest natomiast pouczające zbadać (symulacyjnie!) jak dobre jest przybliżenie. Faktycznie bardzo trudno odróżnić próbkę X_1, \dots, X_n wyprodukowaną przez powyższy algorytm od próbki pochodzącej *dokładnie* z rozkładu $N(0, 1)$ (chyba, że n jest ogromne).

Przykład 3.4 (Algorytm Boxa-Müllera). Oto, dla porównania, bardziej współczesna – i całkiem dokładna metoda generowania zmiennych o rozkładzie normalnym.

Listing.

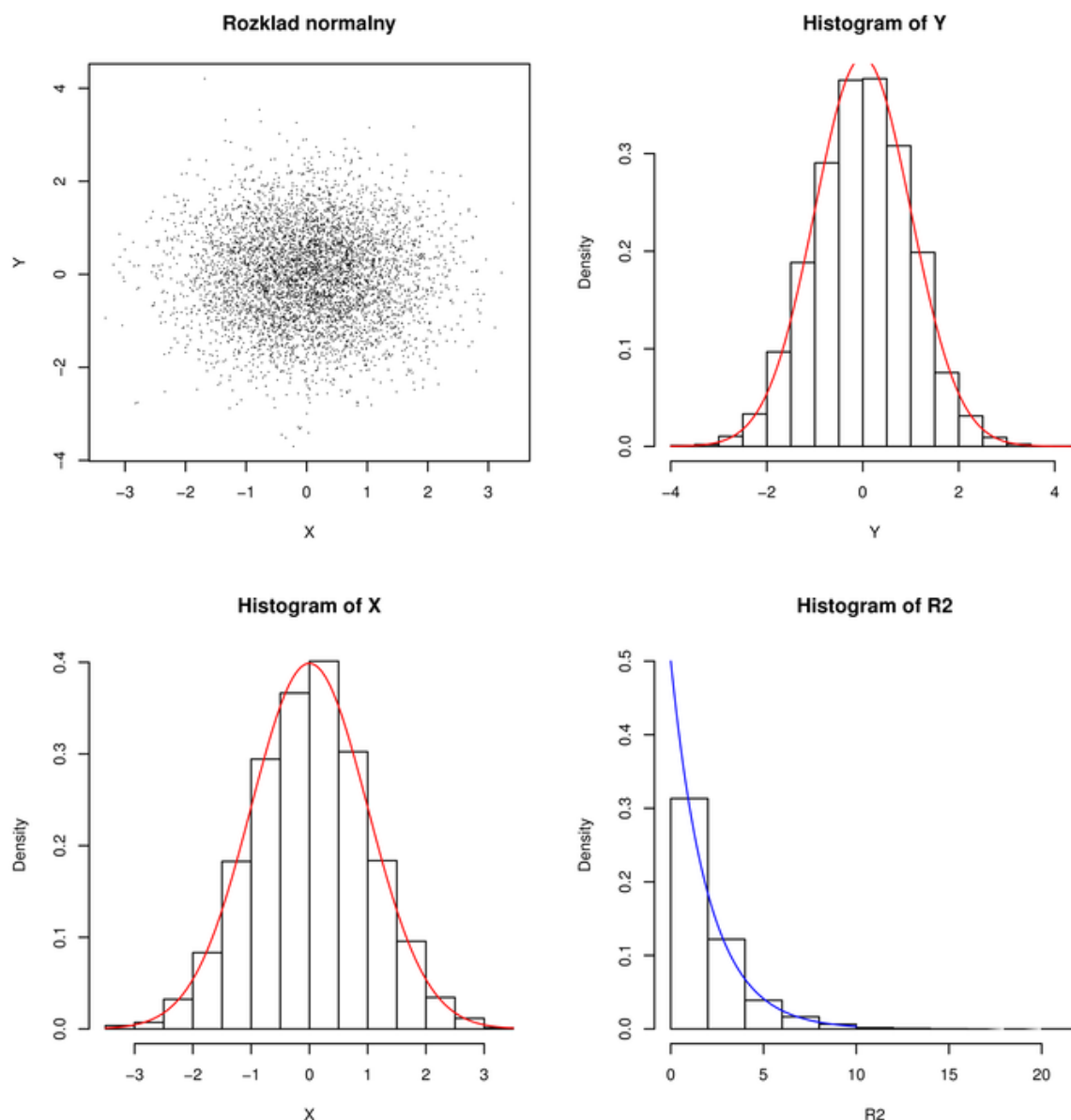
```
Gen  $U_1$ ;  $\Theta := 2\pi U_1$ ,  
Gen  $U_2$ ;  $R := \sqrt{-2 \ln U_2}$ ,  
Gen  $X := R \cos \Theta$ ;  $Y := R \sin \Theta$ 
```

Na wyjściu *obie* zmienne X i Y mają rozkład $N(0, 1)$ i w dodatku są niezależne. Uzasadnienie poprawności algorytmu Boxa-Müllera opiera się na dwu faktach: zmienna $R^2 = X^2 + Y^2$ ma rozkład $\chi^2(2) = \text{Ex}(1/2)$, zaś kąt Θ między osią i promieniem wodzącym punktu (X, Y) ma rozkład $U(0, 2\pi)$.

Ciekawe, że łatwiej jest generować zmienne losowe normalne „parami”.

Doświadczenie, polegające na wygenerowaniu zmiennych losowych X i Y powtórzyłem 10000 razy. Na Rysunku 3.1 widać 10000 wylosowanych w ten sam sposób i niezależnie punktów (X, Y) , histogramy i gęstości brzegowe X i Y (każda ze współrzędnych ma rozkład $N(0, 1)$) oraz histogram i gęstość $R^2 = X^2 + Y^2$ (rozkład wykładniczy $\text{Ex}(1/2)$).

Histogram jest „empirycznym” (może w obecnym kontekście należałoby powiedzieć „symulacyjnym”) odpowiednikiem gęstości: spośród wylosowanych wyników zliczane są punkty należące do poszczególnych przedziałów.



Rysunek 3.1. Dwuwymiarowy rozkład normalny i rozkłady brzegowe.

Przykład 3.5 (Rozkład Poissona). Algorytm przedstawiony niżej nie jest efektywny, ale jest prosty do zrozumienia.

Listing.

```

c := e-λ
Gen U; P := U; N := 0
while P ≥ c do
  begin Gen U; P := PU; N := N + 1 end

```

Na wyjściu $N \sim \text{Pois}(\lambda)$. Istotnie, niech E_1, \dots, E_n, \dots będą i.i.d. $\sim \text{Ex}(1)$ i $S_n = E_1 + \dots + E_n$. Wykorzystamy znany fakt, że $N = \max\{n : S_n \leq \lambda\}$ ma rozkład $\text{Pois}(\lambda)$. Zmienne o rozkładzie wykładniczym przedstawiamy jako $E_i = -\ln U_i$ – patrz Przykład 3.2. Zamiast

dodawac zmienne E_i możemy mnożyć zmienne U_i . Mamy $N = \max\{n : P_n \geq e^{-\lambda}\}$, gdzie $P_n = U_1 \cdots U_n$.

3.2. Metoda przekształceń

Zmienna losowa Y , która ma postać $Y = h(X)$, a więc jest pewną funkcją zmiennej X , w naturalny sposób „dziedziczy” rozkład prawdopodobieństwa zgodnie z ogólnym schematem $\mathbb{P}(Y \in \cdot) = \mathbb{P}(h(X) \in \cdot) = \mathbb{P}(X \in h^{-1}(\cdot))$ („wykropkowany” argument jest zbiorem). Przy tym zmienne X i Y nie muszą być jednowymiarowe. Jeśli obie zmienne mają *ten sam wymiar* i przekształcenie h jest dyfeomorfizmem, to dobrze znany wzór wyraża gęstość rozkładu Y przez gęstość X (Twierdzenie 5.1). Odpowiednio dobierając funkcję h możemy „przetwarzać” jedne rozkłady prawdopodobieństwa na inne, nowe.

Prawie wszystkie algorytmy generowania zmiennych losowych zawierają przekształcenia zmiennych losowych jako część składową. W niemal „czystej” postaci metoda przekształceń pojawiła się w algorytmie Boxa-Müllera, Przykładzie 3.4. Najważniejszym być może szczególnym przypadkiem metody przekształceń jest odwracanie dystrybucji.

3.2.1. Odwrócenie dystrybucji

Faktycznie, ta metoda została już wykorzystana w Przykładzie 3.2. Opiera się ona na prostym fakcie. Jeżeli F jest ciągłą i ściśle rosnącą dystrybucją, $U \sim U(0,1)$ i $X = F^{-1}(U)$, to $X = F^{-1}(U) \sim F$. Następująca definicja funkcji „pseudo-odwrotnej” pozwala pozbyć się kłopotliwych założeń.

Definicja 3.1. Jeżeli $F : \mathbb{R} \rightarrow [0,1]$ jest dowolną dystrybucją, to funkcję $F^- :]0,1[\rightarrow \mathbb{R}$ określamy wzorem:

$$F^-(u) = \inf\{x : F(x) \geq u\}.$$

Stwierdzenie 3.1. Nierówność $F^-(u) \leq x$ jest równoważna $u \leq F(x)$, dla dowolnych $u \in]0,1[$ i $x \in \mathbb{R}$.

Dowód. Z prawostronnej ciągłości dystrybucji F wynika, że kres dolny w Definicji 3.1 jest osiągnięty, czyli

$$F(F^-(u)) \geq u.$$

Z drugiej strony,

$$F^-(F(x)) = \min\{y : F(y) \geq F(x)\} \leq x,$$

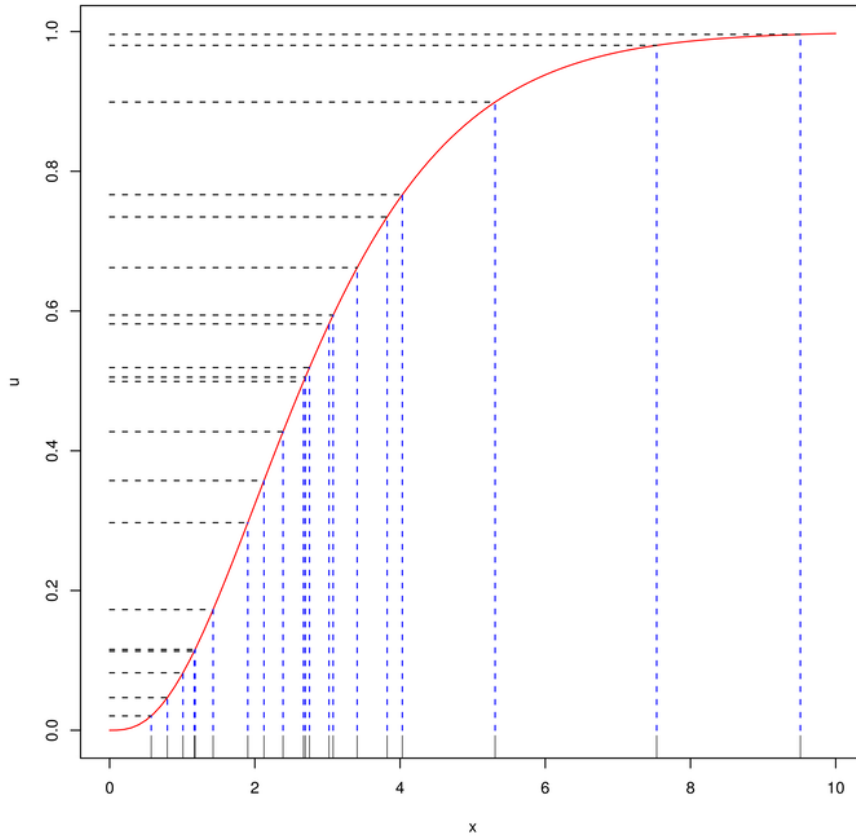
po prostu dlatego, że $x \in \{y : F(y) \geq F(x)\}$. Teza stwierdzenia natychmiast wynika z dwóch nierówności powyżej. \square

Wniosek 3.1 (Ogólna metoda odwrócenia dystrybucji). Jeżeli $U \sim U(0,1)$ i $X = F^-(U)$, to $\mathbb{P}(X \leq x) = F(x)$. W skrócie, $X \sim F$.

Na Rysunku 3.2 widać 20 punktów U_1, \dots, U_{20} , które wylosowałem z rozkładu $U(0,1)$ (na osi pionowej) i odpowiadające im punkty $X_i = F^{-1}(U_i)$ (na osi poziomej). W tym przypadku, F jest dystrybucją rozkładu Gamma(3,1) (linia krzywa). Najważniejszy fragment kodu w R jest taki:

```
curve(pgamma(x,shape=3,rate=1), from=0,to=10) # rysowanie F
U <- runif(20); X <- qgamma(U,shape=3,rate=1)
```

Zauważmy, że ta metoda działa również dla rozkładów dyskretnych i sprowadza się wtedy do metody „oczywistej”.



Rysunek 3.2. Odwracanie dystrybuanty.

Przykład 3.6 (Rozkłady dyskretne). Załóżmy, że $\mathbb{P}(X = i) = p_i$ dla $i = 1, 2, \dots$ i $\sum p_i = 1$. Niech $s_0 = 0$, $s_k = \sum_{i=1}^k p_i$. Jeżeli F jest dystrybuantą zmiennej losowej X , to

$$F^{-}(u) = i \quad \text{wtedy i tylko wtedy gdy} \quad s_{i-1} < u \leq s_i.$$

Odwracanie dystrybuanty ma ogromne znaczenie teoretyczne, bo jest całkowicie ogólną metodą generowania dowolnych zmiennych losowych jednowymiarowych. Może się to wydać dziwne, ale w praktyce ta metoda jest używana stosunkowo rzadko, z dwóch względów:

— Obliczanie F^{-} bywa trudne i nieefektywne.

— Stosowalność metody ogranicza się do zmiennych losowych jednowymiarowych.

Podam dwa przykłady, w których zastosowanie metody odwracania dystrybuanty jest rozsądne

Przykład 3.7 (Rozkład Weibulla). Z definicji, $X \sim \text{Weibull}(\beta)$, jeśli

$$F(x) = 1 - \exp(-x^\beta)$$

dla $x \geq 0$. Odwrócenie dystrybuanty i generacja X są łatwe:

$$X = (-\ln U)^{1/\beta}, \quad U \sim U(0, 1).$$

Przykład 3.8 (Rozkład Cauchy’ego). Gęstość i dystrybuanta zmiennej $X \sim \text{Cauchy}(0, 1)$ są następujące:

$$f(x) = \frac{1}{\pi} \frac{1}{1+x^2}, \quad F(x) = \frac{1}{2} + \frac{1}{\pi} \arctan(x).$$

Można tę zmienną generować korzystając z wzoru:

$$X = \tan \left(\pi \left(U - \frac{1}{2} \right) \right), \quad U \sim U(0, 1).$$

3.3. Metoda eliminacji

To jest najważniejsza, najczęściej stosowana i najbardziej uniwersalna metoda. Zacznę od raczej oczywistego faktu, który jest w istocie probabilistycznym sformułowaniem definicji prawdopodobieństwa warunkowego.

Stwierdzenie 3.2. *Przypuśćmy, że $Z = Z_1, \dots, Z_n, \dots$ jest ciągiem niezależnych zmiennych losowych o jednakowym rozkładzie, o wartościach w przestrzeni \mathcal{Z} . Niech $C \subseteq \mathcal{Z}$ będzie takim zbiorem, że $\mathbb{P}(Z \in C) > 0$. Niech*

$$N = \min\{n : Z_n \in C\}.$$

Zmienne losowe N i Z_N są niezależne, przy tym

$$\mathbb{P}(Z_N \in B) = \mathbb{P}(Z \in B | Z \in C) \quad \text{dla dowolnego } B \subseteq \mathcal{Z},$$

zaś

$$\mathbb{P}(N = n) = pq^{n-1}, \quad (n = 1, 2, \dots), \quad \text{gdzie } p = \mathbb{P}(Z \in C).$$

Dowód. Wystarczy zauważyć, że

$$\begin{aligned} \mathbb{P}(X_N \in B, N = n) &= \mathbb{P}(Z_1 \notin C, \dots, Z_{n-1} \notin C, Z_n \in C \cap B) \\ &= \mathbb{P}(Z_1 \notin C) \cdots \mathbb{P}(Z_{n-1} \notin C) \mathbb{P}(Z_n \in C \cap B) \\ &= q^{n-1} \mathbb{P}(Z \in C \cap B) = q^{n-1} \mathbb{P}(Z \in B | Z \in C) p. \end{aligned}$$

□

W tym Stwierdzeniu \mathcal{Z} może być dowolną przestrzenią mierzalną, zaś C i B – dowolnymi zbiorami mierzalnymi. Stwierdzenie mówi po prostu, że prawdopodobieństwo warunkowe odpowiada doświadczeniu losowemu *powtarzanemu aż do momentu spełnienia warunku*, przy czym rezultaty poprzednich doświadczeń się ignoruje (stąd nazwa: eliminacja).

3.3.1. Ogólny algorytm

Zakładamy, że umiemy generować zmienne losowe o gęstości g , a chcielibyśmy otrzymać zmienną o gęstości *proporcjonalnej* do funkcji f . Zakładamy, że $0 \leq f \leq g$.

Listing.

```
repeat
  Gen  $Y \sim g$ ;
  Gen  $U \sim U(0, 1)$ 
until  $U \leq \frac{f(Y)}{g(Y)}$ ;
 $X := Y$ 
```

Dowód poprawności algorytmu. Na mocy Stwierdzenia 3.2, zastosowanego do zmiennych losowych $Z = (Y, U)$ wystarczy pokazać, że

$$\mathbb{P}(Y \in B | U \leq f(Y)/g(Y)) = \frac{\int_B f(y)dy}{\int_{\mathcal{X}} f(y)dy}, \quad (3.1)$$

gdzie \mathcal{X} jest przestrzenią wartości zmiennej losowej Y (i docelowej zmiennej losowej X). Warunkujemy teraz przez wartości zmiennej Y , otrzymując

$$\begin{aligned} \mathbb{P}(Y \in B, U \leq f(Y)/g(Y)) &= \int_B \mathbb{P}(U \leq f(Y)/g(Y) | Y = y) g(y) dy \\ &= \int_B \frac{f(y)}{g(y)} g(y) dy = \int_B f(y) dy. \end{aligned}$$

Skorzystaliśmy z niezależności Y i U oraz ze wzoru na prawdopodobieństwo całkowite. Otrzymaliśmy licznik we wzorze (3.1). Mianownik dostaniemy zastępując B przez \mathcal{X} . \square

Uwaga 3.1. Ważną zaletą algorytmu jest to, że nie trzeba znać „stałej normującej” gęstości $f/\int f$, czyli liczby $\int f$.

Uwaga 3.2. W istocie \mathcal{X} może być ogólną przestrzenią z miarą μ . Żeby poczuć się pewniej założę, że \mathcal{X} jest przestrzenią polską i miara μ jest σ -skończona. Możemy rozważać gęstości i całki względem miary μ – dowód poprawności algorytmu eliminacji nie ulega zmianie. Nie widzę wielkiego sensu w przeładowaniu notacji, można się umówić, że symbol $\int \cdots dx$ jest *skrót*em $\int \cdots \mu(dx)$. Dociekliwy Czytelnik powinien zastanowić się, w jakiej sytuacji potrafi uzasadnić wszystkie przejścia w dowodzie, powołując się na odpowiednie własności prawdopodobieństwa warunkowego (np. twierdzenie o prawdopodobieństwie całkowitym itp.). W każdym razie, algorytm eliminacji działa poprawnie, gdy

- $\mathcal{X} = \mathbb{R}^d$, gęstości są względem miary Lebesgue’a;
- \mathcal{X} dyskretna, gęstości są względem miary liczącej.

Uwaga 3.3. Efektywność algorytmu zależy od doboru gęstości g tak, aby majoryzowała funkcję f ale nie była dużo od niej większa. Istotnie, liczba prób N do zaakceptowania $X := Y$ ma rozkład geometryczny z prawdopodobieństwem sukcesu $p = \int f/\int g$, zgodnie ze Stwierdzeniem 3.2, zatem $\mathbb{E}N = 1/p$. Iloraz p powinien być możliwie bliski jedynki, co jest możliwe jeśli „kształt funkcji g jest podobny do f ”.

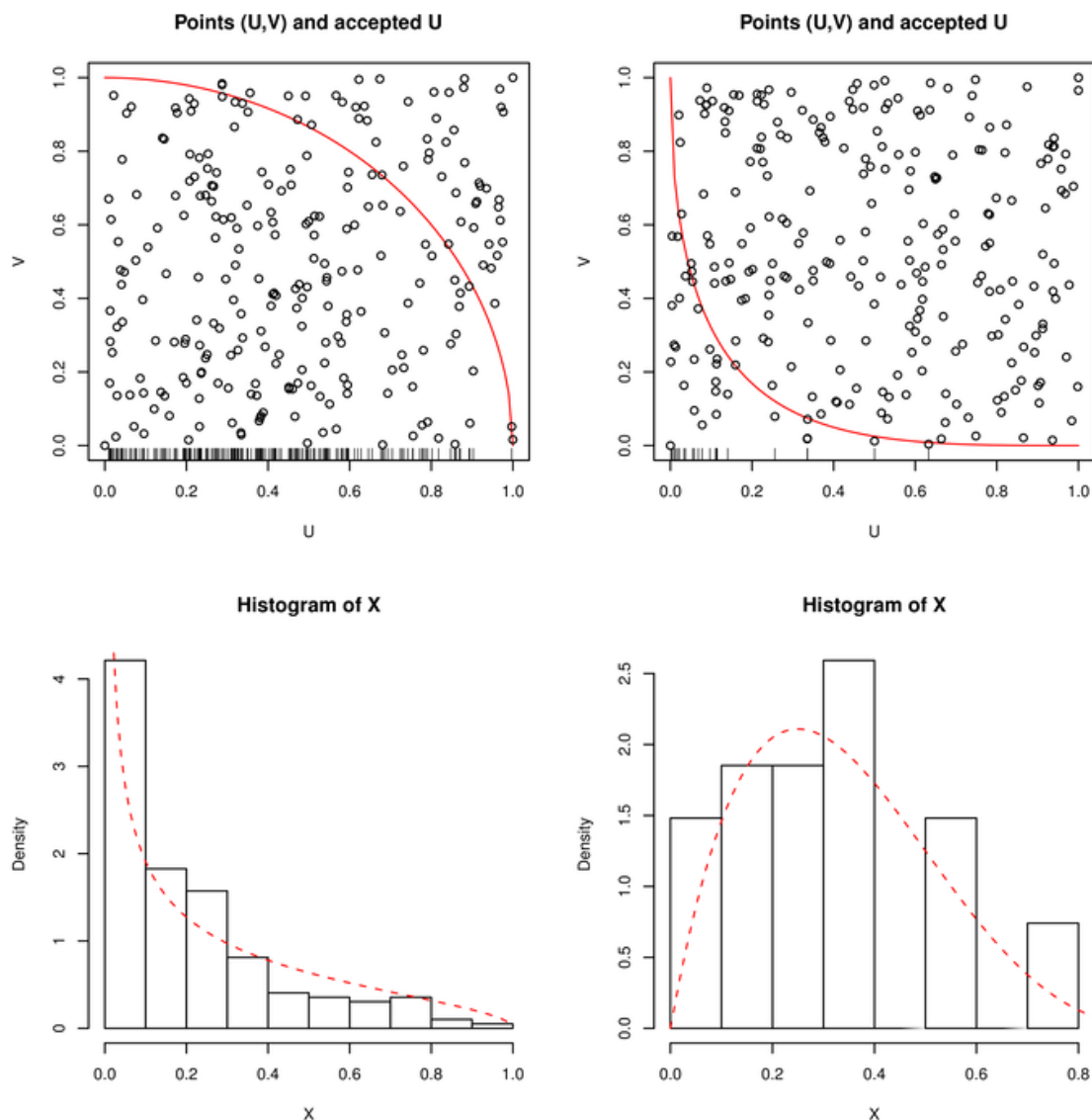
Zwykle metoda eliminacji stosuje się w połączeniu z odpowiednio dobranymi przekształceniami. Zilustrujemy to na poniższym przykładzie.

Przykład 3.9 (Rozkład Beta, algorytm Bermiana). Niech U i V będą niezależnymi zmiennymi losowymi o rozkładzie jednostajnym $U(0, 1)$. Pod warunkiem, że $U^{1/\alpha} + V^{1/\beta}$ zmienna losowa U ma gęstość $f_U(u)$ proporcjonalną do funkcji $(1 - u^{1/\alpha})^\beta$. Zmienna losowa $X = U^{1/\alpha}$ ma gęstość $f_X(x) = f_U(x^\alpha)(dx^\alpha/dx) \propto (1 - x)^\beta x^\alpha$. Zatem X ma rozkład $\text{Beta}(\alpha, \beta + 1)$. Algorytm jest więc następujący:

Listing.

```
repeat
  Gen  $U, V$ 
   $X := U^{1/\alpha}; Y := V^{1/\beta}$ 
until  $X + Y \leq 1$ ;
return  $X$ 
```

Efektywność tego algorytmu jest tym większa, im mniejsze są parametry α i β (frakcja zaakceptowanych U jest wtedy duża). Dla ilustracji rozpatrzmy dwa przypadki.



Rysunek 3.3. Algorytm Bermiana: losowanie z rozkładu beta.

Na rysunku 3.3, po lewej stronie $\alpha = \beta = 0.5$. Spośród $n = 250$ punktów zaakceptowano $N = 197$. Przerywana krzywa jest wykresem gęstości rozkładu $\text{Beta}(0.5, 1, 5)$. Po prawej stronie $\alpha = 2, \beta = 3$. Spośród $n = 250$ punktów zaakceptowano tylko $N = 24$. Tym razem wykres pokazuje gęstość rozkładu $\text{Beta}(2, 4)$. Zaakceptowane punkty U (o gęstości f_U) są widoczne w postaci „kresczek” na górnych rysunkach. Ciągłą linią jest narysowana funkcja $(1 - u^{1/\alpha})^\beta \propto f_U(u)$.

3.3.2. Eliminacja w R

Zasadniczy algorytm eliminacji opisany w Stwierdzeniu 3.2 i w Punkcie 3.3.1 polega na powtarzaniu generacji tak długo, aż zostanie spełnione kryterium akceptacji. Liczba prób jest losowa, a rezultatem jest jedna zmienna losowa o zadanym rozkładzie. Specyfika języka R narzuca inny sposób przeprowadzania eliminacji. Działamy na wektorach, a więc od razu produkujemy

n niezależnych zmiennych X_1, \dots, X_n o rozkładzie g , następnie poddajemy je wszystkie procedurze eliminacji. Przez sito eliminacji, powiedzmy $U_i < f(X_i)/g(X_i)$, przechodzi pewna część zmiennych. Otrzymujemy *losową* liczbę zmiennych o rozkładzie proporcjonalnym do f . Liczba zaakceptowanych zmiennych ma oczywiście rozkład dwumianowy $\text{Bin}(n, p)$ z $p = \int f / \int g$. To nie jest zbyt eleganckie. W sztuczny sposób można zmusić R do wyprodukowania *zadanej*, nie-losowej liczby zaakceptowanych zmiennych, ale jeśli nie ma specjalnej konieczności, lepiej tego nie robić (przymus rzadko prowadzi do pozytywnych rezultatów).

Dla przykładu, generowanie zmiennych (X, Y) o rozkładzie jesnostajnym na kole $\{x^2 + y^2 \leq 1\}$ może w praktyce wyglądać tak:

```
> n <- 1000
> X <- runif(n,min=-1,max=1)  \# generowanie z rozkładu jednostajnego na [-1,1]
> Y <- runif(n,min=-1,max=1)
> Accept <- X^2+Y^2<1        \# wektor logiczny
> X <- X[Accept]
> Y <- Y[Accept]
```

Otrzymujemy pewną losową liczbę (około 785) punktów (X, Y) w kole jednostkowym.

3.4. Metoda kompozycji

Jest to niezwykle prosta technika generowania zmiennych losowych. Załóżmy, że docelowy rozkład jest mieszanką prostszych rozkładów prawdopodobieństwa, czyli jego gęstość jest kombinacją wypukłą postaci

$$f(x) = \sum_{i=1}^k p_i f_i(x), \quad \left(p_i \geq 0, \sum_{i=1}^k p_i = 1 \right).$$

Jeśli umiemy losować z każdej gęstości f_i to możemy uruchomić dwuetapowe losowanie:

Listing.

```
Gen  $I \sim p(\cdot)$ ; { to znaczy  $\mathbb{P}(I = i) = p_i$  }
Gen  $X \sim f_I$ ; { jeśli  $i = i$  to uruchamiamy generator rozkładu  $f_i$  }
return  $X$ 
```

Jasne, że na wyjściu mamy $X \sim F$. W istocie jest to szczególny przypadek metody rozkładów warunkowych, którą omówię później, w Rozdziale 5.

Przykład 3.10 (Rozkład Laplace’a). Rozkład Laplace’a (podwójny rozkład wykładniczy) ma gęstość

$$f(x) = \frac{1}{2\lambda} e^{-\lambda|x|}.$$

Można go „skomponować” z dwóch połówek rozkładu wykładniczego:

Listing.

```
Gen  $W \sim \text{Ex}(\lambda)$ ; { generujemy z rozkładu wykładniczego }
Gen  $U \sim U(0, 1)$ ;
if  $U < 1/2$  then  $X := W$  else  $X := -W$  { losowo zmieniamy znak }
return  $X$ 
```

4. Generowanie zmiennych losowych II. Specjalne metody

4.1. Rozkłady dyskretne

Jak wylosować zmienną losową I o rozkładzie $\mathbb{P}(I = i) = p_i$ mając dane p_1, p_2, \dots ? Metoda odwracanie dystrybucyj w przypadku rozkładów dyskretnych (Przykład 3.6) przyjmuje następującą postać. Obliczamy s_1, s_2, \dots , gdzie $s_k = \sum_{i=1}^k p_i$. Losujemy $U \sim U(0, 1)$ i szukamy przedziału $]s_{I-1}, s_I]$ w którym leży U .

Przykład 4.1 (Algorytm prymitywny). W zasadzie można to zrobić tak:

Listing.

```
Gen  $U, I := 1$ 
while  $s_I \leq U$  do  $I := I + 1$ ;
return  $I$ 
```

Problemem jest mała efektywność tego algorytmu. Jedno z możliwych ulepszeń polega na bardziej „inteligentnym” lokalizowaniu przedziału $]s_{I-1}, s_I] \ni U$, na przykład metodą bisekcji.

Przykład 4.2 (Algorytm podziału odcinka). Załóżmy, że mamy rozkład prawdopodobieństwa (p_i) na przestrzeni skończonej i liczby p_i są uporządkowane: $p_1 > \dots > p_m$.

Listing.

```
Gen  $U$ ;
 $L := 0; R := m$ ;
repeat
   $I := \lfloor (L + R)/2 \rfloor$ ;
  if  $U > s_I$  then  $L := I$  else  $R := I$ ;
until  $L \geq R - 1$ 
return  $I$ 
```

Przedstawię teraz piękną metodę opartą na innym pomysle.

Przykład 4.3 (Metoda „Alias”). Przypuśćmy, że mamy dwa ciągi liczb: q_1, \dots, q_m , gdzie $0 < q_i < 1$ oraz $a(1), \dots, a(m)$, gdzie $a(i) \in \{1, \dots, m\}$ (to są owe „aliasy”). Rozaptrzymy taki algorytm:

Listing.

```
Gen  $K \sim U\{1, \dots, M\}$ ;
Gen  $U$ ;
if  $U < q_K$  then  $I := K$  else  $I := a(K)$ ;
return  $I$ 
```

Jest jasne, że

$$\mathbb{P}(I = i) = \frac{1}{m} \left(q_i + \sum_{j:a(j)=i} (1 - q_j) \right).$$

Jeśli mamy zadany rozkład prawdopodobieństwa (p_1, \dots, p_m) i chcemy, żeby $\mathbb{P}(I = i) = p_i$, to musimy dobrać odpowiednio q_i i $a(i)$. Można zawsze tak zrobić, i to na wiele różnych sposobów. Opracowanie algorytmu dobierania wektorów q i a do zadanego p pozostawiamy jako ćwiczenie.

4.2. Schematy kombinatoryczne

4.2.1. Pobieranie próbki bez zwracania

Spośród r obiektów chcemy wybrać losowo n tak, aby każdy z $\binom{r}{n}$ podzbiorów miał jednakowe prawdopodobieństwo. Oczywiście, można losować tak, jak ze zwracaniem, tylko odrzucać elementy wylosowane powtórnie.

Przykład 4.4 (Losowanie bez zwracania I). W tablicy $c(1), \dots, c(r)$ zaznaczamy, które elementy zostały wybrane.

Listing.

```

for  $i := 1$  to  $r$  do  $c(i) := \text{false}$ ;
 $i := 0$ 
repeat
  repeat
    Gen  $K \sim U\{1, \dots, r\}$ 
  until  $c(K) = \text{false}$ ;
   $c(K) := \text{true}$ ;
   $i := i + 1$ ;
until  $i = n$ 

```

Pewnym ulepszeniem tej prymitywnej metody jest następujący algorytm.

Przykład 4.5 (Losowanie bez zwracania II). Tablica $c(1), \dots, c(r)$ ma takie samo znaczenie jak w poprzednim przykładzie. Będziemy teraz „przeglądać” elementy $1, \dots, r$ po kolei, decydując o zaliczeniu do próbki kolejnego elementu zgodnie z odpowiednim prawdopodobieństwem warunkowym. Niech i oznacza liczbę wybranych, zaś $t - 1$ - liczbę przejrzanych poprzednio elementów.

Listing.

```

for  $i := 1$  to  $r$   $c(i) := \text{false}$ ;
 $t := 1$ ;  $i := 0$ 
repeat
  Gen  $U$ ;
  if  $U \leq \frac{n - i}{r - (t - 1)}$  then
    begin
       $c(t) := \text{true}$ ;
       $i := i + 1$ 
    end;
   $t := t + 1$ ;
until  $i = n$ 

```

Przykład 4.6 (Losowanie bez zwracania III). Tablica $s(1), \dots, s(n)$ będzie teraz zawierała numery wybieranych elementów. Podobnie jak poprzednio, i – oznacza liczbę elementów wybranych, $t - 1$ – liczbę przejranych.

Listing.

```

for  $i := 1$  to  $n$  do  $s(i) := i$ ;
for  $t := n + 1$  to  $r$  do
  begin
    Gen  $U$ ;
    if  $U \leq n/t$  then
      begin
        Gen  $K \sim U\{1, \dots, n\}$ ;
         $s(K) := t$ 
      end
    end
  end

```

Uzasadnienie poprawności tego algorytmu jest rekurencyjne. Załóżmy, że przed t -tym losowaniem każda próbka wybrana ze zbioru $\{1, \dots, t - 1\}$ ma prawdopodobieństwo

$$\binom{t-1}{n}^{-1} = \frac{n!}{(t-1) \cdots (t-n)}.$$

W kroku t ta próbka „przeżywa” czyli pozostaje bez zmiany z prawdopodobieństwem $1 - n/t$ (jest to prawdopodobieństwo, że próbka wylosowana ze zbioru $\{1, \dots, t\}$ nie zawiera elementu t). Zatem po kroku t każda próbka nie zawierająca elementu t ma prawdopodobieństwo

$$\frac{n!}{(t-1) \cdots (t-n)} \cdot \frac{t-n}{t} = \binom{t}{n}^{-1}.$$

Z prawdopodobieństwem n/t „podmieniamy” jeden z elementów próbki (losowo wybrany) na element t . Sprawdzenie, że każda próbka zawierająca element t ma po t -tym losowaniu jednakowe prawdopodobieństwo – pozostawiam jako ćwiczenie.

4.2.2. Permutacje losowe

Przez permutację losową rozumiemy uporządkowanie n elementów wygenerowane zgodnie z rozkładem jednostajnym na przestrzeni wszystkich $n!$ możliwych uporządkowań. Permutację liczb $1, \dots, n$ zapiszemy w tablicy $\sigma(1), \dots, \sigma(n)$. Łatwo sprawdzić poprawność następującego algorytmu.

Listing.

```

for  $i := 1$  to  $n$  do  $\sigma(i) := i$ ;
for  $i := 1$  to  $n - 1$  do
  begin
    Gen  $J \sim U\{i, i + 1, \dots, n\}$ ;
     $Swap(\pi(i), \pi(J))$ 
  end

```

Funkcja $Swap$ zamienia miejscami elementy $\pi(i)$ i $\pi(J)$.

4.3. Specjalne metody eliminacji

4.3.1. Iloraz zmiennych równomiernych

Szczególnie często stosowany jest specjalny przypadek metody eliminacji, znany jako algorytm „ilorazu zmiennych równomiernych” (*Ratio of Uniforms*). Zaczniemy od prostego przykładu, który wyjaśni tę nazwę.

Przykład 4.7 (Rozkład Cauchy’ego). Metodą eliminacji z prostokąta $[0, 1] \times [-1, 1]$ otrzymujemy zmienną losową (U, V) o rozkładzie jednostajnym na półkołu $\{(u, v) : u \geq 0, u^2 + v^2 \leq 1\}$. Kąt Φ pomiędzy osią poziomą i punktem (U, V) ma, oczywiście, rozkład $U(-\pi, \pi)$.

Listing.

```
repeat
  Gen  $U \sim U(0, 1)$ ;
  Gen  $V \sim U(-1, 1)$ 
until  $U^2 + V^2 < 1$ 
 $X := V/U$ 
```

Na wyjściu $X \sim \text{Cauchy}$, bo

$$\mathbb{P}(X \leq x) = \mathbb{P}(V \leq xU) = \frac{1}{\pi}\Phi + \frac{1}{2} = \frac{1}{\pi} \arctan(x) + \frac{1}{2}.$$

Ogólny algorytm metody „ilorazu równomiernych” oparty jest na następującym fakcie.

Stwierdzenie 4.1. Załóżmy o funkcji $h : \mathbb{R} \rightarrow \mathbb{R}$, że

$$h(x) \geq 0, \int h(x) dx < \infty.$$

Niech zbiór $C_h \subset \mathbb{R}^2$ będzie określony następująco:

$$C_h = \left\{ (u, v) : 0 \leq u \leq \sqrt{h\left(\frac{u}{v}\right)} \right\}, \quad |C_h| < \infty.$$

Miara Lebesgue’a (pole) tego zbioru jest skończona, $|C_h| < \infty$, a zatem można mówić o rozkładzie jednostajnym $U(C_h)$.

Jeżeli $(U, V) \sim U(C_h)$ i $X = V/U$, to X ma gęstość proporcjonalną do funkcji h ($X \sim h/\int h$).

Dowód. „Pole figury” C_h jest równe

$$\begin{aligned} |C_h| &= \iint_{C_h} du dv = \iint_{0 \leq u \leq \sqrt{h(x)}} u du dx \\ &= \int_{-\infty}^{\infty} \int_0^{\sqrt{h(x)}} du dx = \frac{1}{2} \int_{-\infty}^{\infty} h(x) dx < \infty \end{aligned}$$

Dokonaliśmy tu zamiany zmiennych:

$$(u, v) \mapsto \left(u, x = \frac{v}{u}\right).$$

Jakobian tego przekształcenia jest równy

$$\frac{\partial(u, x)}{\partial(u, v)} = \det \begin{pmatrix} 1 & 0 \\ v/u^2 & 1/u \end{pmatrix} = \frac{1}{u}$$

W podobny sposób, na mocy znanego wzoru na gęstość przekształconych zmiennych losowych obliczamy łączną gęstość (U, X) :

$$f_{U,X}(u, x) = u f_{U,V}(u, ux) \quad \text{na zbiorze } \{0 \leq u \leq \sqrt{h(x)}\}.$$

Stąd dostajemy gęstość brzegową X :

$$\int_0^{\sqrt{h(x)}} \frac{u \, du}{|C_h|} = \frac{h(x)}{2|C_h|}.$$

□

Żeby wylosować $(U, V) \sim U(C_h)$ stosuje się zazwyczaj eliminację z prostokąta. Użyteczne jest następujące oszacowanie boków tego prostokąta.

Stwierdzenie 4.2. *Jeśli funkcje $h(x)$ i $x^2 h(x)$ są ograniczone, wtedy*

$$C_h \subseteq [0, a] \times [b_-, b_+],$$

gdzie

$$a = \sqrt{\sup_x h(x)},$$

$$b_+ = \sqrt{\sup_{x \geq 0} [x^2 h(x)]}, \quad b_- = -\sqrt{\sup_{x \leq 0} [x^2 h(x)]}.$$

Dowód. Jeśli $(u, v) \in C_h$ to oczywiście $0 \leq u \leq \sqrt{h(v/u)} \leq \sqrt{\sup_x h(x)}$. Załóżmy dodatkowo, że $v \geq 0$ i przejdźmy do zmiennych (v, x) , gdzie $x = v/u$. Nierówność $u \leq \sqrt{h(v/u)}$ jest równoważna $v^2 \leq x^2 h(x)$. Ponieważ $x \geq 0$, więc dostajemy $v^2 \leq b_+^2$. Dla $v \leq 0$ mamy analogicznie $v^2 \leq b_-^2$. □

Ogólny algorytm RU jest następujący:

Listing.

```
repeat
  Gen  $U_1, U_2$ ;
   $U := aU_1$ ;  $V := b_- + (b_+ - b_-)U_2$ 
until  $(U, V) \in C_h$ ;
 $X := \frac{V}{U}$ 
```

Przykład 4.8 (Rozkład normalny). Niech

$$h(x) = \exp\left(-\frac{1}{2}x^2\right).$$

Wtedy

$$C_h = \left\{ (u, v) : 0 \leq u \leq \exp\left(-\frac{1}{4}\frac{v^2}{u^2}\right) \right\} = \left\{ \frac{v^2}{u^2} \leq -4 \ln u \right\}.$$

Zauważmy, że $a = 1$ i $b_+ = -b_- = \sqrt{2e^{-1}}$. Otrzymujemy następujący algorytm:

Listing.

```

repeat
  Gen  $U_1, U_2$ 
   $U := U_1; V := \sqrt{2e^{-1}}(2U_2 - 1);$ 
   $X := \frac{V}{U}$ 
until  $X^2 \leq -4 \ln U$ 

```

4.3.2. Gęstości przedstawione szeregami

Ciekawe, że można skonstruować dokładne algorytmy eliminacji bez konieczności dokładnego obliczania docelowej gęstości f . Podstawowy pomysł jest następujący. Niech f , podobnie jak poprzednio, będzie funkcją proporcjonalną do gęstości (nie jest potrzebna stała normująca) i $f \leq g$. Załóżmy, że mamy dwa ciągi funkcji, przybliżające f z dołu i z góry:

$$\underline{f}_n \leq f \leq \overline{f}_n, \quad \underline{f}_n \rightarrow f, \overline{f}_n \rightarrow f \quad (n \rightarrow \infty).$$

Jeśli umiemy ewaluować funkcje \underline{f}_n i \overline{f}_n to możemy uruchomić następujący algorytm:

Listing.

```

repeat
  Gen  $Y \sim g;$ 
  Gen  $U;$ 
   $W := Ug(Y);$ 
  repeat
     $n := n + 1;$ 
    if  $W \leq \underline{f}_n(Y)$  then
      begin  $X := Y;$  return  $X;$  stop end
    until  $W > \overline{f}_n(Y)$ 
until FALSE

```

Zbieżność przybliżeń dolnych i górnych do funkcji f gwarantuje, że ten algorytm na pewno się kiedyś zatrzyma. Fakt, że na wyjściu $X \sim f/\int f$ jest widoczny.

Poniższy algorytm jest znany jako **metoda szeregów zbieżnych**. Zakładamy, że

$$f(x) = \sum_{i=1}^{\infty} a_n(x),$$

przy czym reszty tego szeregu umiemy oszacować z góry przez znane funkcje,

$$\left| \sum_{i=n+1}^{\infty} a_n(x) \right| \leq r_{n+1}(x).$$

Pseudo-kod algorytmu jest taki:

Listing.

```

repeat
  Gen  $Y \sim g;$ 
  Gen  $U;$ 
   $W := Ug(Y);$ 
   $s := 0;$ 
   $n := 0;$ 

```

```

repeat
   $n := n + 1$ ;
   $s := s + a_n(Y)$ ;
  until  $|s - W| > r_{n+1}(Y)$ ;
until  $W \leq s$ ;
return  $X$ 

```

Rzecz jasna, w tym algorytmie zmienna s przybiera kolejno wartości równe sumom częściowym szeregu, $s_n(x) = \sum_{i=1}^n a_i(x)$.

Metoda szeregów naprzemiennych jest dostosowana do sytuacji gdy

$$f(x) = g(x) \sum_{i=0}^{\infty} (-1)^i a_i(x) = g(x) [1 - a_1(x) + a_2(x) - a_3(x) \dots],$$

gdzie $1 = a_0(x) \geq a_1(x) \geq a_2(x) \geq \dots \geq 0$. Wiadomo, że w takiej sytuacji sumy częściowe przybliżają sumę szeregu na przemian z nadmiarem i z niedomiarem. Ten fakt wykorzystuje algorytm szeregów naprzemiennych:

Listing.

```

repeat
  Gen  $Y \sim g$ ;
  Gen  $U$ ;
   $W := Ug(Y)$ ;
   $s := 0$ ;
   $n := 0$ ;
  repeat
     $n := n + 1$ ; {  $n$  nieparzyste }
     $s := s + a_n(Y)$ ;
    if  $U \geq s$  then  $X := Y$ ; return  $X$ ; stop end
     $n := n + 1$ ; {  $n$  parzyste }
     $s := s - a_n(Y)$ ;
    until  $U > s$ ;
until FALSE

```

Przykład 4.9. Rozkład Kołmogorowa-Smirnowa ma dystrybuantę

$$F(x) = \sum_{n=-\infty}^{\infty} (-1)^n e^{-2n^2 x^2}, \quad (x \geq 0)$$

i gęstość

$$f(x) = 8 \sum_{n=1}^{\infty} (-1)^{n+1} n^2 x e^{-2n^2 x^2}, \quad (x \geq 0).$$

Możemy zastosować metodę szeregów naprzemiennych, kładąc

$$g(x) = 4x e^{-2x^2} \quad (x \geq 0)$$

oraz

$$a_n(x) = n^2 x e^{-2(n^2-1)x^2}.$$

5. Generowanie zmiennych losowych III. Rozkłady wielowymiarowe

5.1. Ogólne metody

Wiele spośród ogólnych metod generowania zmiennych losowych jest całkowicie niezależnych od wymiaru. W szczególności, metody *eliminacji*, *kompozycji* i *przekształceń* są z powodzeniem stosowane do generowania zmiennych losowych wielowymiarowych. Wyjątek stanowi „najbardziej ogólna” metoda *odwracania dystrybucyj*, która nie ma naturalnego odpowiednika dla wymiaru większego niż 1.

5.1.1. Metoda rozkładów warunkowych

Jest to właściwie jedyna metoda „w zasadniczy sposób wielowymiarowa”. Opiera się na przedstawieniu gęstości łącznej zmiennych losowych X_1, \dots, X_d jako iloczynu gęstości brzegowej i gęstości warunkowych (wzór łańcuchowy):

$$f(x_1, x_2, \dots, x_d) = f(x_1)f(x_2|x_1)f(x_3|x_1, x_2) \cdots f(x_d|x_1, \dots, x_{d-1}).$$

Wynika stąd następujący algorytm:

Listing.

```
for i := 1 to d do
  Gen  $X_i \sim f(\cdot | X_1, \dots, X_{i-1})$ 
```

Przykład 5.1 (Wielowymiarowy rozkład normalny). Ograniczmy się do zmiennych losowych X_1, X_2 pochodzących z rozkładu dwuwymiarowego $N(0, 0, \sigma_1^2, \sigma_2^2, \rho)$ o gęstości

$$f(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp \left[-\frac{1}{2(1-\rho^2)} \left(\frac{x_1^2}{\sigma_1^2} - 2\rho\frac{x_1x_2}{\sigma_1\sigma_2} + \frac{x_2^2}{\sigma_2^2} \right) \right].$$

Jak wiadomo (można to sprawdzić elementarnym rachunkiem),

$$f(x_1) = \frac{1}{\sqrt{2\pi}\sigma_1} \exp \left[-\frac{x_1^2}{2\sigma_1^2} \right],$$

$$f(x_2|x_1) = \frac{1}{\sqrt{2\pi}\sigma_2\sqrt{1-\rho^2}} \exp \left[-\frac{1}{2\sigma_2^2(1-\rho^2)} \left(x_2 - \rho\frac{\sigma_2}{\sigma_1}x_1 \right)^2 \right].$$

To znaczy, że $N(0, \sigma_1^2)$ jest rozkładem brzegowym X_1 oraz

$$N\left(\rho\frac{\sigma_2}{\sigma_1}x_1, \sigma_2^2(1-\rho^2)\right)$$

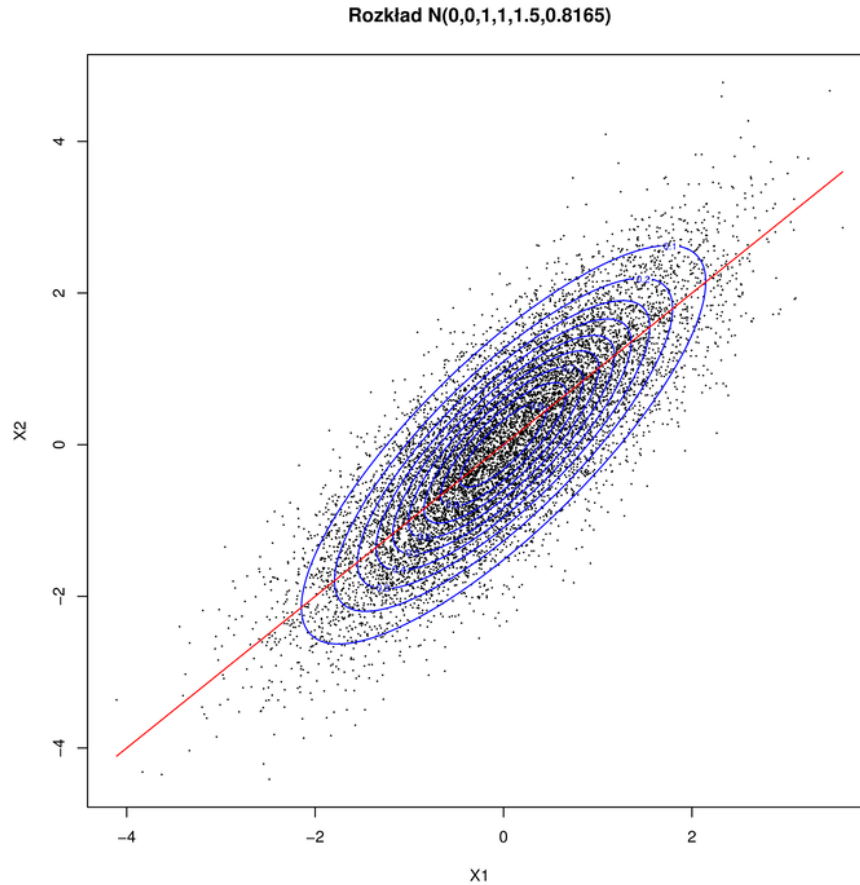
jest rozkładem warunkowym X_2 dla $X_1 = x_1$. Algorytm jest więc następujący.

Listing.

```

Gen  $X_1, X_2 \sim N(0, 1)$ 
 $X_1 := \sigma_1 X_1$ 
 $X_2 := \rho(\sigma_2/\sigma_1)X_1 + \sigma_2\sqrt{1-\rho^2}X_2$ 

```



Rysunek 5.1. Próbkę z dwuwymiarowego rozkładu normalnego, poziomice gęstości i funkcja regresji $x_2 = \mathbb{E}(X_2|X_1 = x_1)$.

Efekt działania powyższego algorytmu widać na Rysunku 5.1. W tym konkretnym przykładzie X_1 ma rozkład brzegowy $N(0, 1)$, zaś X_2 ma rozkład warunkowy $N(X_1, 0.5)$. Zauważmy, że $\text{Var}X_2 = 1.5$ i $\text{Cov}(X_1, X_2) = 1$. Prosta $x_2 = x_1$ jest wykresem funkcji regresji $\mathbb{E}(X_2|X_1 = x_1)$ i jest przedstawiona na wykresie. Pokazane są też poziomice gęstości (elipsy). Warto zwrócić uwagę, że funkcja regresji nie pokrywa się ze wspólną osią tych elips. Dlaczego? Jak obliczyć oś?

Uogólnienie na przypadek rozkładu normalnego $N(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$, z niezerowymi średnimi, jest banalne. Uogólnienie na wyższe wymiary też nie jest skomplikowane. Okazuje się jednak, że metoda rozkładów warunkowych dla rozkładów normalnych prowadzi to do algorytmu identycznego jak otrzymany metodą przekształceń.

5.1.2. Metoda przekształceń

Podstawą jest następujące twierdzenie o przekształcaniu gęstości.

Twierdzenie 5.1. Załóżmy, że X jest d -wymiarową zmienną losową o wartościach w otwartym zbiorze $A \subseteq \mathbb{R}^d$. Jeżeli

$$X \sim f(x)$$

i $h : A \rightarrow \mathbb{R}^d$ jest dyfeomorfizmem, to

$$Y = h(X) \sim g(y) = f(h^{-1}(y)) \left| \det \frac{\partial}{\partial y} h^{-1}(y) \right|.$$

Przykład 5.2 (Wielowymiarowy rozkład normalny). Rozważmy niezależne zmienne losowe $Z_1, \dots, Z_d \sim N(0, 1)$. Wektor $Z = (Z_1, \dots, Z_d)$ ma d -wymiarowy rozkład normalny $N(0, I)$ o gęstości

$$f(z) = (2\pi)^{-d/2} \exp \left[-\frac{1}{2} z^\top z \right].$$

Jeżeli teraz R jest nieosobliwą macierzą ($d \times d$) to przekształcenie $z \mapsto x = Rz$ jest dyfeomorfizmem z jacobianem $\det R$. Z Twierdzenia 5.1 wynika, że wektor losowy $X = RZ$ ma rozkład normalny o gęstości

$$f_X(x) = (2\pi)^{-d/2} (\det R)^{-1/2} \exp \left[-\frac{1}{2} x^\top \Sigma^{-1} x \right],$$

gdzie $\Sigma = RR^\top$. Innymi słowy, $X \sim N(0, \Sigma)$. Algorytm generacji jest oczywisty:

Listing.

```
Gen  $Z \sim N(0, I)$ 
 $X := RZ$ 
```

Jeśli dana jest macierz kowariancji Σ wektora X , to przed uruchomieniem algorytmu trzeba znaleźć taką macierz R , żeby $\Sigma = RR^\top$. Istnieje wiele takich macierzy, ale najlepiej skorzystać z rozkładu Choleskiego i wybrać macierz trójkątną.

5.2. Kilka ważnych przykładów

Różnorodnych rozkładów wielowymiarowych jest więcej, niż gwiazd na niebie – a wśród nich tyle przykładów ciekawych i ważnych! Musiałem wybrać zaledwie kilka z nich. Rzecz jasna, wybrałem te, które mi się szczególnie podobają.

5.2.1. Rozkłady sferyczne i eliptyczne

Najważniejszy, być może, przykład to wielowymiarowy rozkład normalny, omówiony już w 5.2. Rozpatrzmy teraz *rozkłady jednostajne na kuli*

$$B^d = \{x \in \mathbb{R}^d : |x|^2 \leq 1\}$$

i sferze

$$S^{d-1} = \{x \in \mathbb{R}^d : |x|^2 = 1\}.$$

Oczywiście, $|x|$ oznacza normę euklidesową, $|x| = (x_1^2 + \dots + x_d^2)^{1/2} = (x^\top x)^{1/2}$. Rozkład $U(B^d)$ ma po prostu stałą gęstość względem d -wymiarowej miary Lebesgue'a na kuli. Rozkład $U(S^{d-1})$

ma stałą gęstość względem $(d - 1)$ -wymiarowej miary „powierzchniowej” na sferze. Oba te rozkłady są niezmiennicze względem *obrotów* (liniowych przekształceń ortogonalnych) \mathbb{R}^d . Takie rozkłady nazywamy *sferycznie symetrycznymi* lub krócej: **sferycznymi**. Zauważmy, że zmienną losową o rozkładzie $U(S^{d-1})$ możemy interpretować jako losowo wybrany kierunek w przestrzeni $d - 1$ -wymiarowej. Algorytmy „poszukiwań losowych” często wymagają generowania takich losowych kierunków.

Rozkłady jednostajne na kuli i sferze są blisko ze sobą związane.

— Jeśli $V = (V_1, \dots, V_d) \sim U(B^d)$ i $R = |V|$ to

$$Y = \frac{V}{R} = \left(\frac{V_1}{R}, \dots, \frac{V_d}{R} \right) \sim U(S^{d-1}).$$

Latwo też zauważyć, że R jest zmienną losową o rozkładzie $\mathbb{P}(R \leq r) = r^d$ niezależną od Y .

— Jeśli $Y \sim U(S^{d-1})$ i R jest niezależną zmienną losową o rozkładzie $\mathbb{P}(R \leq r) = r^d$ to $V = RY = (RY_1, \dots, RY_d) \sim U(B^d)$.

Zmienną R łatwo wygenerować metodą odwracania dystrybuanty.

Najprostszy wydaje się algorytm eliminacji:

Listing.

repeat

 Gen $V_1, \dots, V_d \sim U(0, 1)$

until $R^2 = V_1^2 + \dots + V_d^2 \leq 1$

Na wyjściu otrzymujemy, zgodnie z żądaniem

$$V = (V_1, \dots, V_d) \sim U(B^d).$$

W istocie, dokładnie ta metoda, dla $d = 2$, była częścią algorytmu biegunowego Marsaglii. Problem w tym, że w wyższych wymiarach efektywność eliminacji gwałtownie maleje. Prawdopodobieństwo akceptacji jest równe stosunkowi „objętości” kuli B^d do kostki $[-1, 1]^d$. Ze znanego wzoru na objętość kuli d -wymiarowej wynika, że

$$\frac{|B^d|}{2^d} = \frac{2\pi^{d/2}}{d\Gamma(d/2)} \cdot \frac{1}{2^d} = \frac{\pi^{d/2}}{d2^{d-1}\Gamma(d/2)} \xrightarrow{d \rightarrow \infty} 0.$$

Zbieżność do zera jest bardzo szybka. Dla dużego d kula jest znikomą częścią opisanej na niej kostki.

Inna metoda, którą z powodzeniem stosuje się dla $d = 2$ jest związana ze współrzędnymi biegunowymi:

Listing.

Gen $\Phi \sim U(0, 2\pi)$;

$Y_1 := \cos \Phi$; $Y_2 := \sin \Phi$;

Gen U ; $R := \sqrt{U}$;

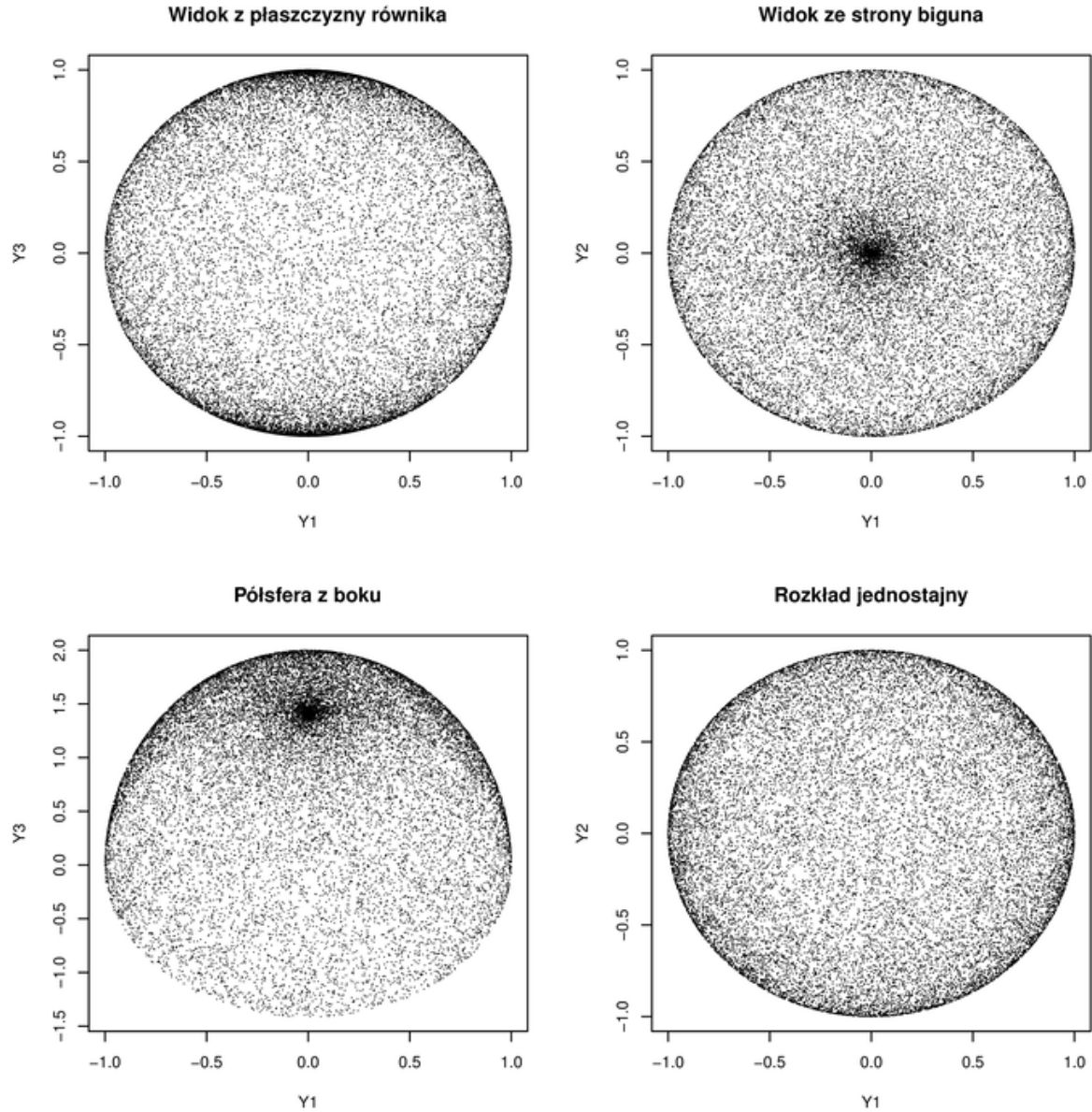
$V_1 := Y_1 \cdot R$; $V_2 := Y_2 \cdot R$;

Na wyjściu $(Y_1, Y_2) \sim S^1$ i $(V_1, V_2) \sim B^2$. Jest to część algorytmu Boxa-Müllera. Uogólnienie na przypadek $d > 2$ nie jest jednak ani proste, ani efektywne. Mechaniczne zastąpienie, współrzędnych biegunowych przez współrzędne sferyczne (dla, powiedzmy $d = 3$) prowadzi do *niepoprawnych wyników*. Popatrzmy na punkty produkowane przez następujący algorytm:

Listing.

Gen $\Phi \sim U(0, 2\pi)$; $\Psi \sim U\left(-\frac{\pi}{2}, \frac{\pi}{2}\right)$

$Y_1 := \cos \Phi \cos \Psi$; $Y_2 := \sin \Phi \cos \Psi$; $Y_3 = \sin \Psi$



Rysunek 5.2. Niepoprawne i poprawne generowanie z rozkładu jednostajnego na sferze.

Efekt działania algorytmu „współrzędnych sferycznych” jest widoczny na trzech początkowych obrazkach na Rysunku 5.2. Górny lewy rysunek przedstawia punkty widoczne „z płaszczyzny równika”, czyli (Y_1, Y_3) . Górny prawy – te same punkty widziane „znad bieguna”, czyli (Y_1, Y_2) . Dolny lewy – górną półsferę widzianą „skośnie”, czyli (Y_2, Y_3^a) , gdzie $Y_3^a = \sqrt{2}Y_3 + \sqrt{2}Y_1$. Dla porównania, na ostatnim rysunku po prawej stronie u dołu – punkty wylosowane rzeczywiście z rozkładu $U(S^2)$ przy pomocy poprawnego algorytmu podanego niżej.

Listing.

```

Gen  $Z_1, \dots, Z_d \sim N(0, 1)$ ;
 $R := (Z_1^2 + \dots + Z_d^2)^{1/2}$ ;
 $Y_1 := Z_1/R, \dots, Y_d := Z_d/R$ ;
Gen  $U$ ;  $R := U^{1/d}$ ;

```

$$V_1 := Y_1 \cdot R, \dots, V_d := Y_d \cdot R;$$

Na wyjściu $Y \sim U(S^{d-1})$ i $V \sim U(B^d)$. Jak widać, algorytm polega na normowaniu punktów wylosowanych ze sferycznie symetrycznego rozkładu normalnego. Jest prosty, efektywny i godny polecenia.

Przykład 5.3 (Wielowymiarowe rozkłady Studenta). Niech $Z = (Z_1, \dots, Z_d)^\top$ będzie wektorem losowym o rozkładzie $N(0, I)$, zaś R^2 – niezależną zmienną losową o rozkładzie $\chi^2(n)$. Wektor

$$(Y_1, \dots, Y_d)^\top = \frac{(Z_1, \dots, Z_d)^\top}{\sqrt{R^2/n}}$$

ma, z definicji, *Sferyczny rozkład t-Studenta z n stopniami swobody*. Gęstość tego rozkładu (z dokładnością do stałej normującej) jest równa

$$f(y) = f(y_1, \dots, y_d) \propto \left[1 + \frac{1}{n} \left(\sum y_i^2\right)\right]^{-(n+d)/2} = \left[1 + \frac{1}{n} |y|^2\right]^{-(n+d)/2}.$$

W przypadku jednowymiarowym, a więc przyjmując $d = 1$, otrzymujemy dobrze znane rozkłady t-Studenta z n stopniami swobody o gęstości

$$f(y) \propto \frac{1}{(1 + y^2/n)^{(n+1)/2}}$$

W szczególnym przypadku, biorąc za liczbę stopni swobody $n = 1$, otrzymujemy rozkłady Cauchy'ego. Na przykład, dwuwymiarowy rozkład Cauchy'ego ma taką gęstość:

$$f(y_1, y_2) \propto \frac{1}{(1 + y_1^2 + y_2^2)^{3/2}}.$$

Użytecznym uogólnieniem rozkładów sferycznych są rozkłady eliptyczne. Są one określone w następujący sposób. Niech Σ będzie macierzą symetryczną i nieosbliwą. Nazwijmy *uogólnionym obrotem* przekształcenie liniowe, które zachowuje uogólnioną normę $|x|_{\Sigma^{-1}} = (x^\top \Sigma^{-1} x)^{1/2}$. Rozkład jest z definicji *eliptycznie konturowany* lub krócej **eliptyczny**, gdy jest niezmienniczy względem uogólnionych obrotów (dla ustalonej macierzy Σ).

5.2.2. Rozkłady Dirichleta

Definicja 5.1. Mówimy, że n -wymiarowa zmienna losowa X ma rozkład Dirichleta,

$$X = (X_1, \dots, X_n) \sim \text{Dir}(\alpha_1, \dots, \alpha_n)$$

jeśli $X_1 + \dots + X_n = 1$ i zmienne X_1, \dots, X_{n-1} mają gęstość

$$f(x_1, \dots, x_{n-1}) = \frac{\Gamma(\alpha_1 + \dots + \alpha_n)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_n)} x_1^{\alpha_1-1} \dots x_{n-1}^{\alpha_{n-1}-1} (1 - x_1 - \dots - x_{n-1})^{\alpha_n-1}.$$

Parametry $\alpha_1, \dots, \alpha_n$ mogą być dowolnymi liczbami dodatnimi.

Uwaga 5.1. Rozkłady Dirichleta dla $d = 2$ są to w istocie rozkłady beta:

$$X_1 \sim \text{Beta}(\alpha_1, \alpha_2) \text{ wtedy i tylko wtedy, gdy } (X_1, X_2) \sim \text{Dir}(\alpha_1, \alpha_2)$$

Wniosek 5.1. Jeśli U_1, \dots, U_{n-1} są niezależnymi zmiennymi o jednakowym rozkładzie jednostajnym $U(0, 1)$ i

$$U_{1:n} < \dots < U_{n-1:n-1}$$

oznaczają statystyki pozycyjne, to spacje

$$X_i = U_{i:n} - U_{i-1:n}, \quad X_n = 1 - U_{n-1:n}$$

mają rozkład $\text{Dir}(1, \dots, 1)$.

Twierdzenie 5.2. Jeśli Y_1, \dots, Y_n są niezależnymi zmiennymi losowymi o rozkładach gamma,

$$Y_i \sim \text{Gamma}(\alpha_i)$$

i $S = Y_1 + \dots + Y_n$ to

$$(X_1, \dots, X_n) = \left(\frac{Y_1}{S}, \dots, \frac{Y_n}{S} \right) \sim \text{Dir}(\alpha_1, \dots, \alpha_n).$$

Wektor losowy X jest niezależny od S .

Dowód. Obliczymy łączną gęstość zmiennych losowych S, X_1, \dots, X_{n-1} . Ze wzoru na przekształcenie gęstości wynika, że

$$\begin{aligned} f_{S, X_1, \dots, X_{n-1}}(s, x_1, \dots, x_{n-1}) &= f_{Y_1, \dots, Y_n}(x_1 s, \dots, x_n s) \\ &\propto (x_1 s)^{\alpha_1-1} e^{-x_1 s} \dots (x_n s)^{\alpha_n-1} e^{-x_n s} \left| \frac{\partial(y_1, \dots, y_n)}{\partial(s, x_1, \dots, x_{n-1})} \right| \\ &\propto x_1^{\alpha_1-1} \dots x_n^{\alpha_n-1} s^{\alpha_1+\dots+\alpha_n-1} e^{-s}, \end{aligned}$$

ponieważ jacobian przekształcenia odwrotnego jest równy s^{n-1} . Wystarczy teraz zauważyć, że

$$\begin{aligned} x_1^{\alpha_1-1} \dots x_n^{\alpha_n-1} &= \text{Dir}, \\ s^{\alpha_1+\dots+\alpha_n-1} e^{-s} &= \text{Gamma}. \end{aligned}$$

□

Wniosek 5.2. Dla niezależnych zmiennych losowych o jednakowym rozkładzie wykładniczym,

$$Y_1, \dots, Y_n \sim \text{Ex}(1),$$

jeśli $S = Y_1 + \dots + Y_n$ to

$$(X_1, \dots, X_n) = \left(\frac{Y_1}{S}, \dots, \frac{Y_n}{S} \right) \sim \text{Dir}(1, \dots, 1)$$

Z 5.1 i 5.2 wynika, że następujące dwa algorytmy:

Listing.

```
Gen  $U_1, \dots, U_{n-1}$ ;
Sort  $(U_1, \dots, U_{n-1})$ ;  $U_0 = 0$ ;  $U_n = 1$ ;
 $X_i := U_i - U_{i-1}$ 
```

oraz

Listing.

```

Gen  $Y_1, \dots, Y_n \sim \text{Ex}(1)$ 
 $S := Y_1 + \dots + Y_n$ 
 $X_i := Y_i/S$ 

```

dają te same wyniki.

Wiele ciekawych własności rozkładów Dirichleta wynika niemal natychmiast z 5.2 (choć nie tak łatwo wyprowadzić je posługując się gęstością 5.2). Mam na myśli przede wszystkim zasadniczą własność „grupowania zmiennych”.

Wniosek 5.3. *Rozważmy rozbitcie zbioru indeksów na sumę rozłącznych podzbiorów:*

$$\{1, \dots, n\} = \bigcup_{j=1}^k I_j.$$

Jeżeli $(X_1, \dots, X_n) \sim \text{Dir}(\alpha_1, \dots, \alpha_n)$ i rozważymy „zgrupowane zmienne”

$$S_j = \sum_{i \in I_j} X_i,$$

to wektor tych zmiennych ma też rozkład Dirichleta,

$$(S_1, \dots, S_k) \sim \text{Dir}(\beta_1, \dots, \beta_k),$$

gdzie

$$\beta_j = \sum_{i \in I_j} \alpha_i.$$

Co więcej, każdy z wektorów $(X_i/S_j)_{i \in I_j}$ ma rozkład Dirichleta $\text{Dir}(\alpha_i)_{i \in I_j}$ i wszystkie te wektory są niezależne od (S_1, \dots, S_k) .

Przykład 5.4. Wniosek 5.3 razem z 5.1 pozwala szybko generować wybrane statystyki pozycyjne. Na przykład łączny rozkład dwóch statystyk pozycyjnych z rozkładu jednostajnego jest wyznaczony przez rozkład trzech „zgrupowanych spacji”:

$$(U_{k:n-1}, U_{l:n-1} - U_{k:n-1}, 1 - U_{l:n-1}) \sim \text{Dir}(k, l - k, n - l)$$

Przykład 5.5 (Rozkład dwumianowy). Aby wygenerować zmienną o rozkładzie dwumianowym $\text{Bin}(n, p)$ wystarczy rozpoznać między którymi statystykami pozycyjnymi z rozkładu $U(0, 1)$ leży liczba p . Nie musimy w tym celu generować *wszystkich* statystyk pozycyjnych, możemy wybierać „najbardziej prawdopodobne”. W połączeniu 5.4 daje to następujący algorytm.

Listing.

```

 $k := n; \theta := p; X := 0;$ 
repeat
   $i := \lfloor 1 + k\theta \rfloor;$ 
  Gen  $V \sim \text{Beta}(i, k + 1 - i);$ 
  if  $\theta < V$  then
    begin  $\theta := \theta/V; k := i - 1$  end
  else
    begin  $X := X + i; \theta = (\theta - V)/(1 - V); k := k - i$  end
until  $k = 0$ 

```

Metoda generowania zmiennych o rozkładzie Dirichleta opiera się na następującym fakcie, który jest w istocie szczególnym przypadkiem „reguły grupowania” 5.3.

Wniosek 5.4. *Jeśli $(X_1, \dots, X_n) \sim \text{Dir}(\alpha_1, \dots, \alpha_n)$ i, dla $k = 1, \dots, n$, określimy $S_k = X_1 + \dots + X_k$ to zmienne*

$$Y_1 = \frac{S_1}{S_2}, Y_2 = \frac{S_2}{S_3}, \dots, Y_n = \frac{S_{n-1}}{S_n}$$

są niezależne i

$$Y_k \sim \text{Beta}(\alpha_1 + \dots + \alpha_k, \alpha_{k+1}).$$

Odwrotnie, jeśli zmienne Y_1, \dots, Y_n są niezależne i każda z nich ma rozkład beta, to wektor (X_1, \dots, X_n) ma rozkład Dirichleta.

Oczywiście, jeśli wygenerujemy niezależne zmienne Y_1, \dots, Y_n o rozkładach beta, to zmienne X_1, \dots, X_n łatwo „odzyskać” przy pomocy wzorów:

$$X_n = 1 - Y_{n-1}$$

$$Z_{n-1} = (1 - Y_{n-2})Y_{n-1}$$

$$\dots$$

$$Z_2 = (1 - Y_1)Y_2 \cdots Y_{n-1}$$

$$Z_1 = Y_1 Y_2 \cdots Y_{n-1}$$

Powyższe równania określają algorytm generowania zmiennych o rozkładzie $\text{Dir}(\alpha_1, \dots, \alpha_n)$.

6. Symulowanie procesów stochastycznych I.

Pełny tytuł tego rozdziału powinien brzmieć „Symulacje Niektórych Procesów Stochastycznych, Bardzo Subiektywnie Wybranych Spośród Mnóstwa Innych”. Nie będę szczegółowo tłumaczył, skąd pochodzi mój subiektywny wybór. Zrezygnowałem z próby przedstawienia procesów z czasem ciągłym i równocześnie ciągłą przestrzenią stanów, bo to temat oddzielny i obszerny.

6.1. Stacjonarne procesy Gaussowskie

Ograniczymy się do dwóch klas procesów, często używanych do modelowania różnych zjawisk. Będą to procesy z czasem dyskretnym i przestrzenią stanów \mathbb{R} , to znaczy ciągi (zależnych) zmiennych losowych $X_0, X_1, \dots, X_n, \dots$ o wartościach rzeczywistych. Niech $\dots, W_{-1}, W_0, W_1, \dots, W_n, \dots$ będzie ciągiem *niezależnych zmiennych losowych* o jednakowym rozkładzie $N(0, v^2)$ (wygodnie posłużyć się tutaj ciągiem indeksowanym wszystkimi liczbami całkowitymi).

Definicja 6.1. Proces **ruchomych średnich** rzędu q , w skrócie $MA(q)$ jest określony równaniem

$$X_n = \beta_1 W_{n-1} + \dots + \beta_q W_{n-q}, \quad (n = 0, 1, \dots),$$

gdzie β_1, \dots, β_q jest ustalonym ciągiem współczynników.

Sposób generowania takiego procesu jest oczywisty i wynika wprost z definicji. Co więcej widać, że proces $MA(q)$ jest *stacjonarny*, to znaczy łączny rozkład prawdopodobieństwa zmiennych X_0, X_1, \dots, X_n jest taki sam jak zmiennych $X_k, X_{k+1}, \dots, X_{k+n-1}$, dla dowolnych n i k . Intuicyjnie, proces nie zmienia się po „przesunięciu czasu” o k jednostek.

Definicja 6.2. Proces **autoregresji** rzędu p , w skrócie $AR(p)$ jest określony równaniem rekurencyjnym

$$X_n = \alpha_1 X_{n-1} + \dots + \alpha_p X_{n-p} + W_n, \quad (n = p, p+1, \dots),$$

gdzie $\alpha_1, \dots, \alpha_p$ jest ustalonym ciągiem współczynników.

Procesy autoregresji wydają się bardzo odpowiednie do modelowania „szeregów czasowych”: stan układu w chwili n zależy od stanów przeszłych i dotatkowo jeszcze od przypadku. Procesy $AR(1)$, w szczególności są łańcuchami Markowa. Sposób generowania procesów $AR(p)$ jest też bezpośrednio widoczny z definicji. Pojawia się jednak pewien problem. Jak znaleźć X_0, \dots, X_{p-1} na początku algorytmu w taki sposób, żeby proces był stacjonarny? Jest to o tyle istotne, że rzeczywiste procesy (na przykład czeregi czasowe w zastosowaniach ekonomicznych) specjaliści uznają za stacjonarne, przynajmniej w przybliżeniu.

Rozważmy dla wygody oznaczeń podwójnie nieskończony proces

$$\dots, X_{-1}, X_0, X_1, \dots$$

spełniający równanie autoregresji rzędu p . Załóżmy, że ten proces jest stacjonarny i wektor X_0, \dots, X_{p-1} rozkład normalny, $N(0, \Sigma)$. Stacjonarność implikuje, że elementy macierzy Σ muszą być postaci $\text{Cov}(X_i, X_j) = \sigma^2 \rho_{i-j}$. Mamy przy tym $\rho_{-k} = \rho_k$, co może być traktowane jako

wygodna konwencja (po to właśnie „rozszerzamy” proces w obie strony). Zastosowanie równania definiującego autoregresję prowadzi do wniosku, że

$$\begin{aligned}\sigma^2 \rho_k &= \text{Cov}(X_0, X_k) = \text{Cov}(X_0, \alpha_1 X_{k-1} + \dots + \alpha_p X_{k-p} + W_k) \\ &= \sigma^2 \alpha_1 \rho_{k-1} + \sigma^2 \alpha_2 \rho_{k-2} + \dots + \sigma^2 \alpha_p \rho_{k-p}.\end{aligned}$$

Podobnie,

$$\begin{aligned}\sigma^2 &= \text{Var}(X_0) = \text{Cov}(X_0, \alpha_1 X_{-1} + \dots + \alpha_p X_{-p} + W_0) \\ &= \sigma^2 \alpha_1 \rho_1 + \sigma^2 \alpha_2 \rho_2 + \dots + \sigma^2 \alpha_p \rho_p + v^2,\end{aligned}$$

gdzie $v^2 = \text{Var}(W_0)$. Otrzymujemy następujący układ równań na współczynniki autokorelacji ρ_k :

$$\begin{cases} \rho_1 = \alpha_1 + \alpha_2 \rho_1 + \alpha_3 \rho_2 + \dots + \alpha_p \rho_{p-1}, \\ \rho_2 = \alpha_1 \rho_1 + \alpha_2 + \alpha_3 \rho_1 + \dots + \alpha_p \rho_{p-2}, \\ \dots, \\ \rho_p = \alpha_1 \rho_{p-1} + \alpha_2 \rho_{p-2} + \alpha_3 \rho_2 + \dots + \alpha_p,\end{cases} \quad (6.1)$$

Można pokazać, że ten układ ma rozwiązanie, jeśli wielomian charakterystyczny $A(z) = 1 - \alpha_1 z - \dots - \alpha_p z^p$ nie ma zer w kole $\{|z| \leq 1\}$. Ponadto mamy równanie na wariancję stacjonarną:

$$\sigma^2 = \frac{v^2}{1 - \rho_1 \alpha_1 - \dots - \alpha_p \rho_p} \quad (6.2)$$

Metoda generowania *stacjonarnego* procesu $\text{AR}(p)$, $X_0, X_1, \dots, X_p, \dots$ jest następująca. Znajdujemy rozwiązanie układu równań (6.1), wariancję obliczamy ze wzoru (6.2) i tworzymy macierz $\Sigma = (\sigma^2 \rho_{i-j})_{i,j=0,\dots,p-1}$. Generujemy wektor losowy $(X_0, X_1, \dots, X_{p-1}) \sim N(0, \Sigma)$ i dalej generujemy rekurencyjnie X_p, X_{p+1}, \dots używając równania autoregresji. Aby się przekonać, że tak generowany proces jest stacjonarny, wystarczy sprawdzić że identyczne są rozkłady wektorów $(X_0, X_1, \dots, X_{p-1})$ i (X_1, X_2, \dots, X_p) . Mamy

$$\begin{pmatrix} X_1 \\ \cdot \\ \cdot \\ \cdot \\ X_p \end{pmatrix} = \begin{pmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ 0 & \ddots & \ddots & \ddots & 0 \\ 0 & \dots & \dots & 0 & 1 \\ \alpha_p & \dots & \dots & \dots & \alpha_1 \end{pmatrix} \begin{pmatrix} X_0 \\ \cdot \\ \cdot \\ \cdot \\ X_{p-1} \end{pmatrix} + \begin{pmatrix} 0 \\ \cdot \\ \cdot \\ \cdot \\ W_p \end{pmatrix}.$$

Niech R oznacza „dużą macierz” w tym wzorze. Z własności wielowymiarowych rozkładów normalnych wynika, że wystarczy sprawdzić równość

$$\Sigma = R \Sigma R^T + \begin{pmatrix} 0 & 0 & \ddots & 0 \\ 0 & \ddots & \ddots & 0 \\ 0 & \ddots & \ddots & 0 \\ 0 & 0 & \dots & v^2 \end{pmatrix}.$$

Macierz Σ została tak wybrana, że ta równość jest spełniona.

6.2. Procesy Poissona

6.2.1. Jednorodny proces Poissona na półprostej

Definicja 6.3. Rozważmy niezależne zmienne losowe W_1, \dots, W_k, \dots o jednakowym rozkładzie wykładniczym, $X_i \sim \text{Ex}(\lambda)$ i utwórzmy kolejne sumy

$$T_0 = 0, T_1 = W_1, T_2 = W_1 + W_2, \dots, T_k = W_1 + \dots + W_k, \dots$$

Niech, dla $t \geq 0$,

$$N(t) = \max\{k : T_k \leq t\}.$$

Rodzinę zmiennych losowych $N(t)$ nazywamy *procesem Poissona*.

Proces Poissona dobrze jest wyobrażać sobie jako *losowy zbiór punktów* na półprostej: $\{T_1, T_2, \dots, T_k, \dots\}$. Zmienna $N(t)$ oznacza liczbę punktów, które „wpadły” w odcinek $]0, t]$. Wygodnie będzie używać symbolu

$$N(s, t) = N(t) - N(s)$$

dla oznaczenia liczby punktów, które „wpadły” w odcinek $]s, t]$.

Stwierdzenie 6.1. *Jeśli $N(t)$ jest procesem Poissona, to*

$$\mathbb{P}(N(t) = k) = e^{-\lambda t} \frac{(\lambda t)^k}{k!}.$$

Dowód. Zauważmy, że

$$T_k \sim \text{Gamma}(k, \lambda).$$

Wobec tego ze wzoru na prawdopodobieństwo całkowite wynika, że

$$\begin{aligned} \mathbb{P}(N(t) = k) &= \mathbb{P}(T_k \leq t, T_{k+1} > t) \\ &= \int_0^t \mathbb{P}(T_{k+1} > t | T_k = s) f_{T_k}(s) ds \\ &= \int_0^t \mathbb{P}(W_{k+1} > t - s | T_k = s) f_{T_k}(s) ds \\ &= \int_0^t e^{-\lambda(t-s)} \frac{\lambda^k}{\Gamma(k)} s^{k-1} e^{-\lambda s} ds \\ &= e^{-\lambda t} \frac{\lambda^k}{\Gamma(k)} \int_0^t s^{k-1} ds = e^{-\lambda t} \frac{\lambda^k}{(k-1)!} \frac{t^k}{k} = e^{-\lambda t} \frac{(\lambda t)^k}{k!}. \end{aligned}$$

□

Oczywiście $\mathbb{E}N(t) = \lambda t$. Liczba

$$\lambda = \frac{1}{\mathbb{E}W_i} = \frac{\mathbb{E}N(t)}{t}$$

jest nazywana *intensywnością procesu*.

Stwierdzenie 6.2. *Jeśli*

$$0 < t_1 < t_2 < \dots < t_i < \dots,$$

to zmienne losowe $N(t_1), N(t_1, t_2), \dots$ są niezależne i każda z nich ma rozkład Poissona:

$$N(t_{i-1}, t_i) \sim \text{Poiss}(\lambda(t_i - t_{i-1})).$$

Dowód. Pokażemy, że warunkowo, dla $N(t_1) = k$, ciąg zmiennych losowych

$$T_{k+1} - t_1, W_{k+2}, W_{k+3} \dots, \text{ jest iid } \sim \text{Ex}(\lambda).$$

Wynika to z własności braku pamięci rozkładu wykładniczego. W istocie, dla ustalonych $N(t_1) = k$ i $S_k = s$ mamy

$$\begin{aligned} & \mathbb{P}(T_{k+1} - t_1 > t | N(t_1) = k, T_k = s) \\ &= \mathbb{P}(T_{k+1} - t_1 > t | T_k = s, T_{k+1} > t_1) \\ &= \mathbb{P}(W_{k+1} > t_1 + t - s | W_{k+1} > t_1 - s) \\ &= \mathbb{P}(W_{k+1} > t) = e^{-\lambda t}. \end{aligned}$$

Fakt, że zmienne $W_{k+2}, W_{k+3} \dots$ są niezależne od zdarzenia $N(t_1) = k$ jest oczywisty. Pokazaliśmy w ten sposób, że losowy zbiór punktów $\{T_{k+1} - t_1, T_{k+2} - t_1, \dots\}$ ma (warunkowo, dla $N(t_1) = k$) taki sam rozkład prawdopodobieństwa, jak $\{T_1, T_2, \dots\}$. Proces Poissona *obserwowany od momentu t_1 jest kopią wyjściowego procesu*. Wynika stąd w szczególności, że zmienna losowa $N(t_2, t_1)$ jest niezależna od $N(t_1)$ i $N(t_2, t_1) \sim \text{Poiss}(\lambda(t_2 - t_1))$. Dalsza część dowodu przebiega analogicznie i ją pominiemy. \square

Metoda generowania procesu Poisson oparta na Definicji 6.3 jest raczej oczywista. Zauważmy jednak, że nie jest to *jedyna* metoda. Inny sposób generowania (i inny sposób patrzenia na proces Poissona) jest związany z następującym faktem.

Stwierdzenie 6.3. *Warunkowo, dla $N(t) = n$, ciąg zmiennych losowych*

$$T_1, \dots, T_n$$

ma rozkład taki sam, jak ciąg statystyk pozycyjnych

$$U_{1:n}, \dots, U_{n:n}$$

z rozkładu $U(0, t)$.

Dowód. Z Wniosku 5.1 wynika, że warunkowo, dla $T_n = s$,

$$\left(\frac{W_1}{s}, \dots, \frac{W_n}{s} \right) \sim \text{Dir}(\underbrace{1, \dots, 1}_n).$$

Innymi słowy wektor losowy (T_1, \dots, T_{n-1}) ma taki rozkład, jak $n - 1$ statystyk pozycyjnych z $U(0, s)$.

Obliczmy warunkową gęstość zmiennej losowej T_n , jeśli $N(t) = n$:

$$\begin{aligned} f_{T_n}(s | N(t) = n) &= \frac{\mathbb{P}(N(t) = n | T_n = s) f_{T_n}(s)}{\mathbb{P}(N(t) = n)} \\ &= \frac{\mathbb{P}(W_{n+1} > t - s) f_{T_n}(s)}{\mathbb{P}(N = n)} \\ &= \frac{e^{-\lambda(t-s)} (\lambda^n / (n-1)!) s^{n-1} e^{-\lambda s}}{e^{-\lambda t} (\lambda t)^n / n!} \\ &= \frac{n s^{n-1}}{t^n}. \end{aligned}$$

A zatem warunkowo, dla $N(t) = n$, zmienna losowa T_n/t ma rozkład Beta($n, 1$) i w konsekwencji (przypomnijmy sobie algorytm generowania zmiennych o rozkładzie Dirichleta)

$$\left(\frac{W_1}{t}, \dots, \frac{W_n}{t}, \frac{t - T_n}{t}\right) \sim \text{Dir}(\underbrace{1, \dots, 1}_{n+1}).$$

Innymi słowy wektor losowy $(T_1, \dots, T_{n-1}, T_n)$ ma taki rozkład, jak n statystyk pozycyjnych z $U(0, t)$. \square

Wynika stąd następujący sposób generowania procesu Poissona na przedziale $[0, t]$.

Listing.

```

Gen  $N \sim \text{Pois}(\lambda t)$ 
for  $i = 1$  to  $N$  do Gen  $U_i \sim U(0, t)$ ;
Sort  $(U_1, \dots, U_N)$ 
 $(T_1, \dots, T_N) := (U_{1:N}, \dots, U_{N:N})$ 

```

Co ważniejsze, Stwierdzenia 6.1, 6.2 i 6.3 wskazują, jak powinny wyglądać uogólnienia procesu Poissona i jak generować takie ogólniejsze procesy. Zanim tym się zajmiemy, zrobmy dygresję i podajmy twierdzenie charakteryzujące „zwykły” proces Poissona. Nie jest ono bezpośrednio używane w symulacjach, ale wprowadza pewien ważny sposób określania procesów, który ułatwia zrozumienie łańcuchów Markowa z czasem ciągłym i jest bardzo użyteczny w probabilistycznym modelowaniu zjawisk.

Twierdzenie 6.1. *Załóżmy, że $(N(t) : t \geq 0)$ jest procesem o wartościach w $\{0, 1, 2, \dots\}$, stacjonarnych i niezależnych przyrostach (to znaczy $N(t) - N(s)$ jest niezależne od $(N(u), u \leq s)$ i ma rozkład zależny tylko od $t - s$ dla dowolnych $0 < s < t$) oraz, że trajektorie $N(t)$ są prawostronnie ciągłymi funkcjami mającymi lewostronne granice (prawie na pewno). Jeżeli $N(0) = 0$ i spełnione są następujące warunki:*

$$(i) \quad \lim_{t \rightarrow 0} \frac{\mathbb{P}(N(t) = 1)}{t} = \lambda,$$

$$(ii) \quad \lim_{t \rightarrow 0} \frac{\mathbb{P}(N(t) \geq 2)}{t} = 0,$$

to $N(\cdot)$ jest jednorodnym procesem Poissona z intensywnością λ

Bardzo prosto można zauważyć, że proces Poissona $(N(t) : t \geq 0)$ o intensywności λ ma własności wymienione w Twierdzeniu 6.1. Ciekawe jest, że te własności w pełni charakteryzują proces Poissona.

Dowód Tw. 6.1 – szkic. Pokażemy tylko, że

$$p_n(t) := \mathbb{P}(N(t) = n) = e^{-\lambda t} \frac{(\lambda t)^n}{n!}.$$

Najpierw zajmiemy się funkcją $p_0(t) = \mathbb{P}(N(t) = 0)$. Z niezależności i jednorodności przyrostów wynika tożsamość

$$p_0(t + h) = p_0(h)p_0(t).$$

Stąd

$$\frac{p_0(t+h) - p_0(t)}{h} = \frac{p_0(h) - 1}{h} p_0(t) = \left[-\frac{p_1(h)}{h} - \frac{\sum_{i=2}^{\infty} p_i(h)}{h} \right] p_0(h).$$

Przejdźmy do granicy z $h \rightarrow 0$ i skorzystajmy z własności (i) i (ii). Dostajemy proste równanie różniczkowe:

$$p'_0(t) = -\lambda p_0(t).$$

Rozwiązanie tego równania z warunkiem początkowym $p_0(0) = 1$ jest funkcja

$$p_0(t) = e^{-\lambda t}.$$

Bardzo podobnie obliczamy kolejne funkcje p_n . Postępujemy rekurencyjnie: zakładamy, że znamy p_{n-1} i układamy równanie różniczkowe dla funkcji p_n . Podobnie jak poprzednio,

$$p_n(t+h) = p_n(t)p_0(h) + p_{n-1}(t)p_1(h) + \sum_{i=2}^n p_{n-1}(t)p_i(h),$$

a zatem

$$\frac{p_n(t+h) - p_n(t)}{h} = \frac{p_0(h) - 1}{h} p_n(t) + \frac{p_1(h)}{h} p_{n-1}(t) + \frac{1}{h} \sum_{i=2}^n p_{n-1}(t)p_i(h).$$

Korzystając z własności (i) i (ii) otrzymujemy równanie

$$p'_n(t) = -\lambda p_n(t) + \lambda p_{n-1}(t).$$

To równanie można rozwiązać metodą uźmiennienia stałej: poszukujemy rozwiązania postaci $p_n(t) = c(t)e^{-\lambda t}$. Zakładamy przy tym indukcyjnie, że $p_{n-1}(t) = (\lambda t)^{n-1}e^{-\lambda t}/(n-1)!$ i mamy oczywisty warunek początkowy $p_n(0) = 0$. Stąd już łatwo dostać dowodzony wzór na p_n .

Na koniec zauważmy, że z postaci funkcji p_0 łatwo wywnioskować jaki ma rozkład zmienna $T_1 = \inf\{t : N(t) > 0\}$. Istotnie, $\mathbb{P}(T_1 > t) = \mathbb{P}(N(t) = 0) = p_0(t) = e^{-\lambda t}$. \square

Własności (i) i (ii), w połączeniu z jednorodnością przyrostów można przepisać w następującej sugestywnej formie:

$$\begin{aligned} \mathbb{P}(N(t+h) = n+1 | N(t) = n) &= \lambda h + o(h), \\ \mathbb{P}(N(t+h) = n | N(t) = n) &= 1 - \lambda h + o(h), \quad h \searrow 0. \end{aligned} \tag{6.3}$$

Jest to o tyle ważne, że w podobnym języku łatwo formułować założenia o ogólniejszych procesach, na przykład tak zwanych procesach urodzin i śmierci. Wrócimy do tego przy okazji omawiania łańcuchów Markowa.

6.2.2. Niejednorodne procesy Poissona w przestrzeni

Naturalne uogólnienia procesu Poissona polegają na tym, że rozważa się losowe zbiory punktów w przestrzeni o dowolnym wymiarze i dopuszcza się różną intensywność pojawiania się punktów w różnych rejonach przestrzeni. Niech \mathcal{X} będzie przestrzenią polską. Czytelnik, który nie lubi abstrakcji może założyć, że $\mathcal{X} \subseteq \mathbb{R}^d$.

Musimy najpierw wprowadzić odpowiednie oznaczenia. Rozważmy ciąg wektorów losowych w \mathcal{X} :

$$X_1, \dots, X_n, \dots$$

(może to być ciąg skończony lub nie, liczba tych wektorów może być zmienną losową). Niech, dla $A \subseteq \mathcal{X}$,

$$N(A) = \#\{i : X_i \in A\}$$

oznacza liczbę wektorów, które „wpadły do zbioru A ” (przy tym dopuszczamy wartość $N(A) = \infty$ i umawiamy się liczyć powtarzające się wektory tyle razy, ile razy występują w ciągu).

Niech teraz μ będzie miarą na (σ -ciele borelowskich podzbiorów) przestrzeni \mathcal{X} .

Definicja 6.4. $N(\cdot)$ jest procesem Poissona z miarą intensywności μ , jeśli

- dla parami rozłącznych zbiorów $A_1, \dots, A_i \subseteq \mathcal{X}$, odpowiadające im zmienne losowe $N(A_1), \dots, N(A_i)$ są niezależne;
- $\mathbb{P}(N(A) = n) = e^{-\mu(A)} \frac{\mu(A)^n}{n!}$, dla każdego $A \subseteq \mathcal{X}$ takiego, że $\mu(A) < \infty$ i dla $n = 0, 1, \dots$

Z elementarnych własności rozkładu Poissona wynika następujący wniosek.

Wniosek 6.1. Rozważmy rozbitcie zbioru A o skończonej mierze intensywności, $\mu(A) < \infty$, na rozłączną sumę $A = A_1 \cup \dots \cup A_k$. Wtedy

$$\begin{aligned} \mathbb{P}(N(A_1) = n_1, \dots, N(A_k) = n_k | N(A) = n) \\ = \frac{n!}{n_1! \dots n_k!} \mu(A_1)^{n_1} \dots \mu(A_k)^{n_k}, \quad (n_1 + \dots + n_k = n). \end{aligned}$$

Jeśli natomiast $\mu(A) = \infty$ to łatwo zauważyć, że $N(A) = \infty$ z prawdopodobieństwem 1.

Z Definicji 6.4 i Wniosku 6.1 natychmiast wynika następujący algorytm generowania procesu Poissona. Załóżmy, że interesuje nas „fragment” procesu w zbiorze $A \subseteq \mathcal{X}$ o skończonej mierze intensywności. W praktyce zawsze symulacje muszą się do takiego fragmentu ograniczać. Zauważmy, że *unormowana* miara $\mu(\cdot)/\mu(A)$ jest rozkładem prawdopodobieństwa (zmienna losowa X ma ten rozkład, jeśli $\mathbb{P}(X \in B) = \mu(B)/\mu(A)$, dla $B \subseteq A$).

Listing.

```
Gen  $N \sim \text{Poiss}(\mu(A))$ 
for  $i = 1$  to  $N$  do Gen  $X_i \sim \mu(\cdot)/\mu(A)$ 
```

Należy rozumieć, że formalnie definiujemy $N(B) = \#\{i : X_i \in B\}$ dla $B \subseteq A$, w istocie jednak za realizację procesu uważamy zbiór punktów $\{X_1, \dots, X_N\} \subset A$ („zapominamy” o uporządkowaniu punktów). Widać, że to jest proste uogólnienie analogicznego algorytmu dla „zwykłego” procesu Poissona, podanego wcześniej.

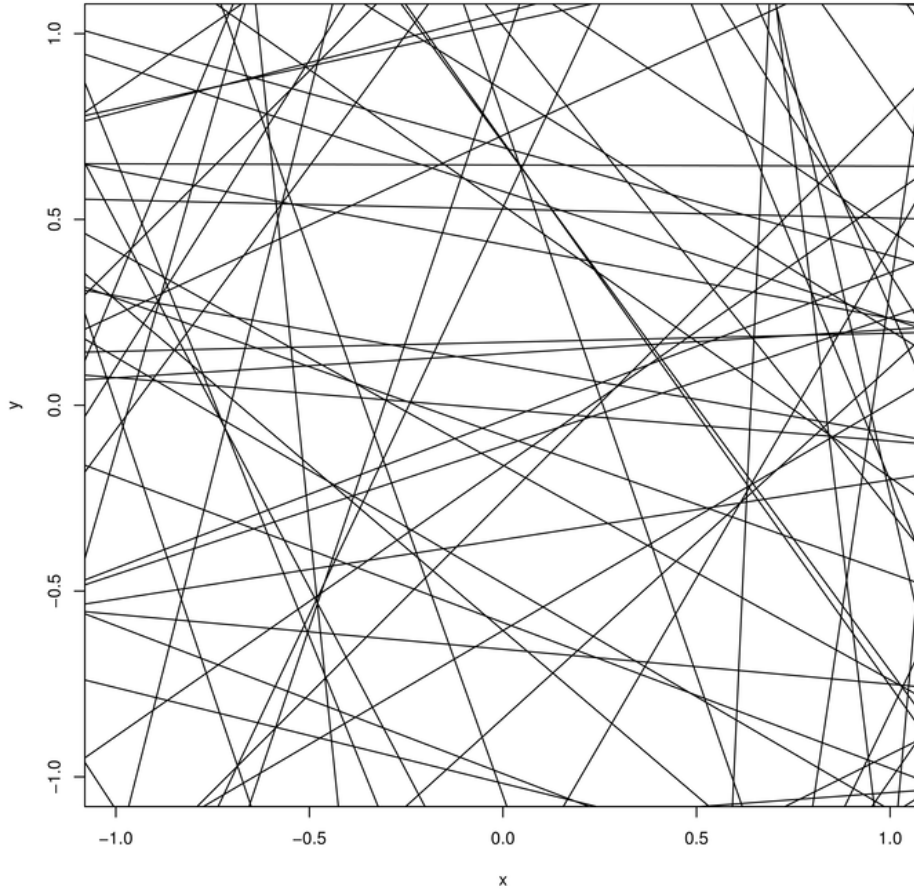
Można również rozbić zbiór A na rozłączną sumę $A = A_1 \cup \dots \cup A_k$ i generować *niezależnie* fragmenty procesu w każdej części A_j .

Przykład 6.1. Jednorodny proces Poissona na kole $B^2 = \{x^2 + y^2 \leq 1\}$ można wygenerować w następujący sposób. Powiedzmy, że intensywność (punktów na jednostkę pola) jest λ , to znaczy $\mu(B) = \lambda|B|$, gdzie $|B|$ jest polem (miarą Lebesgue’a) zbioru B .

Najpierw generujemy $N \sim \text{Poiss}(\lambda\pi)$, następnie punkty $(X_1, Y_1), \dots, (X_N, Y_N)$ niezależnie z rozkładu $U(B^2)$. To wszystko.

Ciekawszy jest następny przykład.

Przykład 6.2 (Jednorodny proces Poissona w przestrzeni prostych). Prostą na płaszczyźnie można sparametryzować podając kąt $\varphi \in [0, 2\pi[$ jako tworzy prostopadła do prostej z osią poziomą oraz odległość $r \geq 0$ prostej od początku układu. Każda prosta jest więc opisana przez parę liczb (φ, r) czyli punkt przestrzeni $\mathcal{L} = [0, 2\pi[\times [0, \infty[$. Jeśli teraz wygenerujemy jednorodny proces Poissona na tej przestrzeni, to znaczy proces o intensywności $\mu(B) = \lambda|B|$, $B \subseteq \mathcal{L}$, to można się spodziewać zbioru „losowo położonych prostych”. To widać na Rysunku 6.1. W istocie, wybór parametryzacji zapewnia, że rozkład prawdopodobieństwa procesu nie zależy od wyboru układu współrzędnych. Dowód, czy nawet precyzyjne sformułowanie tego stwierdzenia przekracza ramy tego skryptu. Intuicyjnie chodzi o to, że na obrazku „nie ma wyróżnionego



Rysunek 6.1. Proces Poissona w przestrzeni prostych.

kierunku ani wyróżnionego punktu”. Nie można sensownie zdefiniować pojęcia „losowej prostej” ale każdy przyzna, że proces Poissona o którym można uznać za uściślenie intuicyjnie rozumianego pojęcia „losowego zbioru prostych”.

Ciekawe, że podstawowe metody generowania zmiennych losowych mają swoje odpowiedniki dla procesów Poissona. Rolę rozkładu prawdopodobieństwa przejmuje miara intensywności.

Przykład 6.3 (Odwracanie dystrybuanty). Dla miary intensywności λ na przestrzeni jednowymiarowej można zdefiniować *dystrybuantę* tej miary. Dla uproszczenia rozważmy przestrzeń $\mathcal{X} = [0, \infty[$ i założymy, że każdy zbiór ograniczony ma miarę skończoną. Niech $\Lambda(t) = \lambda([0, t])$. Funkcję $\Lambda : [0, \infty[\rightarrow [0, \infty[$ nazwiemy dystrybuantą. Jest ona niemalejąca, prawostronnie ciągła, ale granica w nieskończoności $\Lambda(\infty) = \lim_{x \rightarrow \infty} \Lambda(x)$ może być dowolnym elementem z $[0, \infty]$. Dla procesu Poissona na $[0, \infty[$ wygodnie wrócić do prostszych oznaczeń, pisząc $N(t) = N([0, t])$ jak we wcześniej rozpatrywanym przypadku jednorodnym.

Niech $J(t)$ będzie jednorodnym procesem Poissona na $[0, \infty[$ z intensywnością równą 1. Wtedy

$$N(t) = J(\Lambda(t))$$

jest niejednorodnym procesem Poissona z miarą intensywności λ . W istocie, jeśli $0 < t_1 < t_2 < \dots$ to $N(t_1), N(t_2) - N(t_1), \dots$ są niezależne i mają rozkłady odpowiednio $\text{Poiss}(\Lambda(t_1)), \text{Poiss}(\Lambda(t_2) - \Lambda(t_1)), \dots$

$\Lambda(t_1), \dots$. Zauważmy, że jeśli $R_1 < R_2 < \dots$ oznaczają punkty skoku procesu $J(\cdot)$ to $N(t) = \max\{k : R_k \leq \Lambda(t)\} = \max\{k : \Lambda^-(R_k) \leq t\}$, gdzie Λ^- jest uogólnioną funkcją odwrotną do dystrybucyj. Wobec tego punktami skoku procesu N są $T_k = \Lambda^-(R_k)$. Algorytm jest taki:

Listing.

```
Gen  $N \sim \text{Poiss}(\Lambda(t))$ ;
for  $i := 1$  to  $N$  do Gen  $R_i \sim U(0, \Lambda(t))$ ;  $T_i := \Lambda^-(R_i)$ 
```

Pominęliśmy tu sortowanie skoków i założyliśmy, że symulacje ograniczamy do odcinka $[0, t]$.

Przykład 6.4 (Przerzedzanie). To jest odpowiednik metody eliminacji. Załóżmy, że mamy dwie miary intensywności: μ o gęstości m i λ o gęstości l . To znaczy, że $\lambda(B) = \int_B l(x)dx$ i $\mu(B) = \int_B m(x)dx$ dla dowolnego zbioru $B \subseteq \mathcal{X}$. Załóżmy, że $l(x) \leq m(x)$ i przypuśćmy, że umiemy generować proces Poissona o intensywności μ . Niech X_1, \dots, X_N będą punktami tego procesu w zbiorze A o skończonej mierze μ (wiemy, że $N \sim \text{Poiss}(\mu(A))$). Punkt X_i *akceptujemy* z prawdopodobieństwem $l(X_i)/m(X_i)$ (pozostawiamy w zbiorze) lub odrzucamy (usuwanie ze zbioru) z prawdopodobieństwem $1 - l(X_i)/m(X_i)$. Liczba pozostawionych punktów L ma rozkład $\text{Poiss}(\lambda(A))$, zaś każdy z tych punktów ma rozkład o gęstości $l(x)/\lambda(A)$, gdzie $\lambda(A) = \int_A l(x)dx$. Te punkty tworzą proces Poissona z miarą intensywności λ .

Listing.

```
Gen  $\mathbb{X} = \{X_1, \dots, X_N\} \sim \text{Poiss}(\mu(\cdot))$ ;
for  $i := 1$  to  $N$  Gen  $U_i \sim U(0, 1)$ ;
  if  $U_i > l(X_i)/m(X_i)$  then  $\mathbb{X} := \mathbb{X} \setminus X_i$ ;
return  $\mathbb{X} = \{X'_1, \dots, X'_L\} \sim \text{Poiss}(\lambda(\cdot))$ 
```

Przykład 6.5 (Superpozycja). To jest z kolei odpowiednik metody *kompozycji*. Metoda opiera się na następującym prostym fakcie. Jeżeli $N_1(\cdot), \dots, N_k(\cdot)$ są niezależnymi procesami Poissona z miarami intensywności odpowiednio $\mu_1(\cdot), \dots, \mu_k(\cdot)$, to $N(\cdot) = \sum_i N_i(\cdot)$ jest procesem Poissona z intensywnością $\mu(\cdot) = \sum_i \mu_i(\cdot)$. Dodawanie należy tu rozumieć w dosłowny sposób, to znaczy $N(A)$ jest określone jako $\sum_i N_i(A)$ dla każdego zbioru A . Jeśli utożsamimy procesy z losowymi zbiorami punktów to odpowiada temu operacja brania sumy mnogościowej (złączenia zbiorów). Niech $X_{j,1}, \dots, X_{j,N_j}$ będą punktami j -tego procesu w zbiorze A o skończonej mierze μ_j .

Listing.

```
 $\mathbb{X} = \emptyset$ ;
for  $j := 1$  to  $k$  do
  begin
    Gen  $\mathbb{X}_j = \{X_{j,1}, \dots, X_{j,N_j}\} \sim \text{Poiss}(\mu_j(\cdot))$ ;
     $\mathbb{X} := \mathbb{X} \cup \mathbb{X}_j$ ;
  end
return  $\mathbb{X} = \{X'_1, \dots, X'_N\} \sim \text{Poiss}(\mu(\cdot))$ 
  { mamy tu  $N = \sum_j N_j$  }
```

Wygodnie jest utożsamiać procesy Poissona z losowymi zbiorami punktów, jak uczyniliśmy w ostatnich przykładach (i mniej jawnie w wielu miejscach wcześniej). Te zbiory można rozumieć w zwykłym sensie, dodawać, odejmować tak jak w teorii mnogości pod warunkiem, że ich elementy się *nie powtarzają*. W praktyce mamy najczęściej do czynienia z intensywnościami, które mają gęstości „w zwykłym sensie”, czyli względem miary Lebesgue’a. Wtedy, z prawdopodobieństwem 1, punkty procesu Poissona nie powtarzają się.

7. Symulowanie procesów stochastycznych II.

Procesy Markowa

7.1. Czas dyskretny, przestrzeń dyskretna

Zacznijmy od najprostszej sytuacji. Rozważmy skończony zbiór \mathcal{X} , który będziemy nazywali *przestrzenią stanów*. Czasem wygodnie przyjąć, że $\mathcal{X} = \{1, 2, \dots, d\}$. Jest to tylko umowne ponumerowanie stanów.

Definicja 7.1. (i) Ciąg $X_0, X_1, \dots, X_n, \dots$ zmiennych losowych o wartościach w \mathcal{X} nazywamy **łańcuchem Markowa**, jeśli dla każdego $n = 1, 2, \dots$ i dla każdego ciągu $x_0, x_1, \dots, x_n, x_{n+1}$ punktów przestrzeni \mathcal{X} ,

$$\begin{aligned} \mathbb{P}(X_{n+1} = x_{n+1} | X_n = x_n, X_{n-1} = x_{n-1}, \dots, X_1 = x_1, X_0 = x_0) \\ = \mathbb{P}(X_{n+1} = x_{n+1} | X_n = x_n), \end{aligned}$$

(o ile tylko $\mathbb{P}(X_n = x_n, X_{n-1} = x_{n-1}, \dots, X_1 = x_1, X_0 = x_0) > 0$, czyli prawdopodobieństwo warunkowe w tym wzorze ma sens).

(ii) Łańcuch Markowa nazywamy **jednorodnym**, jeśli dla dowolnych stanów x i y i każdego n możemy napisać

$$\mathbb{P}(X_{n+1} = y | X_n = x) = P(x, y),$$

to znaczy prawdopodobieństwo warunkowe w powyższym wzorze zależy tylko od x i y , ale nie zależy od n .

Jeśli $X_n = x$, to mówimy, że łańcuch w chwili n znajduje się w stanie $x \in \mathcal{X}$. Warunek (i) w Definicji 7.1 znaczy tyle, że przyszła ewolucja łańcucha zależy od stanu obecnego, ale nie zależy od przeszłości. Łańcuch jest jednorodny, jeśli prawo ewolucji łańcucha nie zmienia się w czasie. W dalszym ciągu rozpatrywać będziemy głównie łańcuchy jednorodne. Macierz

$$P = (P(x, y))_{x, y \in \mathcal{X}} = \begin{pmatrix} P(1, 1) & \cdots & P(1, y) & \cdots & P(1, d) \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ P(x, 1) & \cdots & P(x, y) & \cdots & P(x, d) \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ P(d, 1) & \cdots & P(d, y) & \cdots & P(d, d) \end{pmatrix}$$

nazywamy *macierzą* (prawdopodobieństw) *przejścia* łańcucha. Jest to macierz stochastyczna, to znaczy $P(x, y) \geq 0$ dla dowolnych stanów $x, y \in \mathcal{X}$ oraz $\sum_y P(x, y) = 1$ dla każdego $x \in \mathcal{X}$. Jeśli $q(x) = \mathbb{P}(X_0 = x)$, to wektor wierszowy

$$q^\top = (q(1), \dots, q(x), \dots, q(d))$$

nazywamy *rozkładem początkowym* łańcucha (oczywiście, $\sum_x q(x) = 1$). Jest jasne, że

$$\begin{aligned} \mathbb{P}(X_0 = x_0, X_1 = x_1, \dots, X_{n-1} = x_{n-1}, X_n = x_n) \\ = q(x_0)P(x_0, x_1) \cdots P(x_{n-1}, x_n). \end{aligned} \tag{7.1}$$

Ten wzór określa jednoznacznie łączny rozkład prawdopodobieństwa zmiennych $X_0, X_1, \dots, X_n, \dots$ i może być przyjęty za definicję jednorodnego łańcucha Markowa. W tym sensie możemy utożsamiać łańcuch z parą (q, P) : łączny rozkład prawdopodobieństwa jest wyznaczony przez podanie rozkładu początkowego i macierzy przejścia.

Opiszemy teraz bardzo ogólną konstrukcję, która jest podstawą algorytmów generujących łańcuchy Markowa. Wyobraźmy sobie, jak zwykle, że mamy do dyspozycji ciąg $U = U_0, U_1, \dots, U_n, \dots$ „liczb losowych”, produkowanych przez komputerowy generator, czyli z teoretycznego punktu widzenia ciąg *niezależnych zmiennych losowych o jednakowym rozkładzie*, $U(0, 1)$. Niech $\phi : [0, 1] \rightarrow \mathcal{X}$ i $\psi : \mathcal{X} \times [0, 1] \rightarrow \mathcal{X}$ będą takimi funkcjami, że

$$\begin{aligned}\mathbb{P}(\psi(U) = x) &= q(x) \quad \text{dla każdego } x \in \mathcal{X}, \\ \mathbb{P}(\phi(x, U) = y) &= P(x, y) \quad \text{dla dowolnych } x, y \in \mathcal{X}.\end{aligned}\tag{7.2}$$

Jeśli określimy $X_0, X_1, \dots, X_n, \dots$ rekurencyjnie w następujący sposób:

$$X_0 = \psi(U_0), \quad X_{n+1} = \phi(X_n, U_{n+1}),\tag{7.3}$$

to jest jasne, że tak otrzymany ciąg zmiennych losowych jest łańcuchem Markowa z rozkładem początkowym q i macierzą przejścia P .

Na zakończenie tego podrozdziału zrobmy kilka prostych uwag.

- Wszystko, co w tym podrozdziale zostało powiedziane można niemal mechanicznie uogólnić na przypadek *nieskończonej ale przeliczalnej* przestrzeni \mathcal{X} .
- Można sobie wyobrażać, że dla ustalonego x , zbiór $\{u : \phi(x, u) = y\}$ jest *odcinkiem* długości $P(x, y)$, ale nie musimy tego żądać. Nie jest istotne, że zmienne U_i mają rozkład $U(0, 1)$. Moglibyśmy założyć, że są określone na dość dowolnej przestrzeni \mathcal{U} . Istotne jest, że te zmienne są niezależne, mają jednakowy rozkład i zachodzi wzór (7.2). W praktyce bardzo często do zrealizowania jednego „kroku” łańcucha Markowa generujemy więcej niż jedną „liczbę losową”.

7.2. Czas dyskretny, przestrzeń ciągła

Przypadek ciągłej lub raczej *ogólnej* przestrzeni stanów jest w istocie równie prosty, tylko oznaczenia trzeba trochę zmienić i sumy zamienić na całki. Dla prostoty przyjmijmy, że $\mathcal{X} \subseteq \mathbb{R}^d$. Poniżej sformułujemy definicję w taki sposób, żeby podkreślić analogię do przypadku dyskretnego i pominiemy subtelności teoretyczne. Zwróćmy tylko uwagę, że prawdopodobieństwo warunkowe nie może tu być zdefiniowane tak elementarnie jak w Definicji 7.1, bo zdarzenie warunkujące może mieć prawdopodobieństwo zero.

Definicja 7.2. (i) Ciąg $X_0, X_1, \dots, X_n, \dots$ zmiennych losowych o wartościach w \mathcal{X} nazywamy **łańcuchem Markowa**, jeśli dla każdego $n = 1, 2, \dots$, dla każdego ciągu x_0, x_1, \dots, x_n punktów przestrzeni \mathcal{X} oraz dowolnego (borelowskiego) zbioru $B \subseteq \mathcal{X}$,

$$\begin{aligned}\mathbb{P}(X_{n+1} \in B | X_n = x_n, X_{n-1} = x_{n-1}, \dots, X_1 = x_1, X_0 = x_0) \\ = \mathbb{P}(X_{n+1} \in B | X_n = x_n),\end{aligned}$$

(dla prawie wszystkich punktów $(x_n, x_{n-1}, \dots, x_1, x_0)$).

(ii) Łańcuch Markowa nazywamy **jednorodnym**, jeśli dla dowolnego zbioru $B \subseteq \mathcal{X}$ i każdego n możemy napisać

$$\mathbb{P}(X_{n+1} \in B | X_n = x) = P(x, B)$$

(dla prawie wszystkich punktów x).

Funkcja $P(\cdot, \cdot)$, której argumentami są punkt x i zbiór B , jest nazywana *jądrem przejścia*. Ważne dla nas jest to, że dla ustalonego $x \in \mathcal{X}$, jądro $P(x, \cdot)$ rozważane jako funkcja zbioru jest rozkładem prawdopodobieństwa. W wielu przykładach jest to rozkład zadany przez gęstość,

$$P(x, B) = \int_B p(x, y) dy.$$

Wtedy odpowiednikiem wzoru (7.1) jest następujący wzór na łączną gęstość:

$$p(x_0, x_1, \dots, x_{n-1}, x_n) = q(x_0)p(x_0, x_1) \cdots p(x_{n-1}, x_n), \quad (7.4)$$

gdzie $q(\cdot)$ jest gęstością rozkładu początkowego, zaś $p(x_i, x_{i+1}) = p(x_{i+1}|x_i)$ jest gęstością warunkową. Sformułowaliśmy Definicję 7.2 nieco ogólniej głównie dlatego, że w symulacjach ważną rolę odgrywają łańcuchy dla których prawdopodobieństwo przejścia nie ma gęstości. Przykładem może być łańcuch, który z niezerowym prawdopodobieństwem potrafi „stać w miejscu”, to znaczy $P(x, \{x\}) =: \alpha(x) > 0$, i ma prawdopodobieństwo przejścia postaci, powiedzmy

$$P(x, B) = \alpha(x)\mathbb{I}(x \in B) + (1 - \alpha(x)) \int_B p(x, y) dy.$$

Łańcuch Markowa na ogólnej przestrzeni stanów można rekurencyjnie generować w sposób opisany wzorem (7.3), to znaczy $X_0 = \psi(U_0)$ i $X_{n+1} = \phi(X_n, U_{n+1})$ dla ciągu i.i.d. $U_0, U_1, \dots, U_n, \dots$ z rozkładu $U(0, 1)$. Niech $Q(\cdot)$ oznacza rozkład początkowy. Funkcje ψ i ϕ muszą spełniać warunki

$$\begin{aligned} \mathbb{P}(\psi(U) \in B) &= Q(B) \quad \text{dla każdego (mierzalnego) } B \subseteq \mathcal{X}, \\ \mathbb{P}(\phi(x, U) \in B) &= P(x, B) \quad \text{dla dowolnych } x \in \mathcal{X}, B \subseteq \mathcal{X}. \end{aligned} \quad (7.5)$$

Różnica między (7.2) i (7.5) jest tylko taka, że w ogólnej sytuacji rozkłady prawdopodobieństwa zmiennych losowych $\psi(U)$ i $\phi(x, U)$, czyli odpowiednio $Q(\cdot)$ i $P(x, \cdot)$ nie muszą być dyskretne. Ale można zastosować w zasadzie każdą metodę generacji zmiennych o zadanym rozkładzie, przy czym ten zadany rozkład zależy od x .

7.3. Czas ciągły, przestrzeń dyskretna

Rozważmy proces stochastyczny $(X(t), t \geq 0)$, czyli rodzinę zmiennych losowych o wartościach w skończonej przestrzeni stanów \mathcal{X} , indeksowaną przez parameter t , który interpretujemy jako *ciągły czas*. Pojedyncze stany oznaczamy przez $x, y, z, \dots \in \mathcal{X}$ lub podobnie. Załóżmy, że trajektorie są prawostronnie ciągłymi funkcjami mającymi lewostronne granice (prawie na pewno).

Definicja 7.3. (i) Mówimy, że $X(t)$ jest **procesem Markowa**, jeśli dla dowolnych $x, y \in \mathcal{X}$ oraz $s, t \geq 0$,

$$\mathbb{P}(X(s+t) = y | X(s) = x, X(u), 0 \leq u < s) = \mathbb{P}(X(t+s) = y | X(s) = x).$$

(ii) Proces Markowa nazywamy **jednorodnym**, jeśli dla dowolnych stanów x i y i każdego t i s możemy napisać

$$\mathbb{P}(X(s+t) = y | X(s) = x) = P^t(x, y),$$

to znaczy prawdopodobieństwo warunkowe w powyższym wzorze zależy tylko od x, y i t , ale nie zależy od s .

Sens założenia (i) jest taki sam jak w Definicji 7.1. Żądamy mianowicie, żeby zmienna $X(t+s)$ była *warunowo niezależna* od zmiennych $X(u)$, $0 \leq u < s$ pod warunkiem zdarzenia $X(s) = x$. Oczywiście, $P^t(x, y)$ jest *prawdopodobieństwem przejścia w ciągu czasu t* . Rozkład początkowy oznaczmy przez $q(x) = \mathbb{P}(X(0) = x)$. Ponieważ \mathcal{X} jest skończona, q może być przedstawiony jako wektor a P^t jako macierz.

Odpowiednikiem macierzy przejścia jest *macierz intensywności przejść* zdefiniowana następująco.

$$Q(x, y) = \lim_{h \rightarrow 0} \frac{1}{h} [P^h(x, y) - I(x, y)],$$

gdzie $I = P^0$ jest macierzą identycznościową. Można udowodnić, że granice istnieją. Tak więc $Q = \frac{d}{dt} P^t|_{t=0}$ i $P^h = I + hQ + o(h)$ przy $h \searrow 0$. Innymi słowy, $Q(x, y)$ jest rzeczywiście „intensywnością skoków z x do y ” (na jednostkę czasu). Mamy wzór analogiczny do (6.3):

$$\begin{aligned} \mathbb{P}(X(t+h) = y | X(t) = x) &= hQ(x, y) + o(h) \text{ dla } x \neq y; \\ \mathbb{P}(X(t+h) = x | X(t) = x) &= 1 + hQ(x, x) + o(h), \quad h \searrow 0. \end{aligned}$$

Zauważmy, że mamy $\sum_y Q(x, y) = 0$. Wspomnijmy jeszcze mimochodem, że macierze przejścia wyrażają się przez macierz intensywności Q przy pomocy „macierzowej funkcji wykładniczej”: $P^t = \exp[tQ]$. Dla uproszczenia notacji, niech $Q(x) = -Q(x, x) = \sum_{y \neq x} Q(x, y)$ oznacza „intensywność wszystkich skoków ze stanu x ”.

Z tego co zostało powiedziane łatwo się domyślić, że generowanie procesu Markowa można zorganizować podobnie jak dla procesu Poissona, poprzez *czasy skoków*. Niech $0 < T_1 < T_2 < \dots < T_n < \dots$ będą kolejnymi momentami skoków,

$$\begin{aligned} T_1 &= \inf \{t > 0 : X(t) \neq X(0)\}; \\ T_{n+1} &= \inf \{t > T_n : X(t) \neq X(T_n)\}. \end{aligned}$$

Jeśli $X(T_n) = x$, to czas oczekiwania na następny skok, $W_{n+1} = T_{n+1} - T_n$ zależy tylko od x i ma rozkład *wykładniczy* z parametrem $Q(x)$. W szczególności, mamy $\mathbb{E}(W_{n+1} | X(T_n) = x) = 1/Q(x)$. Jeśli obserwujemy proces w momentach skoków, czyli skupimy uwagę na ciągu $\hat{X}_n = X(T_n)$ to otrzymujemy łańcuch Markowa z czasem dyskretnym, nazywany czasem „szkieletem”. Jego prawdopodobieństwa przejścia są takie:

$$\hat{P}(x, y) = \mathbb{P}(\hat{X}_{n+1} = y | \hat{X}_n = x) = \begin{cases} Q(x, y)/Q(x) & \text{jeśli } x \neq y; \\ 0 & \text{jeśli } x = y. \end{cases} \quad (7.6)$$

Oczywiście, cała trajektoria procesu $X(t)$ jest w pełni wyznaczona przez momenty skoków T_1, \dots, T_n, \dots i kolejne stany łańcucha szkieletowego chain $X(0) = \hat{X}_0, \hat{X}_1, \dots, \hat{X}_n, \dots$. Po prostu, $X(t) = \hat{X}_n$ dla $T_n \leq t < T_{n+1}$. Następujący algorytm formalizuje opisany powyżej sposób generacji.

Listing.

```

Gen  $\hat{X}_0 \sim q(\cdot)$ ;
 $T_0 := 0$ ;
for  $i := 1$  to  $\infty$  do
  begin
    Gen  $W_i \sim \text{Ex}(Q(\hat{X}_{i-1}))$ ;
     $T_i := T_{i-1} + W_i$ ;
    Gen  $\hat{X}_i \sim \hat{P}(\hat{X}_{i-1}, \cdot)$ ; {  $\hat{P}$  dane wzorem powyżej }
  end

```

Rzecz jasna, w praktyce trzeba ustalić sobie skończony horyzont czasowy h i przerwać symulację gdy $T_i > h$. Zauważmy, że ten algorytm bez żadnych modyfikacji nadaje się do symulowania procesu Markowa na *nieskończonej* ale *przeliczalnej* przestrzeni stanów, na przykład $\mathcal{X} = \{0, 1, 2, \dots\}$.

Czasami warto zmodyfikować algorytm w następujący sposób. Generowanie procesu Markowa można rozbić na dwie fazy: najpierw wygenerować „potencjalne czasy skoków” a następnie symulować „szkielet”, który będzie skakał wyłącznie w poprzednio otrzymanych momentach (ale nie koniecznie we wszystkich spośród nich). Opiszemy dokładniej tę konstrukcję, która bazuje na własnościach procesu Poissona. Wyobraźmy sobie najpierw, że dla każdego stanu $x \in \mathcal{X}$ symulujemy niezależnie jednorodny proces Poissona \mathbb{R}^x o punktach skoku $R_1^x < \dots < R_k^x < \dots$ przy czym ten proces ma intensywność $Q(x)$. Jeśli teraz, w drugiej fazie $\hat{X}_{i-1} = X(T_{i-1}) = x$, to za następny moment skoku wybierzemy najbliższy punkt procesu \mathbb{R}^x na prawo od T_{i-1} , czyli $T_i = \min\{R_k^x : R_k^x > T_{i-1}\}$. Z własności procesu Poissona (patrz Stwierdzenie 6.2 i jego dowód) wynika, że zmienna $T_i - T_{i-1}$ ma rozkład wykładniczy z parametrem $Q(x)$ i metoda jest poprawna. Naprawdę nie ma nawet potrzeby generowania wszystkich procesów Poissona \mathbb{R}^x . Warto posłużyć się metodą *przerzedzania* i najpierw wygenerować jeden proces o dostatecznie dużej intensywności a następnie „w miarę potrzeby” go przerzedzać. Niech \mathbb{R}^* będzie procesem Poissona o intensywności $Q^* \geq \max_x Q(x)$ i oznaczmy jego punkty skoków przez $\{R_1 < \dots < R_k < \dots\}$. Jeśli w drugiej fazie algorytmu mamy $X(R_{i-1}) = x$ w pewnym momencie $R_{i-1} \in \mathbb{R}^*$, to zaczynamy przerzedzać proces \mathbb{R}^* w taki sposób, aby otrzymać proces o intensywności $Q(x) \leq Q^*$. Dla każdego punktu R_j , $j \geq i$ powinniśmy rzucić monetą i z prawdopodobieństwem $(1 - Q(x))/Q^*$ ten punkt usunąć. Ale jeśli usuniemy punkt R_i to znaczy, że w tym punkcie *nie ma skoku*, czyli $X(R_i) = x$. Jeśli punkt R_i zostawimy, to wykonujemy skok, losując kolejny stan z prawdopodobieństwem $\tilde{P}(x, \cdot)$. W rezultacie, stan $X(R_i)$ jest wylosowany z rozkładu prawdopodobieństwa

$$\tilde{P}(x, y) = \begin{cases} Q(x, y)/Q^* & \text{jeśli } x \neq y; \\ 1 - Q(x)/Q^* & \text{jeśli } x = y. \end{cases} \quad (7.7)$$

Niech $\tilde{X}_n = X(R_i)$ będzie „nadmiarowym szkieletem” procesu. Jest to łańcuch Markowa, który (w odróżnieniu od „cieńszego szkieletu” \hat{X}_n) *może* w jednym kroku pozostać w tym samym stanie i ma prawdopodobieństwa przejścia $\mathbb{P}(\tilde{X}_n = y | \tilde{X}_{n-1} = x) = \tilde{P}(x, y)$.

Odpowiada temu następujący dwufazowy algorytm.

— **Faza I**

Listing.

```
Gen  $\mathbb{R} \sim \text{Poiss}(Q^*) \{ \text{proces Poissona} \}$ 
```

— **Faza II**

Listing.

```
Gen  $\tilde{X}_0 \sim q(\cdot)$ ;  
for  $i := 1$  to  $\infty$  do Gen  $\tilde{X}_i \sim \tilde{P}(\hat{X}_{i-1}, \cdot)$ ; {  $\tilde{P}$  dane wzorem powyżej }
```

Algorytm jest bardzo podobny do metody wynalezionej przez Gillespie’go do symulowania wielowymiarowych procesów urodzin i śmierci (głównie w zastosowaniach do chemii). Te procesy mają najczęściej przestrzeń stanów postaci $\mathcal{X} = \{0, 1, 2, \dots\}^d$; stanem układu jest układ liczb $x = (x(1), \dots, x(d))$, gdzie $x(i)$ jest liczbą osobników (na przykład cząstek) typu i . Przestrzeń jest nieskończona, ale większa część rozważań przenosi się bez trudu. Pojawia się tylko problem

w ostatnim algorytmie jeśli $\max_x Q(x) = \infty$, czego nie można wykluczyć. Zeby sobie z tym poradzić, można wziąć Q^* dostatecznie duże, żeby „na ogół” wystarczyło, a w mało prawdopodobnym przypadku zawitania do stanu x z $Q(x) > Q^*$ przerzucać się na pierwszy, jednofazowy algorytm.

8. Algorytmy Monte Carlo I. Obliczanie całek

Duża część zastosowań metod MC wiąże się z obliczaniem całek lub sum. Typowe zadanie polega na obliczeniu wartości oczekiwanej

$$\theta = \mathbb{E}_\pi f(X) = \int_{\mathcal{X}} f(x) \pi(dx),$$

gdzie X jest zmienną losową o rozkładzie prawdopodobieństwa π na przestrzeni \mathcal{X} , zaś $f : \mathcal{X} \rightarrow \mathbb{R}$. Zazwyczaj \mathcal{X} jest podzbiorem wielowymiarowej przestrzeni euklidesowej lub jest zbiorem skończonym ale bardzo licznym (na przykład zbiorem pewnych obiektów kombinatorycznych). Jeśli $\mathcal{X} \subseteq \mathbb{R}^d$ i rozkład π jest opisany przez gęstość p to całka określająca wartość oczekiwaną jest zwykłą całką Lebesgue'a. W tym przypadku zatem możemy napisać

$$\theta = \int_{\mathcal{X}} f(x) p(x) dx.$$

Równie ważny w zastosowaniach jest przypadek *dyskretnej* przestrzeni \mathcal{X} . Jeśli $p(x) = \mathbb{P}(X = x)$, to

$$\theta = \sum_{x \in \mathcal{X}} f(x) p(x).$$

W dalszym ciągu tego rozdziału utożsamiamy rozkład π z funkcją p . Zwróćmy uwagę, że przedstawiamy tu θ jako wartość oczekiwaną tylko dla wygody oznaczeń; w istocie *każda* całka, suma, (a także prawdopodobieństwo zdarzenia losowego) jest wartością oczekiwaną.

Na pierwszy rzut oka nie widać czym polega problem! Sumowanie wykonuje każdy kalkulator, całki się sprawnie oblicza numerycznie. Ale nie zawsze. Metody MC przydają się w sytuacjach gdy spotyka się z następującymi trudnościami (lub przynajmniej którąś z nich).

- Przestrzeń \mathcal{X} jest ogromna. To znaczy wymiar d jest bardzo duży lub skończona przestrzeń zawiera astronomicznie dużą liczbę punktów.
- Rozkład prawdopodobieństwa π jest „skupiony” w małej części ogromnej przestrzeni \mathcal{X} .
- Nie ma podstaw do zakładania, że funkcja f jest w jakimkolwiek sensie „gładka” (co jest warunkiem stosowania standardowych numerycznych metod całkowania).
- Gęstość p rozkładu π jest znana tylko z dokładnością do stałej normującej. Innymi słowy, umiemy obliczać $p'(x) = Zp(x)$ ale nie znamy stałej $Z = \int p'(x) dx$. Celem zadanie polega właśnie na obliczeniu tej stałej (Z jest nazwane „funkcją podziału” lub sumą statystyczną).

Prosta metoda MC (po angielsku nazywana bardziej brutalnie: *Crude Monte Carlo*, czyli CMC) nasuwa się sama. Należy wygenerować n niezależnych zmiennych losowych X_1, \dots, X_n o jednakowym rozkładzie π i za estymator wartości oczekiwanej wziąć średnią z próbki,

$$\hat{\theta}_n = \hat{\theta}_n^{\text{CMC}} = \sum_{i=1}^n f(X_i).$$

Mocne Prawo Wielkich Liczb gwarantuje, że $\hat{\theta}_n \rightarrow \theta$ prawie na pewno, gdy $n \rightarrow \infty$. W terminologii statystycznej, $\hat{\theta}_n$ jest *mocno zgodnym* estymatorem obliczanej wielkości. Zadanie wydaje się rozwiązane. Są jednak dwa zasadnicze kłopoty.

- Estymator $\hat{\theta}_n$ skonstruowany metodą CMC może się zbliżać do θ przerażająco wolno.
- Co zrobić, jeśli nie umiemy losować z rozkładu π ?

8.1. Losowanie istotne

Zadziwiająco skutecznym sposobem na *oba* przedstawione wyżej kłopoty jest **losowanie istotne** (*Importance Sampling*, w skrócie IS). Przypuśćmy, że umiemy losować z rozkładu o gęstości q . Zauważmy, że

$$\theta = \int_{\mathcal{X}} \frac{p(x)}{q(x)} f(x) q(x) dx = \mathbb{E}_q \frac{p(X)}{q(X)} f(X) = \mathbb{E}_q w(X) f(X),$$

gdzie $w(x) = p(x)/q(x)$. Piszemy tu wzory dla całek ale oczywiście dla sum jest tak samo. Niech X_1, \dots, X_n będą niezależnymi zmiennymi losowymi o jednakowym rozkładzie g ,

$$\hat{\theta}_n = \hat{\theta}_n^{\text{IS1}} = \sum_{i=1}^n W_i f(X_i), \quad (8.1)$$

gdzie

$$W_i = w(X_i) = \frac{p(X_i)}{q(X_i)}$$

traktujemy jako *wagi* wylosowanych punktów X_i . Z tego co wyżej powiedzieliśmy wynika, że $\mathbb{E}_q \hat{\theta}_n = \theta$ oraz $\hat{\theta}_n \rightarrow \theta$ prawie na pewno, gdy $n \rightarrow \infty$. Mamy więc estymator *nieobciążony* i *zgodny*. Milcząco założyliśmy, że $q(X_i) > 0$ w każdym wylosowanym punkcie, czyli, że $\{x : p(x) > 0\} \subseteq \{x : q(x) > 0\}$. Z tym zwykle nie ma wielkiego kłopotu. Ponadto musimy założyć, że umiemy obliczać funkcję w w każdym wylosowanym punkcie. Jeśli znamy tylko $p'(x) = Zp(x)$ a nie znamy stałej Z to jest kłopot. Możemy tylko obliczyć $W'_i = p'(X_i)/q(X_i) = ZW_i$. Używamy zatem nieco innej postaci estymatora IS, mianowicie

$$\hat{\theta}_n = \hat{\theta}_n^{\text{IS2}} = \frac{\sum_{i=1}^n W_i f(X_i)}{\sum_{i=1}^n W_i}. \quad (8.2)$$

Ten estymator wygląda dziwnie, bo w mianowniku występuje *estymator jedynki* ale chodzi o to, że we wzorze (8.2) w przeciwieństwie do (8.1) możemy zastąpić W_i przez W'_i , bo nieznamy czynnik Z skraca się. Możemy jeszcze zapisać nasz estymator w zgrabnej formie

$$\hat{\theta}_n^{\text{IS2}} = \sum_{i=1}^n \tilde{W}_i f(X_i),$$

gdzie $\tilde{W}_i = W_i / \sum_j W_j$ są *unormowanymi wagami*. Trzeba jednak pamiętać, że to „unormowanie” polega na podzieleniu wag przez zmienną losową. W rezultacie estymator IS2 jest *obciążony*. Jest jednak mocno zgodny, bo licznik we wzorze (8.2) po podzieleniu przez n zmierza do θ a mianownik do 1, prawie na pewno. Estymator IS2 jest znacznie częściej używany niż IS1. Oprócz uniezależnienia się od stałej normującej, jest jeszcze inny powód. Okazuje się, że (pomimo obciążenia) estymator IS2 może być bardziej *efektywny*. Sens tego pojęcia wyjaśnimy w następnym podrozdziale. Zarówno dobór przykładów jak i ogólny tok rozważań jest zapożyczony z monografii Liu [15].

8.2. Efektywność estymatorów MC

Naturalne jest żądanie, aby estymator był mocno zgodny, $\hat{\theta}_n \rightarrow \theta$ prawie na pewno przy $n \rightarrow \infty$. Znaczy to, że zbliżymy się dowolnie blisko aproksymowanej wielkości, gdy tylko dostatecznie długo przedłużamy symulacje. Chcielibyśmy jednak wiedzieć coś konkretniejszego, oszacować jak szybko *błąd aproksymacji* $\hat{\theta}_n - \theta$ maleje do zera i jakie n wystarczy do osiągnięcia

wystarczającej dokładności. Przedstawimy teraz najprostsze i najczęściej stosowane w praktyce podejście, oparte na tzw. asymptotycznej wariancji i konstrukcji asymptotycznych przedziałów ufności. W typowej sytuacji estymator $\hat{\theta}_n$ ma następującą własność, nazywaną *asymptotyczną normalnością*:

$$\sqrt{n}(\hat{\theta}_n - \theta) \rightarrow N(0, \sigma^2), \quad (n \rightarrow \infty)$$

w sensie zbieżności według rozkładu. W konsekwencji, dla ustalonej liczby $a > 0$,

$$\mathbb{P}\left(|\hat{\theta}_n - \theta| \leq \frac{z\sigma}{\sqrt{n}}\right) \rightarrow \Phi(z) - \Phi(-z) = 2\Phi(z) - 1, \quad (n \rightarrow \infty),$$

gdzie Φ jest dystrybuantą rozkładu $N(0, 1)$. Często nie znamy *asymptotycznej wariancji* σ^2 ale umiemy skonstruować jej zgodny estymator $\hat{\sigma}_n^2$. Dla ustalonej „małej” dodatniej liczby α łatwo dobrać kwantyl rozkładu normalnego $z = z_{1-\alpha/2}$ tak, żeby $2\Phi(z) - 1 = 1 - \alpha$. Typowo $\alpha = 0.05$ i $z = z_{0.975} = 1.9600 \approx 2$. Otrzymujemy

$$\mathbb{P}\left(|\hat{\theta}_n - \theta| \leq \frac{z_{1-\alpha/2}\hat{\sigma}_n}{\sqrt{n}}\right) \rightarrow 1 - \alpha. \quad n \rightarrow \infty.$$

Ten wzór interpretuje się w praktyce tak: estymator $\hat{\theta}_n$ ma błąd nie przekraczający $z\hat{\sigma}_n/\sqrt{n}$ z prawdopodobieństwem około $1 - \alpha$. W żargonie statystycznym $1 - \alpha$ jest nazywane asymptotycznym poziomem ufności.

Chociaż opisane powyżej podejście ma swoje słabe strony, to prowadzi do prostego kryterium porównywania estymatorów. Przypuśćmy, że mamy dwa estymatory $\hat{\theta}_n^I$ i $\hat{\theta}_n^{II}$, oba asymptotycznie normalne, o asymptotycznej wariancji σ_I^2 i σ_{II}^2 odpowiednio. Przypuśćmy dalej, że dla obliczenia pierwszego z tych estymatorów generujemy n_I punktów, zaś dla drugiego n_{II} . Błędy obu estymatorów na tym samym poziomie istotności są ograniczone przez, odpowiednio, $z\sigma_I/\sqrt{n_I}$ i $z\sigma_{II}/\sqrt{n_{II}}$. Przyrównując te wyrażenia do siebie dochodzimy do wniosku, że oba estymatory osiągają podobną dokładność, jeśli $n_I/n_{II} = \sigma_{II}^2/\sigma_I^2$. Liczbę

$$\text{eff}(\hat{\theta}_n^I, \hat{\theta}_n^{II}) = \frac{\sigma_{II}^2}{\sigma_I^2}$$

nazywamy *względną efektywnością* (asymptotyczną). Czasami dobrze jest wybrać za „naturalny punkt odniesienia” estymator CMC i zdefiniować *efektywność* estymatora $\hat{\theta}_n$ o asymptotycznej wariancji σ^2 jako

$$\text{eff}(\hat{\theta}_n) = \text{eff}(\hat{\theta}_n, \hat{\theta}_n^{\text{CMC}}) = \frac{\sigma_{\text{CMC}}^2}{\sigma^2}.$$

Mówi się też, że jeśli wygenerujemy próbkę n punktów i obliczymy $\hat{\theta}_n$ to „efektywna liczność próbek” (ESS, czyli *effective sample size*) jest $n/\text{eff}(\hat{\theta}_n)$. Tyle bowiem należałoby wygenerować punktów stosując CMC, żeby osiągnąć podobną dokładność.

Asymptotyczna normalność prostego estymatora CMC wynika wprost z Centralnego Twierdzenia Granicznego (CTG). Istotnie, jeśli generujemy niezależnie X_1, \dots, X_n o jednakowym rozkładzie π , to

$$\begin{aligned} \sqrt{n}(\hat{\theta}_n^{\text{CMC}} - \theta) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n (f(X_i) - \theta) \\ &\rightarrow N(0, \sigma_{\text{CMC}}^2), \quad (n \rightarrow \infty), \end{aligned}$$

gdzie

$$\sigma_{\text{CMC}}^2 = \text{Var}_\pi f(X) = \int (f(x) - \theta)^2 p(x) dx.$$

Zupełnie podobnie, dla losowania istotnego w formie (8.1) otrzymujemy asymptotyczną normalność i przy tym

$$\begin{aligned}\sigma_{\text{IS1}}^2 &= \text{Var}_q w(X) f(X) = \mathbb{E}_q \frac{(f(X)p(X) - \theta q(X))^2}{q(X)^2} \\ &= \int \frac{(f(x)p(x) - \theta q(x))^2}{q(x)} dx.\end{aligned}$$

Z tego wzorku widać, że estymator może mieć wariancję zero jeśli $q(x) \propto f(x)p(x)$. Niestety, żeby obliczyć $q(x)$ potrzebna jest znajomość współczynnika proporcjonalności, który jest równy... θ . Nie wszystko jednak stracone. Pozostaje ważna reguła heurystyczna:

Gęstość $q(x)$ należy tak dobrać, aby jej „kształt” był zbliżony do funkcji $f(x)p(x)$.

Dla drugiej wersji losowania istotnego, (8.2), asymptotyczna normalność wynika z następujących rozważań:

$$\begin{aligned}\sqrt{n} (\hat{\theta}_n^{\text{IS2}} - \theta) &= \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n (W_i f(X_i) - \theta W_i)}{\frac{1}{n} \sum_{i=1}^n W_i} \\ &\rightarrow N(0, \sigma_{\text{IS2}}^2), \quad (n \rightarrow \infty),\end{aligned}$$

gdzie

$$\begin{aligned}\sigma_{\text{IS2}}^2 &= \text{Var}_q w(X) (f(X) - \theta) \\ &= \text{Var}_q w(X) f(X) - 2\theta \text{Cov}_q(w(X), w(X) f(X)) + \theta^2 \text{Var}_q w(X) \\ &= \sigma_{\text{IS1}}^2 + \theta [-2\text{Cov}_q(w(X), w(X) f(X)) + \theta \text{Var}_q w(X)].\end{aligned}$$

Wyrażenie w kwadratowym nawiasie może być ujemne, jeśli jest duża dodatnia korelacja zmiennych $w(X)$ i $w(X)f(X)$. W tej sytuacji estymator IS2 jest lepszy od IS1. Okazuje się więc rzecz na pozór paradoksalna: dzielenie przez estymator jedynki może poprawić estymator. Poza tym oba estymatory IS1 i IS2 mogą mieć mniejszą (asymptotyczną) wariancję, niż CMC. Jeśli efektywność jest większa niż 100%, to używa się czasem określenia „estymator *superefektywny*”.

Jeśli interesuje nas obliczenie wartości oczekiwanej dla wielu różnych funkcji f , to warto wprowadzić uniwersalny, niezależny od f wskaźnik efektywności. Niezłym takim wskaźnikiem jest $\text{Var}_q w(X)$. Istotnie, „odchylenie χ^2 ” pomiędzy gęstością instrumentalną q i docelową p , jest określone poniższym wzorem:

$$\begin{aligned}\chi^2 &= \chi^2(q, p) = \int \frac{(q(x) - p(x))^2}{q(x)} dx \\ &= \mathbb{E}_q \left(\frac{q(X) - p(X)}{q(X)} \right)^2 = \mathbb{E}_q (1 - w(X))^2 \\ &= \text{Var}_q w(X).\end{aligned}$$

Niestety, to „odchylenie” nie jest odległością, bo nie spełnia warunku symetrii. Niemniej widać, że małe wartości $\text{Var}_q w(X)$ świadczą o bliskości q i p . Zauważmy, jeszcze, że $\text{Var}_q w(X)$ jest wariancją asymptotyczną estymatora IS1 dla funkcji $f(x) = 1$.

Zgodnym estymatorem wielkości χ^2 jest

$$\hat{\chi}^2 = \frac{1}{n-1} \sum_{i=1}^n \frac{(\bar{W}_n - W_i)^2}{\bar{W}^2}.$$

Dzielenie przez $n-1$ (a nie n) wynika z tradycji. Dzielenie przez \bar{W}^2 powoduje, że estymator nie wymaga znajomości czynnika normującego gęstość p . Jest to w istocie estymator typu IS2.

Estymator typu IS1 jest równy po prostu $\sum(1 - W_i)^2/n$, ale można go stosować tylko gdy znamy czynnik normujący.

Wprowadźmy teraz bardziej formalnie pojęcie ważonej próbki. Niech (X, W) będzie parą zmiennych losowych o wartościach w $\mathcal{X} \times [0, \infty[$. Jeśli rozkład brzegowy X ma gęstość q zaś docelowa gęstość jest postaci $p(x) = p'(x)/Z$ to żądamy aby

$$\mathbb{E}(W|X = x) = \frac{p'(x)}{q(x)}$$

dla prawie wszystkich x . Mówimy wtedy, że próbka jest *dopasowana* do rozkładu p . Idea jest taka jak w losowaniu istotnym IS2: dopasowanie gwarantuje, że dla każdej funkcji f mamy $\mathbb{E}_q f(X)W = Z\mathbb{E}_p f(X)$. Nie jest konieczne, aby W było deterministyczną funkcją X . Pozostawiamy sobie możliwość generowania wag losowo.

Stała normująca Z jest na ogół nieznana, a więc waga W jest określona z dokładnością do proporcjonalności. W poprzednich podrozdziałach dla podkreślenia opatrywaliśmy takie nieunormowane wagi znacznikiem „prim” ale teraz to pomijamy. Gdy generujemy ciąg ważonych zmiennych losowych $(X_i, W_i), \dots, (X_n, W_n)$, to oczywiście wymaga się aby $\mathbb{E}(W_i|X_i = x) = Zp(x)/q(x)$, gdzie stała Z musi być jednakowa dla wszystkich $i = 1, \dots, n$. Jeśli założymy się niezależność par (X_i, W_i) , to zgodność i asymptotyczną normalność estymatora $\hat{\theta}_n^{\text{IS2}}$ danego równaniem (8.2) wykazuje się tak samo jak poprzednio. Asymptotyczna wariancja jest równa $Z^{-2}\text{Var}_q W(f(X) - \theta)$; czynnik Z^{-2} pojawia się dlatego, że operujemy nieunormowanymi wagami.

Jest jednak dużo ciekawych zastosowań, w których punkty X_i są *zależne*. Wtedy nawet zgodność estymatora $\hat{\theta}_n^{\text{IS2}}$ nie jest automatyczna. Asymptotyczna normalność może zachodzić z graniczną wariancją zupełnie inną niż w przypadku niezależnym lub nie zachodzić wcale.

8.3. Ważona eliminacja

Interesujące jest powiązanie idei eliminacji z ważeniem próbek. Czysta metoda eliminacji pracuje w sytuacji, gdy gęstość „instrumentalna” majoryzuje funkcję proporcjonalną do gęstości docelowej, czyli $Zp(x) = p'(x) \leq q(x)$. Jeśli ten warunek *nie jest* spełniony, to można naprawić odpowiednio *ważąc* wylosowaną próbkę. Dokładniej, algorytm ważonej eliminacji (*Rejection Control*, w skrócie RC) jest następujący.

Listing.

```

repeat
  Gen  $Y \sim q$ ;
   $W := p'(Y)/q(Y)$ ;
  if  $W \leq 1$  then
    begin
       $W := 1$ ;
      Gen  $U \sim U(0, 1)$ 
    end;
until  $(W > 1 \text{ or } U < W)$ 
 $X := Y$ 

```

Dowód poprawności algorytmu. Zauważmy, że prawdopodobieństwo akceptacji Y w tym algorytmie jest równe $a(Y)$, gdzie

$$a(x) = \frac{p'(x) \wedge q(x)}{q(x)} = \begin{cases} \frac{p'(x)}{q(x)} & \text{jeśli } p'(x) \leq q(x); \\ 1 & \text{jeśli } p'(x) > q(x). \end{cases}$$

Rozkład X na wyjściu algorytmu ma zatem gęstość $\propto p'(x) \wedge q(x)$. Waga $W = w(Y)$ jest więc „dopasowana” do rozkładu p w sensie zdefiniowanym w poprzednim podrozdziale, bo

$$w(x) = \frac{p'(x)}{p'(x) \wedge q(x)} = \begin{cases} 1 & \text{jeśli } p'(x) \leq q(x); \\ \frac{p'(x)}{q(x)} & \text{jeśli } p'(x) > q(x). \end{cases}$$

□

Zauważmy jeszcze, jak należy *modyfikować* wagi w RC, jeśli na wejściu mamy próbkę ważoną, dopasowaną do p . Jeżeli punkt Y , pochodzący z rozkładu q , ma na wejściu wagę $W_Y = p'(Y)/q(Y)$, to na wyjściu zaakceptowany punkt $X = Y$ ma rozkład $\propto p'(x) \wedge q(x)$ i powinien mieć wagę $W_X = p'(Y)/(p'(Y) \wedge q(Y))$. Widać stąd, że

$$W_Y = \frac{W_X}{a(Y)},$$

gdzie $a(\cdot)$ jest napisanym wyżej prawdopodobieństwem akceptacji w RC.

Przykład 8.1 (Nie-samo-przecinające się błędzenia). Po angielsku nazywają się *Self Avoiding Walks*, w skrócie SAW. Niech \mathbb{Z}^d będzie d -wymiarową kratą całkowitoliczbową. Mówimy, że ciąg $s = (0 = s_0, s_1, \dots, s_k)$ punktów kraty jest SAW-em jeśli

- każde dwa kolejne punkty s_{i-1} i s_i sąsiadują ze sobą, czyli różnią się o ± 1 na dokładnie jednej współrzędnej,
- żadne dwa punkty nie zajmują tego samego miejsca, czyli $s_i \neq s_j$ dla $i \neq j$.

Zbiór wszystkich SAW-ów o k ogniwach w \mathbb{Z}^d oznaczmy SAW_k^d , a dla d i k ustalonych w skrócie SAW. Przykład $s \in \text{SAW}_{15}^2$ widać na rysunku. Natychmiast nasuwa się bardzo proste pytanie:

- Jak policzyć SAW-y, czyli obliczyć $L = L_{k,d} = |\text{SAW}_k^d|$?

Zainteresujmy się teraz „losowo wybranym SAW-em”. Rozumiemy przez to zmienną losową S o rozkładzie jednostajnym w zbiorze SAW_k^d , czyli taką, że $\mathbb{P}(S = s) = 1/L_{k,d}$ dla każdego $s \in \text{SAW}_k^d$. W skrócie, $S \sim \text{U}(\text{SAW})$. Niech $\delta(s_k)$ oznacza odległość euklidesową końca SAW-a, czyli punktu s_k , od początku, czyli punktu 0. Na przykład dla łańcuszka widocznego na rysunku mamy $\delta(s_{15}) = \sqrt{5}$. Można zadać sobie pytanie, jaka jest średni kwadrat takiej odległości, czyli

- Jak obliczyć $\overline{\delta^2} = \overline{\delta_{d,k}^2} = \mathbb{E}\delta(S_k)^2$, przy założeniu, że $S \sim \text{U}(\text{SAW})$?

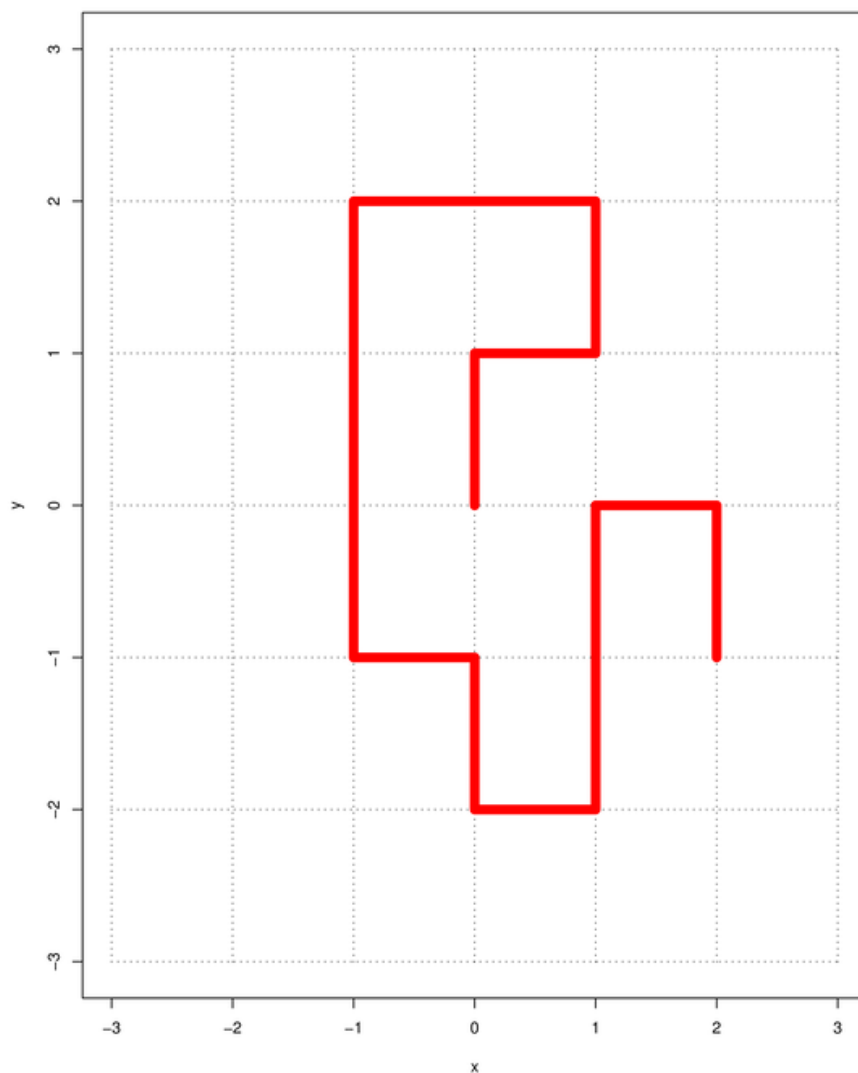
Można zastosować prostą metodę MC, czyli eliminację. Niech WALK_k^d oznacza zbiór wszystkich „błądzeń”, czyli ciągów $s = (0 = s_0, s_1, \dots, s_k)$ niekoniecznie spełniających warunek „ $s_i \neq s_j$ dla $i \neq j$ ”. Oczywiście $|\text{WALK}_k^d| = (2d)^k$ i metoda generowania „losowego błędzenia” (z rozkładu $\text{U}(\text{WALK})$) jest bardzo prosta: kolejno losujemy pojedyncze kroki, wybierając zawsze jedną z $2d$ możliwości. Żeby otrzymać „losowy SAW”, stosujemy eliminację. Ten sposób pozwala w zasadzie estymować $\overline{\delta^2}$ (przez uśrednianie długości zaakceptowanych błędzeń) oraz SAW/WALK (przez zanotowanie frakcji zaakceptowanych błędzeń). Niestety, metoda jest bardzo nieefektywna, bo dla dużych k prawdopodobieństwo akceptacji szybko zbliża się do zera.

Metoda „wzrostu” zaproponowana przez Rosenbluthów polega na losowaniu kolejnych kroków błędzenia spośród „dopuszczalnych punktów”, to znaczy punktów wcześniej nie odwiedzonych. W każdym kroku, z wyjątkiem pierwszego mamy co najwyżej $2d - 1$ możliwości. W błędzeniu widocznym na rysunku kolejne kroki wybieraliśmy spośród:

$$4, 3, 3, 3, \quad 2, 3, 2, 2, \quad 3, 2, 3, 3, \quad 2, 1, 3, 2$$

możliwych. Nasz SAW został zatem wylosowany z prawdopodobieństwem

$$\frac{1}{4} \cdot \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{2} \cdot \frac{1}{3} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{3} \cdot \frac{1}{2} \cdot \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{2} \cdot \frac{1}{1} \cdot \frac{1}{3} \cdot \frac{1}{2}$$



Rysunek 8.1. Przykład SAW-a.

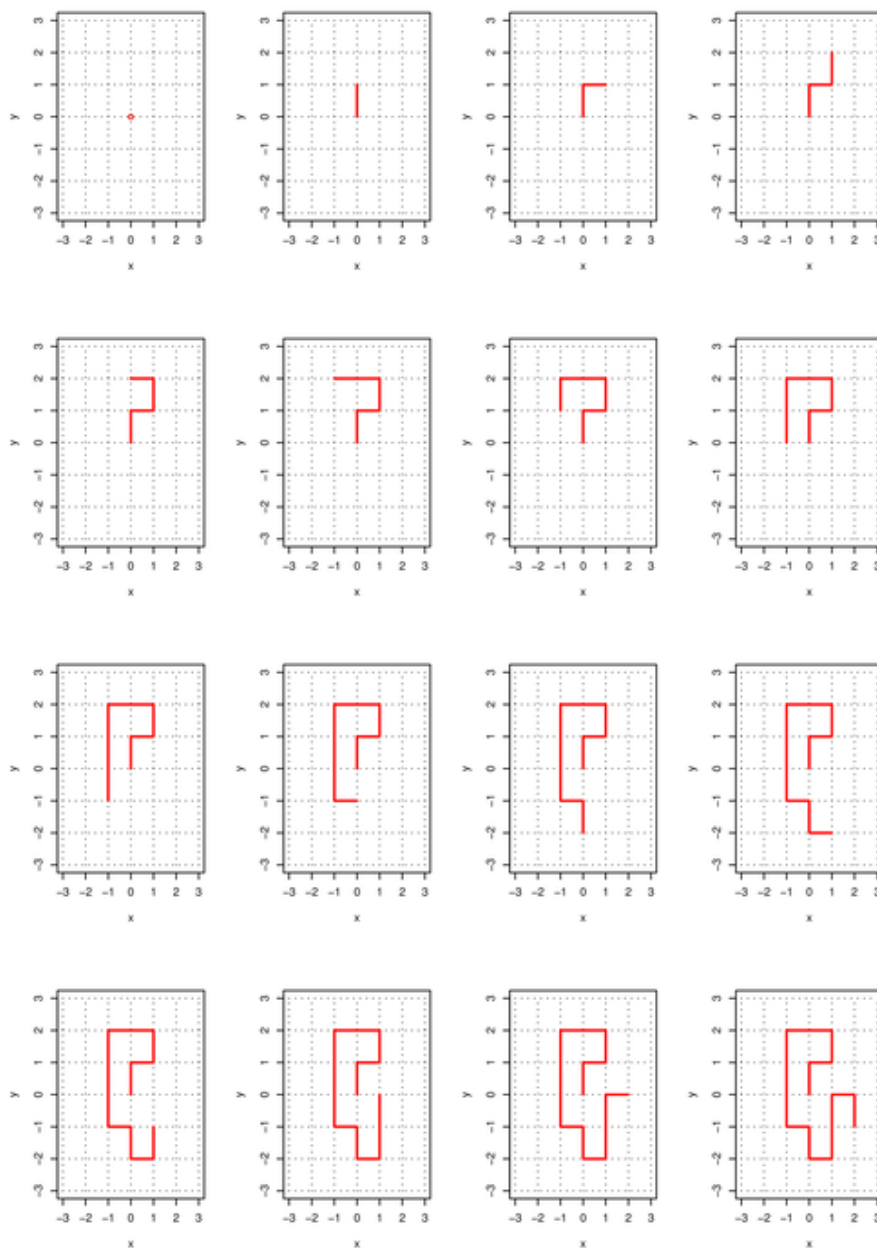
Powiedzmy ogólniej, że przy budowaniu SAW-a $s = (0 = s_0, s_1, \dots, s_k)$ mamy kolejno

$$n_1 = 2d, n_2, \dots, n_k$$

możliwości (nie jest przy tym wykluczone, że w pewnym kroku *nie mamy żadnej możliwości*, $n_i = 0$). Używając terminologii i oznaczeń związanych z losowaniem istotnym powiemy, że

$$q(s) = \frac{1}{n_1} \cdot \frac{1}{n_2} \cdots \frac{1}{n_k}.$$

jest gęstością instrumentalną dla $s \in \text{SAW}$, gęstość docelowa jest stała, równa $p(s) = 1/Z$, zatem funkcja podziału jest po prostu liczbą SAW-ów: $Z = |\text{SAW}|$. Wagi przypisujemy zgodnie ze wzorem $w(s) = n_1 \cdot n_2 \cdots n_k$ (jeśli wygenerowanie SAW-a się nie udało, $n_i = 0$ dla pewnego



Rysunek 8.2. Generowanie SAW-a metodą „wzrostu”.

i , to waga jest zero). Niech teraz $S(1), \dots, S(n)$ będą niezależnymi błędzeniami losowanymi metodą Rosenbluthów. Zgodnie z ogólnymi zasadami losowania istotnego,

$$\frac{\sum w(S(i))\delta(S(i))}{\sum w(S(i))} \text{ jest estymatorem } \overline{\delta^2},$$

$$\sum w(S(i)) \text{ jest estymatorem } |SAW|.$$

Należy przy tym *uwzględnić* w tych wzorach błędzenia o wadze zero, czyli „nieudane SAW-y”.

Przykład 8.2 (Prawdopodobieństwo ruiny i wykładnicza zamiana miary). Rozpatrzmy najprostszy model procesu opisującego straty i przychody w ubezpieczeniowej „teorii ryzyka”. Niech

Y_1, Y_2, \dots będą zmiennymi losowymi oznaczającymi *straty netto* (straty – przychody) towarzystwa ubezpieczeniowego w kolejnych okresach czasu. Założymy (co jest dużym uproszczeniem), że te zmienne są niezależne i mają jednakowy rozkład o gęstości $p(y)$. Tak zwana „nadwyżka ubezpieczyciela” na koniec n -go roku jest równa

$$u - S_n = u - \sum_{i=1}^n Y_i,$$

gdzie u jest rezerwą początkową. Interesuje nas prawdopodobieństwo zdarzenia, polegającego na tym, że $u - S_n < 0$ dla pewnego n . Mówimy wtedy (znowu w dużym uproszczeniu) o „ruinie ubezpieczyciela”. Wygodnie jest przyjąć następujące oznaczenia i konwencje. Zmienna losowa

$$R = \begin{cases} \min\{n : S_n > u\} & \text{jeśli takie } n \text{ istnieje;} \\ \infty & \text{jeśli } S_n \leq u \text{ dla wszystkich } n \end{cases}$$

oznacza czas oczekiwania na ruinę, przy czym jeśli ruina nigdy nie nastąpi to ten czas uznajemy za nieskończony. Przy takiej umowie, prawdopodobieństwo ruiny możemy zapisać jako $\psi = \mathbb{P}_p(R < \infty)$. Wskaźnik p przy symbolu wartości oczekiwanej przypomina, że chodzi tu o „oryginalny” proces, dla którego $Y_i \sim p$.

Obliczenie ψ analitycznie jest możliwe tylko w bardzo specjalnych przykładach. Metody numeryczne istnieją, ale też nie są łatwe. Pokażemy sposób obliczania ψ metodą Monte Carlo, który stanowi jeden z najpiękniejszych, klasycznych, przykładów losowania istotnego. Przyjmujemy bardzo rozsądne założenie, że $\mathbb{E}_p Y_i < 0$. Funkcję tworzącą momenty, która odpowiada gęstości p określamy wzorem

$$M_p(t) = \mathbb{E}_p e^{tY_i} = \int_{-\infty}^{\infty} e^{ty} p(y) dy.$$

Założymy, że ta funkcja przyjmuje wartości skończone przynajmniej w pewnym otoczeniu zera i istnieje takie $r > 0$, że

$$M_p(t) = 1.$$

Liczba r jest nazywana *współczynnikiem dopasowania* i odgrywa tu kluczową rolę.

Metoda *wykładniczej zamiany miary* jest specjalnym przypadkiem losowania istotnego. Żeby określić instrumentalny rozkład prawdopodobieństwa, połóżmy

$$q(y) = e^{ry} p(y).$$

Z definicji współczynnika dopasowania wynika, że q jest gęstością prawdopodobieństwa, to znaczy $\int q(y) dy = 1$. Generuje się ciąg Y_1, Y_2, \dots jednakowo rozłożonych zmiennych losowych o gęstości q . Dla utworzonego w ten sposób „instrumentalnego procesu” używać będziemy symboli \mathbb{P}_q i \mathbb{E}_q . Zauważmy, że

$$\begin{aligned} \mathbb{E}_q Y_i &= \int y q(y) dy = \int y e^{ry} p(y) dy \\ &= \frac{d}{dr} \int e^{ry} p(y) dy \\ &= M_p'(r) > 0, \end{aligned}$$

ponieważ funkcja tworząca momenty $M_p(\cdot)$ jest wypukłą, $M_p'(0) < 0$ i $M_p(0) = M_p(r) = 1$. Po zamianie miary, proces $u - S_n$ na mocy Prawa Wielkich Liczb zmierza prawie na pewno do minus nieskończoności, a zatem ruina następuje z prawdopodobieństwem jeden, $\mathbb{P}_q(R < \infty) = 1$.

$\infty) = 1$. Pokażemy, jak wyrazić prstwo ruiny dla oryginalnego procesu w terminach procesu instrumentalnego. Niech

$$\mathcal{R}_n = \{(y_1, \dots, y_n) : y_1 \leq u, \dots, y_1 + \dots + y_{n-1} \leq u, y_1 + \dots + y_{n-1} + y_n > u\},$$

Innymi słowy, zdarzenie $\{R = n\}$ zachodzi gdy $(Y_1, \dots, Y_n) \in \mathcal{R}_n$. Mamy zatem

$$\begin{aligned} \mathbb{P}_p(R = n) &= \int \dots \int_{\mathcal{R}_n} p(y_1) \dots p(y_n) dy_1 \dots dy_n \\ &= \int \dots \int_{\mathcal{R}_n} e^{-ry_1} q(y_1) \dots e^{-ry_n} q(y_n) dy_1 \dots dy_n \\ &= \int \dots \int_{\mathcal{R}_n} e^{-r(y_1 + \dots + y_n)} q(y_1) \dots q(y_n) dy_1 \dots dy_n \\ &= \mathbb{E}_q e^{rS_n} \mathbb{I}(R = n). \end{aligned}$$

Weźmy sumę powyższych równości dla $n = 1, 2, \dots$ i skorzystajmy z faktu, że $\sum_{n=1}^{\infty} \mathbb{P}_q(R = n) = 1$. Dochodzimy do wzoru

$$\mathbb{P}_p(R < \infty) = \mathbb{E}_q \exp\{-rS_R\} = e^{-ru} \mathbb{E}_q \exp\{-r(S_R - u)\}. \quad (8.3)$$

Ten fakt jest podstawą algorytmu Monte Carlo:

Listing.

$\hat{\psi} := 0$; [$\hat{\psi}$ będzie estymatorem prawdopodobieństwa ruiny]

for $m := 1$ **to** k **do**

begin

$n := 0$; $S := 0$;

repeat

 Gen $Y \sim q$; $S := S + Y$;

until $S > u$;

$Z := \exp\{-r(S - u)\}$; $\hat{\psi} := \hat{\psi} + Z$

end

$\hat{\psi} := e^{-ru} \hat{\psi} / k$

Algorytm jest prosty i efektywny. Trochę to zadziwiające, że w celu obliczenia prawdopodobieństwa ruiny generuje się proces dla którego ruina jest pewna. Po chwili zastanowienia można jednak zauważyć, że wykładnicza zamiana miary realizuje podstawową ideę losowania istotnego: rozkład instrumentalny „naśladuje” proces docelowy *ograniczony do zdarzenia* $\{R < \infty\}$.

Ciekawe, że wykładnicza zamiana miary nie tylko jest techniką Monte Carlo, ale jest też techniką *dowodzenia twierdzeń*! Aby się o tym przekonać, zauważmy, że „po drodze” udowodniłmy nierówność $\psi \leq e^{-ru}$ (wynika to z podstawowego wzoru (8.3) gdyż $\exp\{-r(S_R - u)\} \leq 1$). Jest to sławna nierówność Lundberga i wcale nie jest ona oczywista.

9. Algorytmy Monte Carlo II. Redukcja wariancji

W tym rozdziale omówię niektóre metody redukcji wariancji dla klasycznych algorytmów Monte Carlo, w których losujemy próbki niezależnie, z jednakowego rozkładu. Wiemy, że dla takich algorytmów wariancja estymatora zachowuje się jak const/n . Jedyne, co możemy zrobić – to konstruować takie algorytmy, dla których stała „const” jest możliwie mała. Najważniejszą z tych metod faktycznie już poznaliśmy: jest to *losowanie istotne*, omówione w Rozdziale 8. Odpowiedni wybór „rozkładu instrumentalnego” może zmniejszyć wariancję...setki tysięcy razy! Istnieje jeszcze kilka innych, bardzo skutecznych technik. Do podstawowych należą: losowanie warstwowe, metoda zmiennych kontrolnych, metoda zmiennych antytetycznych. Możliwe są niezliczone modyfikacje i kombinacje tych metod. Materiał zawarty w tym rozdziale jest w dużym stopniu zaczerpnięty z monografii Ripleya [18].

9.1. Losowanie warstwowe

Tak jak w poprzednim rozdziale, zedanie polega na obliczeniu wielkości

$$\theta = \mathbb{E}_\pi f(X) = \int_{\mathcal{X}} f(x)\pi(x)dx,$$

gdzie rozkład prawdopodobieństwa i jego gęstość dla uproszczenia oznaczamy tym samym symbolem π . **Losowanie warstwowe** polega na tym, że rozbijamy przestrzeń X na sumę k rozłącznych podzbiorów (warstw),

$$\mathcal{X} = \bigcup_{h=1}^k A_h, \quad A_h \cap A_g = \emptyset \quad (h \neq g),$$

i losujemy k niezależnych próbek, po jednej z każdej warstwy. Niech π_h będzie gęstością rozkładu *warunkowego* zmiennej X przy $X \in A_h$, czyli

$$\pi_h(x) = \frac{\pi(x)}{p_h} \mathbb{I}(x \in A_h), \quad \text{gdzie} \quad p_h = \pi(A_h) = \int_{A_h} \pi(x)dx.$$

Widać natychmiast, że $\int_B \pi_h(x)dx = \pi(X \in B | X \in A_h)$. Rozbijamy teraz całkę, którą chcemy obliczyć:

$$\theta = \sum_h p_h \int f(x)\pi_h(x)dx.$$

Możemy użyć następującego estymatora warstwowego:

$$\hat{\theta}_n^{\text{stra}} = \sum_h \frac{p_h}{n_h} \sum_{i=1}^{n_h} f(X_{hi}), \tag{9.1}$$

gdzie

$$X_{h1}, \dots, X_{hn_h} \sim_{\text{i.i.d}} \pi_h,$$

jest próbką rozmiaru n_h wylosowaną z h -tej warstwy ($h = 1, \dots, k$). Jest to estymator nieobciążony,

$$\text{Var} \hat{\theta}_n^{\text{stra}} = \sum_h \frac{p_h^2}{n_h} \sigma_h^2, \quad (9.2)$$

gdzie $\sigma_h^2 = \text{Var}_\pi(f(X)|x \in A_h)$.

Porównajmy estymator warstwowy (9.1) ze „zwykłym” estymatorem

$$\hat{\theta}_n^{\text{CMC}} = \frac{1}{n} \sum_{i=1}^n f(X_i).$$

opartym na jednej próbce

$$X_1, \dots, X_n \sim_{\text{i.i.d}} \pi.$$

Oczywiście

$$\text{Var} \hat{\theta}_n^{\text{CMC}} = \frac{1}{n} \sigma^2,$$

gdzie $\sigma^2 = \text{Var}_\pi f(X)$. Żeby porównanie było „uczciwe” rozważmy próbkę liczności $n = \sum_h n_h$. Jeśli decydujemy się na sumaryczną licznosc próbki n , to możemy w różny sposób „rozdzielić” to n pomiędzy warstwami. To się nazywa *alokacja* próbki.

Stwierdzenie 9.1 (Alokacja proporcjonalna). *Jeżeli $n_h = p_h n$ dla $n = 1, \dots, k$, to*

$$\text{Var} \hat{\theta}_n^{\text{stra}} = \frac{1}{n} \sum_h p_h \sigma_h^2 \leq \frac{1}{n} \sigma^2 = \text{Var} \hat{\theta}_n^{\text{CMC}}.$$

Dowód. Wyrażenie na wariancję wynika znatychmiast z podstawienia $n_h = p_h n$ w ogólnym wzorze (9.2). Nierówność wynika z następującej tożsamości:

$$\sigma^2 = \sum_h p_h \sigma_h^2 + \sum_h p_h (\theta_h - \theta)^2, \quad (9.3)$$

gdzie $\theta_h = \int f(x) \pi_h(x) dx = \mathbb{E}_\pi(f(X)|X \in A_h)$. Jest to klasyczny wzór na dekompozycję wariancji na „wariancję wewnątrz warstw” (pierwszy składnik w (9.3)) i „wariancję pomiędzy warstwami” (drugi składnik w (9.3)). Zdefiniujemy zmienną losową H jako „numer warstwy do której wpada $X \sim \pi$ ”, czyli $\{H = h\} = \{X \in A_h\}$. Możemy teraz wzór (9.3) przepisać w dobrze znanej postaci

$$\text{Var} f(X) = \mathbb{E} \text{Var}(f(X)|H) + \text{Var} \mathbb{E}(f(X)|H).$$

□

Wzór (9.3) podpowiada, jak dzielić przestrzeń na warstwy. Największy zysk w porównaniu z losowaniem nie-warstwowym jest wtedy, gdy „wariancja międzywarstwowa” jest dużo większa od „wewnątrzwarstwowej”. Warstwy należy więc wybierać tak, żeby funkcja $\pi(c)f(x)$ była możliwie bliska stałej na każdym zbiorze A_h i różniła się jak najbardziej pomiędzy poszczególnymi zbiorami.

Stwierdzenie 9.1 pokazuje, że zawsze zyskujemy na losowaniu warstwowym, jeśli zastosujemy najprostszą, proporcjonalną alokację. Okazuje się, że nie jest to alokacja najlepsza. Jerzy Sława-Neyman odkrył prostą regułę wyznaczania alokacji optymalnej. Wychodzimy od wzoru (9.2) i staramy się zoptymalizować prawą stronę przy warunku $n = \sum_h n_h$. Poszukujemy zatem rozwiązania zadania

$$\sum_h \frac{p_h^2}{n_h} \sigma_h^2 = \min ! \quad \left(\sum_h n_h - n = 0 \right) \quad (9.4)$$

(względem zmiennych n_h). Zastosujemy metodę mnożników Lagrange’a. Szukamy minimum

$$\mathcal{L} = \sum_h \frac{p_h^2}{n_h} \sigma_h^2 + \lambda (\sum_h n_h - n) = \min ! \quad (9.5)$$

Obliczamy pochodną i przyrównujemy do zera:

$$\frac{\partial}{\partial n_h} \mathcal{L} = -\frac{p_h^2}{n_h^2} \sigma_h^2 + \lambda = 0. \quad (9.6)$$

Stąd natychmiast otrzymujemy rozwiązanie: $n_h \propto \sigma_h p_h$.

Stwierdzenie 9.2 (Alokacja optymalna, J. Neyman). *Estymator warstwowy ma najmniejszą wariancję jeśli alokacja n losowanych punktów jest dana wzorem*

$$n_h = \frac{\sigma_h p_h}{\sum_g \sigma_g p_g}, \quad (h = 1, \dots, k).$$

Zignorowaliśmy tutaj niewygodne wymaganie, że licznosci próbek n_h muszą być całkowite. Gdyby to wziąć pod uwagę, rozwiązanie stałoby się skomplikowane, a zysk praktyczny z tego byłby znikomy. W praktyce rozwiązanie neymanowskie zaokrągla się do liczb całkowitych i koniec. Ważniejszy jest inny problem. Żeby wyznaczyć alokację optymalną, trzeba znać nie tylko prawdopodobieństwa p_h ale i wariancje warstwowe σ_h^2 . W praktyce często opłaca się wylosować wstępne próbki, na podstawie których *estymuje* się wariancje σ_h^2 . Dopiero potem alokuje się dużą, roboczą próbkę rozmiaru n , która jest wykorzystana do obliczania docelowej całki.

9.2. Zmienne kontrolne

Idea zmiennych kontrolnych polega na rozbiciu docelowej całki (wartości oczekiwanej) na dwa składniki, z których jeden umiemy obliczyć analitycznie. Metodę Monte Carlo stosujemy do drugiego składnika. Przedstawmy wielkość obliczaną w postaci

$$\theta = \mathbb{E}_\pi f(X) = \int_{\mathcal{X}} f(x) \pi(x) dx = \int_{\mathcal{X}} [f(x) - k(x)] \pi(x) dx + \int_{\mathcal{X}} k(x) \pi(x) dx.$$

Dążymy do tego, żeby całka funkcji k była obliczona analitycznie (lub numerycznie) a różnica $f - k$ była możliwie bliska stałej, bo wtedy wariancja metody Monte Carlo jest mała. Funkcję k lub zmienną losową $k(X)$ nazywamy zmienną kontrolną.

Dla uproszczenia połóżmy $Y = f(X)$, gdzie $X \sim \pi$. Przypuśćmy, że zmiennej kontrolnej będziemy szukać pośród kombinacji liniowych funkcji k_1, \dots, k_d o znanych całkach. Niech $Z_j = k_j(X)$ i $Z = (Z_1, \dots, Z_d)^\top$. Przy tym stale pamiętajmy, że $X \sim \pi$ i będziemy pomijać indeks π przy wartościach oczekiwanych i wariancjach. Zatem

$$k(x) = \sum_{j=1}^d \beta_j k_j(x).$$

Innymi słowy poszukujemy wektora współczynników $\beta = (\beta_1, \dots, \beta_d)^\top$ który minimalizuje wariancję $\text{Var}(Y - \beta^\top Z)$. Wykorzystamy następujący standardowy wynik z teorii regresji liniowej.

Stwierdzenie 9.3. *Niech Y i Z będą zmiennymi losowymi o skończonych drugich momentach, wymiaru odpowiednio 1 i d . Zakładamy dodatkowo odwracalność macierzy wariancji-kowariancji $\text{VAR}(Z)$. Kowariancję pomiędzy Y i Z traktujemy jako wektor wierszowy i oznaczamy przez*

$\text{COV}(Y, Z)$ Spośród zmiennych losowych postaci $Y - \beta^\top Z$, najmniejszą wariancję otrzymujemy dla

$$\beta_*^\top = \text{COV}(Y, Z) \text{VAR}(Z)^{-1}.$$

Ta najmniejsza wartość wariancji jest równa

$$\text{Var}(Y - \beta_*^\top Z) = \text{Var}Y - \text{COV}(Y, Z) \text{VAR}(Z)^{-1} \text{COV}(Z, Y).$$

Przepiszmy ten wynik w bardziej sugestywnej formie. Niech $\text{Var}Y = \sigma^2$. Można pokazać, że β_* maksymalizuje korelację pomiędzy Y i $\beta^\top Z$. Napiszmy

$$\varrho_{Y,Z} = \max_{\beta} \text{corr}(Y, \beta^\top Z) = \text{corr}(Y, \beta_*^\top Z) = \sqrt{\frac{\text{COV}(Y, Z) \text{VAR}(Z)^{-1} \text{COV}(Z, Y)}{\text{Var}Y}}.$$

Niech teraz $\hat{\theta}_n^{\text{contr}}$ będzie estymatorem zmiennych kontrolnych. To znaczy, że losujemy próbkę $X_1, \dots, X_n \sim \pi$,

$$\hat{\theta}_n^{\text{contr}} = \frac{1}{n} \sum (Y_i - \beta_*^\top Z_i) + \beta_*^\top \mu,$$

gdzie $Z_i^\top = (Z_{i1}, \dots, Z_{in}) = (f_1(X_i), \dots, f_d(X_i))$, zaś $\mu_j = \mathbb{E}f_j(X)$ są obliczone analitycznie lub numerycznie $\mu = (\mu_1, \dots, \mu_d)^\top$. Wariancja estymatora jest wyrażona wzorem

$$\text{Var}\hat{\theta}_n^{\text{contr}} = \frac{1}{n} \sigma^2 (1 - \varrho_{Y,Z}^2),$$

co należy porównać z wariancją σ^2/n „zwykłego” estymatora.

Zróbmy podobną uwagę, jak w poprzednim podrozdziale. Optymalny wybór współczynników regresji wymaga znajomości wariancji i kowariancji, których obliczenie może być ... (i jest zazwyczaj!) trudniejsze niż wyjściowe zadanie obliczenia wartości oczekiwanej. Niemniej, można najpierw wylosować *wstępną* próbkę, wyestymować potrzebne wariancje i kowariancje (nawet niezbyt dokładnie) po to, żeby dla dużej, *roboczej* próbki skonstruować dobre zmienne kontrolne.

Wspomnijmy na koniec, że dobieranie zmiennej kontrolnej metodą regresji liniowej nie jest jedynym sposobem. W konkretnych przykładach można spotkać najróżniejsze, bardzo pomysłowe konstrukcje, realizujące podstawową ideę zmiennych kontrolnych.

9.3. Zmienne antytetyczne

Przypuśćmy, że estymujemy wielkość $\theta = \mathbb{E}_\pi f(X)$. Jeśli mamy dwie zmienne losowe X i X' o jednakowym rozkładzie π ale nie zakładamy ich niezależności, to

$$\text{Var} \frac{f(X) + f(X')}{2} = \frac{1}{2} \text{Var}(X) [1 + \text{corr}(f(X), f(X'))] = \frac{1}{2} \sigma^2 (1 + \varrho).$$

Jeśli $\varrho = \text{corr}(f(X), f(X')) < 0$, to wariancja w powyższym wzorze jest *mniejsza* niż $\sigma^2/2$, czyli mniejsza niż w przypadku niezależności X i X' . To sugeruje, żeby zamiast losować *niezależnie* n zmiennych X_1, \dots, X_n , wygenerować *pary zmiennych ujemnie skorelowanych*. Załóżmy, że $n = 2k$ i mamy k niezależnych par $(X_1, X'_1), \dots, (X_k, X'_k)$, a więc łącznie n zmiennych. Porównajmy wariancję „zwykłego” estymatora $\hat{\theta}_n^{\text{CMC}}$ i estymatora $\hat{\theta}_n^{\text{ant}}$ wykorzystującego ujemne skorelowanie par:

$$\begin{aligned} \hat{\theta}_n^{\text{CMC}} &= \frac{1}{n} \sum_{i=1}^n f(X_i), & \text{Var}\hat{\theta}_n^{\text{CMC}} &= \frac{\sigma^2}{n} \\ \hat{\theta}_n^{\text{ant}} &= \frac{1}{n} \sum_{i=1}^{n/2} [f(X_i) - f(X'_i)], & \text{Var}\hat{\theta}_n^{\text{ant}} &= \frac{\sigma^2}{n} (1 + \varrho). \end{aligned}$$

Wariancja estymatora $\hat{\theta}_n^{\text{ant}}$ jest tym mniejsza, im ϱ bliższe -1 (im bardziej ujemnie skorelowane są pary). Standardowym sposobem generowania ujemnie skorelowanych par jest odwracanie dystrybucji z użyciem „odwróconych losów losowych”.

Stwierdzenie 9.4. *Jeśli $h :]0, 1[\rightarrow \mathbb{R}$ jest funkcją monotoniczną różną od stałej, $\int_0^1 h(u)^2 du < \infty$ i $U \sim U(0, 1)$, to*

$$\text{Cov}(h(U), h(1 - U)) < 0.$$

Dowód. Bez straty ogólności założmy, że h jest niemalejąca. Niech $\mu = \mathbb{E}h(U) = \int_0^1 h(u) du$ i $t = \sup\{u : h(1 - u) > \mu\}$. Łatwo zauważyć, że $0 < t < 1$. Zauważmy, że

$$\begin{aligned} \text{Cov}(h(U), h(1 - U)) &= \mathbb{E}h(U)[h(1 - U) - \mu] \\ &= \int_0^1 h(u)[h(1 - u) - \mu] du \\ &= \int_0^t h(u)[h(1 - u) - \mu] du + \int_t^1 h(u)[h(1 - u) - \mu] du \\ &< \int_0^t h(t)[h(1 - u) - \mu] du + \int_t^1 h(t)[h(1 - u) - \mu] du \\ &= h(t) \int_0^1 [h(1 - u) - \mu] du = 0, \end{aligned}$$

ponieważ dla $0 < u < t$ mamy $h(1 - u) - \mu > 0$ i dla $t < u < 1$ mamy $h(1 - u) - \mu \leq 0$. \square

Przypomnijmy, że dla dowolnej dystrybucji G określamy uogólnioną funkcję odwrotną G^- następującym wzorem:

$$G^-(u) = \inf\{x : G(x) \geq u\}.$$

Natychmiast wnioskujemy ze Stwierdzenia 9.4, że $\text{Cov}(G^-(U), G^-(1 - U)) < 0$. W ten sposób możemy produkować ujemnie skorelowane pary zmiennych o zadanej dystrybucji. Okazuje się, że są to najbardziej ujemnie skorelowane pary. Jeśli $X \sim G$ i $X' \sim G$, to

$$\text{Cov}(X, X') \geq \text{Cov}(G^-(U), G^-(1 - U)).$$

Powyższy fakt ma oczywiste znaczenie z punktu widzenia metod Monte Carlo. Udowodnimy ogólniejsze twierdzenie.

Twierdzenie 9.1. *Jeżeli $X \sim F$, $Y \sim G$ oraz $\mathbb{E}X^2 < \infty$, $\mathbb{E}Y^2 < \infty$ to*

$$\text{Cov}(F^-(U), G^-(1 - U)) \leq \text{Cov}(X, Y) \leq \text{Cov}(F^-(U), G^-(U)).$$

Twierdzenie 9.2 wynika z trzech poniższych faktów. Każdy z nich jest sam w sobie interesujący. Zaczniemy od sławnego wyniku Frecheta.

Twierdzenie 9.2 (Ograniczenia Frecheta). *Jeżeli $\mathbb{P}(X \leq x) = F(x)$, $\mathbb{P}(Y \leq x) = G(x)$ oraz $\mathbb{P}(X \leq x, Y \leq y) = H(x, y)$ oznaczają, odpowiednio, dystrybucje brzegowe oraz łączną dystrybucję dwóch zmiennych losowych to*

$$\max(0, F(x) + G(y) - 1) \leq H(x, y) \leq \min(F(x), G(y)).$$

Istnieją rozkłady łączne o brzegowych F i G , dla których jest osiągane ograniczenie dolne i ograniczenie górne.

Dowód. Ograniczenie górne wynika z oczywistych nierówności

$$\begin{aligned}\mathbb{P}(X \leq x, Y \leq y) &\leq \mathbb{P}(X \leq x) = F(x), \\ \mathbb{P}(X \leq x, Y \leq y) &\leq \mathbb{P}(Y \leq y) = G(y).\end{aligned}$$

Ograniczenie dolne jest równie proste:

$$\begin{aligned}\mathbb{P}(X \leq x, Y \leq y) &= \mathbb{P}(X \leq x) - \mathbb{P}(X \leq x, Y > y) \\ &\geq \mathbb{P}(X \leq x) - \mathbb{P}(Y > y) = F(x) - [1 - G(y)].\end{aligned}$$

□

Pozostała do pokazania osiągalność. Następujący lemat jest jednym z piękniejszych przykładów „symulacyjnego punktu widzenia” w rachunku prawdopodobieństwa.

Lemat 9.1. *Jeśli $U \sim U(0, 1)$ to*

$$\begin{aligned}\mathbb{P}(F^-(U) \leq x, G^-(U) \leq y) &= \min(F(x), G(y)); \\ \mathbb{P}(F^-(U) \leq x, G^-(1 - U) \leq y) &= \max(0, F(x) + G(y) - 1).\end{aligned}$$

Dowód. Pierwsza równość jest oczywista:

$$\begin{aligned}\mathbb{P}(F^-(U) \leq x, G^-(U) \leq y) &= \mathbb{P}(U \leq F(x), U \leq G(y)) \\ &= \min(F(x), G(y)).\end{aligned}$$

Druga równość też jest oczywista:

$$\begin{aligned}\mathbb{P}(F^-(U) \leq x, G^-(1 - U) \leq y) &= \mathbb{P}(U \leq F(x), 1 - U \leq G(y)) \\ &= \mathbb{P}(1 - G(y) \leq U \leq F(x)) \\ &= \max(0, F(x) - [1 - G(y)]).\end{aligned}$$

□

Głębokie twierdzenie Frecheta składa się więc z kilku dość oczywistych spostrzeżeń. Zeby udowodnić Twierdzenie 9.1 potrzeba jeszcze jednego ciekawego lematu.

Lemat 9.2. *Niech $F(x)$, $G(y)$ i $H(x, y)$ oznaczają, odpowiednio, dystrybuanty brzegowe oraz łączną dystrybuantę zmiennych losowych X i Y . Jeśli $\mathbb{E}X^2 < \infty$, $\mathbb{E}Y^2 < \infty$ to*

$$\text{Cov}(X, Y) = \iint [H(x, y) - F(x)F(y)] dx dy.$$

Dowód. Niech (X_1, Y_1) i (X_2, Y_2) będą niezależnymi parami o jednakowym rozkładzie takim jak para (X, Y) . Wtedy

$$\begin{aligned}2\text{Cov}(X, Y) &= \mathbb{E}(X_1 - X_2)(Y_1 - Y_2) \\ &= \mathbb{E} \iint [\mathbb{I}(X_1 \leq x) - \mathbb{I}(X_2 \leq x)][\mathbb{I}(Y_1 \leq y) - \mathbb{I}(Y_2 \leq y)] dx dy \\ &= \iint \mathbb{E}[\mathbb{I}(X_1 \leq x) - \mathbb{I}(X_2 \leq x)][\mathbb{I}(Y_1 \leq y) - \mathbb{I}(Y_2 \leq y)] dx dy \\ &= \iint [\mathbb{P}(X_1 \leq x, Y_1 \leq y) + \mathbb{P}(X_2 \leq x, Y_2 \leq y) \\ &\quad - \mathbb{P}(X_1 \leq x, Y_2 \leq y) - \mathbb{P}(X_2 \leq x, Y_1 \leq y)] dx dy \\ &= \iint [2\mathbb{P}(X \leq x, Y \leq y) - 2\mathbb{P}(X \leq x)\mathbb{P}(Y \leq y)] dx dy.\end{aligned}$$

□

10. Markowskie Monte Carlo I. Wprowadzenie

10.1. Co to jest MCMC ?

W pewnych sytuacjach okazuje się, że wygenerowanie zmiennej losowej X z interesującego nas rozkładu prawdopodobieństwa π jest praktycznie niemożliwe. Wyobraźmy sobie, że π jest bardzo skomplikowanym rozkładem na „ogromnej”, wielowymiarowej przestrzeni \mathcal{X} . Zazwyczaj ten rozkład jest dany poprzez podanie funkcji proporcjonalnej do gęstości, $p' \propto p$, ale bez znajomości stałej normującej $1/\int p'$. Przy tym „ogromna” przestrzeń \mathcal{X} może być zbiorem skończonym (ale bardzo liczny) zaś rozkład π może być... jednostajny, o gęstości $p \propto 1$. Jeśli jednak nośnik rozkładu $U(\mathcal{X})$ jest bardzo „skomplikowanym” zbiorem, to metody typu eliminacji/akceptacji mogą zawieść. Właśnie z takim przypadkiem spotkaliśmy się już we wstępnym Rozdziale 1, w Przykładzie 1.5. Naszkicowany tam algorytm jest przedstawicielem metod, którymi obecnie zajmujemy się dokładniej, czyli algorytmów MCMC. Ten skrót oznacza *Markov Chain Monte Carlo*, czyli po polsku algorytmy MC wykorzystujące łańcuchy Markowa. Podstawowa idea jest taka: jeśli nie umiemy generować zmiennej losowej X o rozkładzie π to zadowolimy się generowaniem ciągu zmiennych losowych $X_0, X_1, \dots, X_n, \dots$, który w pewnym sensie zbliża się, zmierza do rozkładu π . Podobnie jak w cytowanym przykładzie „plecakowym” ciąg X_n ma charakter „losowego błędzenia” po przestrzeni \mathcal{X} .

To wszystko co dotychczas powiedzieliśmy jest bardzo niekonkretne i wymaga uściślenia. W moich wykładach ściśle przedstawienie tych zagadnień będzie możliwe tylko w ograniczonym zakresie. W obecnym rozdziale skupię się na głównych ideach MCMC. Następnie, w Rozdziałach 11 i 12 przedstawię podstawowe algorytmy i kilka motywujących przykładów, pokażę wyniki symulacji, ale niemal nic nie udowodnię. Spróbuję to częściowo naprawić w Rozdziałach 14 i 15, które w całości będą poświęcone łańcuchom Markowa i algorytmom MCMC na skończonej przestrzeni \mathcal{X} . W tej uproszczonej sytuacji podam dowody (a przynajmniej szkice dowodów) podstawowych twierdzeń. Zastosowania MCMC w przypadku „ciągłej” przestrzeni \mathcal{X} (powiedzmy, $\mathcal{X} \subseteq \mathbb{R}^d$) są przynajmniej równie ważne i w gruncie rzeczy bardzo podobne. Zobaczymy to na paru przykładach. Algorytmy MCMC pracujące na przestrzeni ciągłej są w zasadzie takie same jak te w przypadku przestrzeni skończonej, ale ich analiza robi się trudniejsza, wymaga więcej abstrakcyjnej matematyki – i w rezultacie wykracza poza zakres mojego skryptu. Ograniczę się w tej materii do kilku skromnych uwag. Bardziej systematyczne, a przy tym dość przystępne przedstawienie ogólnej teorii można znaleźć w [20] lub [17].

10.2. Łańcuchy Markowa

Klasyczne metody MCMC, jak sama nazwa wskazuje, opierają się na generowaniu łańcucha Markowa. Co prawda, rozwijają się obecnie bardziej wyrafinowane metody MCMC (zwane adaptacyjnymi), które wykorzystują procesy niejednorodne a nawet nie-markowskie. Na razie ograniczymy się do rozpatrzenia sytuacji, gdy generowany ciąg zmiennych losowych X_n jest *jednorodnym łańcuchem Markowa* na przestrzeni \mathcal{X} . Jeśli \mathcal{X} jest zbiorem skończonym, możemy

posługiwać się Definicją 7.1, w ogólnym przypadku – Definicją 7.2. Przypomnijmy oznaczenie *prawdopodobieństw przejścia* łańcucha: dla $x \in \mathcal{X}$ oraz $B \subseteq \mathcal{X}$,

$$\mathbb{P}(X_{n+1} \in B | X_n = x) = P(x, B).$$

W przypadku przestrzeni skończonej wygodniej posługiwać się *macierzą przejścia* o elementach

$$\mathbb{P}(X_{n+1} = y | X_n = x) = P(x, y).$$

Metoda generowania łańcuchów Markowa jest dość oczywista i sprowadza się do wzorów (7.2) w przypadku skończonym i (7.5) w przypadku ogólnym.

10.2.1. Rozkład stacjonarny

Niech π oznacza docelowy rozkład prawdopodobieństwa. Chcemy tak generować łańcuch X_n , czyli tak wybrać prawdopodobieństwa przejścia P , aby uzyskać zbieżność do rozkładu π . Spróbujemy uściślić co to znaczy, w jakim sensie rozumiemy zbieżność.

Definicja 10.1. Mówimy, że π jest **rozkładem stacjonarnym** (lub rozkładem równowagi) łańcucha Markowa o prawdopodobieństwach przejścia P , jeśli dla każdego (mierzalnego) zbioru $B \subseteq \mathcal{X}$ mamy

$$\pi(B) = \int_{\mathcal{X}} \pi(dx) P(x, B).$$

W przypadku skończonej przestrzeni \mathcal{X} równoważne sformułowanie jest takie: dla każdego stanu y ,

$$\pi(y) = \sum_{x \in \mathcal{X}} \pi(x) P(x, y).$$

Jeśli rozkład początkowy jest rozkładem stacjonarnym, $\mathbb{P}(X_0 \in \cdot) = \pi(\cdot)$, to mamy $\mathbb{P}(X_n \in \cdot) = \pi(\cdot)$ dla każdego n . Co więcej, w takiej sytuacji łączny rozkład zmiennych X_n, X_{n+1}, \dots jest taki sam, jak rozkład zmiennych X_0, X_1, \dots . Mówimy, że łańcuch jest w położeniu równowagi lub, że jest *procesem stacjonarnym*. To uzasadnia nazwę rozkładu stacjonarnego. Oczywiście, łańcuchy generowane przez algorytmy MCMC nie są w stanie równowagi, bo z założenia nie umiemy wygenerować $X_0 \sim \pi$. Generujemy X_0 z pewnego innego rozkładu ξ , nazywanego rozkładem początkowym. Gdy zajdzie potrzeba, żeby uwidocznili zależność od rozkładu początkowego, będziemy używali oznaczeń $\mathbb{P}_\xi(\dots)$ i $\mathbb{E}_\xi(\dots)$. Przeważnie start jest po prostu deterministyczny, czyli ξ jest rozkładem skupionym w pewnym punkcie $x \in \mathcal{X}$. Piszemy wtedy $\mathbb{P}_x(\dots)$ i $\mathbb{E}_x(\dots)$.

Jeśli π jest rozkładem stacjonarnym, to przy pewnych dodatkowych założeniach uzyskuje się tak zwane *Słabe Twierdzenie Ergodyczne* (STE). Jego tezą jest zbieżność rozkładów prawdopodobieństwa zmiennych losowych X_n do π w następującym sensie: dla dowolnego (mierzalnego) zbioru $B \subseteq \mathcal{X}$ i dowolnego rozkładu początkowego ξ mamy

$$\mathbb{P}_\xi(X_n \in B) \rightarrow \pi(B) \quad (n \rightarrow \infty). \quad (10.1)$$

Dla skończonej przestrzeni \mathcal{X} równoważne jest stwierdzenie, że dla każdego $y \in \mathcal{X}$,

$$\mathbb{P}_\xi(X_n = y) \rightarrow \pi(y) \quad (n \rightarrow \infty).$$

Pewną wersję STE dla skończonej przestrzeni \mathcal{X} udowodnimy w następnym rozdziale (Twierdzenie 14.5). Na razie poprzestańmy na następującym prostym spostrzeżeniu.

Uwaga 10.1. Jeżeli zachodzi teza STE dla pewnego rozkładu granicznego π_∞ , czyli $P^n(x, B) \rightarrow \pi_\infty(B)$ dla dowolnych $x \in \mathcal{X}$, $B \subseteq \mathcal{X}$, to π_∞ jest rozkładem stacjonarnym. W istocie, wystarczy przejść do granicy w równości

$$\begin{array}{ccc} P^{n+1}(x, B) & = & \int P^n(x, dz) P(z, B) \\ \downarrow & & \downarrow \\ \pi_\infty(B) & & \int \pi_\infty(dz) P(z, B). \end{array}$$

Co więcej, π_∞ jest *jedynym* rozkładem stacjonarnym.

Uwaga 10.2. W teorii łańcuchów Markowa rozważa się różne pojęcia zbieżności rozkładów. Zauważmy, że w powyżej przytoczonej tezie STE oraz w Uwadze 10.1 mamy do czynienia z silniejszym rodzajem zbieżności, niż poznana na rachunku prawdopodobieństwa zbieżność słaba (według rozkładu), oznaczana \rightarrow_d .

10.2.2. Twierdzenia graniczne dla łańcuchów Markowa

Naszukujemy podstawowe twierdzenia graniczne dla łańcuchów Markowa. Dokładniejsze sformułowania i dowody pojawiają się w następnym rozdziale i będą ograniczone do przypadku skończonej przestrzeni \mathcal{X} . Bardziej dociekliwych Czytelników muszę odesłać do przeglądowych prac [20, 13].

Sformułujemy najpierw odpowiednik *Mocnego Prawa Wielkich Liczb* (PWL) dla łańcuchów Markowa. Rozważmy funkcję $f : \mathcal{X} \rightarrow \mathbb{R}$. Wartość oczekiwana funkcji f względem rozkładu π jest określona jako całka

$$\mathbb{E}_\pi f = \int_{\mathcal{X}} f(x) \pi(dx).$$

Jeżeli rozkład π ma gęstość p względem miary Lebesgue’a to jest to „zwyczajna całka”,

$$\mathbb{E}_\pi f = \int_{\mathcal{X}} f(x) p(x) dx.$$

W przypadku *dyskretnej* przestrzeni \mathcal{X} jest to suma

$$\mathbb{E}_\pi f = \sum_{x \in \mathcal{X}} f(x) \pi(x).$$

Jeśli założymy π jest rozkładem stacjonarnym łańcucha Markowa X_n , to możemy oczekiwać, że zachodzi zbieżność średnich do granicznej wartości oczekiwanej,

$$\frac{1}{n} \sum_{i=0}^{n-1} f(X_i) \longrightarrow \mathbb{E}_\pi f \quad (n \rightarrow \infty) \quad (10.2)$$

z prawdopodobieństwem 1. W istocie, można udowodnić (10.2) przy pewnych dodatkowych założeniach. Mówimy wtedy, że zachodzi PWL lub Mocne Twierdzenie Ergodyczne. Ze względu na zastosowania MCMC, wymagamy aby (10.2) zachodziło dla *dowolnego rozkładu początkowego* ξ (a nie dla $\xi = \pi$, czyli dla łańcucha stacjonarnego). Jedną z wersji PWL dla łańcuchów Markowa przedstawimy, wraz z pięknym i prostym dowodem, w Rozdziale 14 rozdziale (Twierdzenie 14.3).

Centralne Twierdzenie Graniczne (CTG) dla łańcuchów Markowa ma tezę następującej postaci. Dla dowolnego rozkładu początkowego ξ zachodzi zbieżność według rozkładu:

$$\frac{1}{\sqrt{n}} \left(\sum_{i=0}^{n-1} [f(X_i) - \mathbb{E}_\pi f] \right) \longrightarrow N(0, \sigma_{\text{as}}^2(f)) \quad (n \rightarrow \infty). \quad (10.3)$$

Liczba $\sigma_{\text{as}}^2(f)$, zwana asymptotyczną wariancją, nie zależy od rozkładu początkowego ξ , zależy zaś od macierzy przejścia P i funkcji f . Ponadto można udowodnić następujący fakt: dla dowolnego rozkładu początkowego ξ ,

$$\frac{1}{n} \text{Var}_{\xi} \left(\sum_{i=0}^{n-1} f(X_i) \right) \longrightarrow \sigma_{\text{as}}^2(f). \quad (10.4)$$

Przy pewnych dodatkowych założeniach, asymptotyczną wariancję można wyrazić w terminach „stacjonarnych kowariancji” jak następuje. Mamy

$$\sigma_{\text{as}}^2(f) = \text{Var}_{\pi} f(X_0) + 2 \sum_{n=1}^{\infty} \text{Cov}_{\pi}(f(X_0), f(X_n)), \quad (10.5)$$

gdzie Var_{π} i Cov_{π} oznaczają, oczywiście, wariancję i kowariancję obliczoną przy założeniu, że łańcuch jest stacjonarny. Niech

$$\begin{aligned} \sigma^2(f) &= \text{Var}_{\pi} f(X_0); \\ \sigma^2(f) \rho_n(f) &= \text{Cov}_{\pi}[f(X_0), f(X_n)], \end{aligned}$$

Przyjmijmy jeszcze, że $\rho_n(f) = \rho_{-n}(f)$. To określenie jest naturalne bo łańcuch stacjonarny można „przedłużyć wstecz”. Wzór (10.5) można przepisać tak:

$$\sigma_{\text{as}}^2(f) = \sigma^2(f) \left(1 + 2 \sum_{n=1}^{\infty} \rho_n(f) \right) = \sigma^2(f) \sum_{n=-\infty}^{\infty} \rho_n(f).$$

Później pojawi się kilka innych wzorów na asymptotyczną wariancję.

W tym miejscu chcę podkreślić różnicę między stacjonarną wariancją $\sigma^2(f) = \text{Var}_{\pi} f = \text{Var}_{\pi} f(X_n)$ i asymptotyczną wariancją $\sigma_{\text{as}}^2(f)$. W większości zastosowań kowariancje we wzorze (10.5) są dodatnie (zmienne losowe $f(X_0)$ i $f(X_k)$ są dodatnio skorelowane). W rezultacie $\sigma_{\text{as}}^2(f)$ jest dużo większa od $\sigma^2(f)$. To jest cena, którą płacimy za używanie łańcucha Markowa zamiast ciągu zmiennych niezależnych, jak w Rozdziale 8.

Zauważmy, że algorytmy Monte Carlo przeważnie mają za zadanie obliczyć pewną wartość oczekiwaną, a więc wielkość postaci $\theta = \mathbb{E}_{\pi} f$. Jeśli potrafimy generować łańcuch Markowa zbieżny do π , to naturalnym estymatorem $\mathbb{E}_{\pi} f$ jest

$$\hat{\theta}_n = \frac{1}{n} \sum_{i=0}^{n-1} f(X_i).$$

Możemy tezy PWL oraz CTG zapisać w skrócie tak:

$$\begin{aligned} \hat{\theta}_n &\longrightarrow_{\text{p.n.}} \theta \quad (n \rightarrow \infty), \\ \sqrt{n} (\hat{\theta}_n - \theta) &\longrightarrow_d N(0, \sigma_{\text{as}}^2(f)), \quad (n \rightarrow \infty). \end{aligned}$$

PWL gwarantuje *zgodność* estymatora, a więc w pewnym sensie *poprawność* metody. Jest to, rzecz jasna, zaledwie wstęp do dokładniejszej analizy algorytmu. Zauważmy, że CTG może służyć do budowania asymptotycznych przedziałów ufności dla estymowanej wielkości θ : jeśli przyjmujemy poziom ufności $1 - \alpha$ i dobierzemy odpowiedni kwantyl z rozkładu normalnego, to

$$\lim_{n \rightarrow \infty} \mathbb{P} \left(|\hat{\theta}_n - \theta| \leq \frac{z \sigma_{\text{as}}(f)}{\sqrt{n}} \right) = \Phi(z) - \Phi(-z) = 1 - \alpha.$$

Oczywiście, $\Phi(z)$ oznacza to dystrybuantę rozkładu $N(0, 1)$ i $\Phi(z) = 1 - \alpha/2$. Potrzebne jest jeszcze oszacowanie asymptotycznej wariancji $\sigma_{\text{as}}^2(f)$, co nie jest wcale łatwe. Co więcej, CTG nie daje żadnej informacji o tym, dla jakich n przybliżenie rozkładem normalnym jest rozsądne.

Graniczne zachowanie wariancji estymatora $\hat{\theta}_n$ wyjaśnia wzór (10.4). Uzupełnijmy to (pomijając chwilowo uzasadnienie) opisem granicznego zachowania *obciążenia*: przy $n \rightarrow \infty$,

$$\begin{aligned}\mathrm{Var}_\xi(\hat{\theta}_n) &= \frac{1}{n} \sigma_{\mathrm{as}}^2(f) + o\left(\frac{1}{n}\right), \\ \mathbb{E}_\xi \hat{\theta}_n - \theta &= O\left(\frac{1}{n}\right).\end{aligned}$$

Oczywiście, naturalną miarą jakości estymatora jest *błąd średniokwadratowy* (BŚK). Ponieważ BŚK jest sumą wariancji i *kwadratu* obciążenia to, przynajmniej w granicy dla $n \rightarrow \infty$, wariancja ma dominujący wpływ, zaś obciążenie staje się zaniedbywalne:

$$\mathbb{E}_\xi (\hat{\theta}_n - \theta)^2 = \frac{1}{n} \sigma_{\mathrm{as}}^2(f) + o\left(\frac{1}{n}\right). \quad (10.6)$$

Powtórzmy jednak zastrzeżenie dotyczące wszystkich wyników zawartych w tym podrozdziale, Wzór (10.6) tylko sugeruje pewne przybliżenie interesującej nas wielkości.

Sławne i ważne twierdzenia graniczne sformułowane w tym podrozdziale nie są, niestety, całkowicie zadowalającym narzędziem analizy algorytmów Monte Carlo. Algorytmy wykorzystujące łańcuchy Markowa są użyteczne wtedy, gdy osiągają wystarczającą dokładność dla liczby kroków *n znikomo małej* w porównaniu z rozmiarem przestrzeni stanów. W przeciwnym przypadku można po prostu deterministycznie „przejrzeć wszystkie stany” i dokładnie obliczyć interesującą nas wielkość.

Niemniej, twierdzenia graniczne są interesujące z jakościowego punktu widzenia. Ponadto, ważne dla nas oszacowania dotyczące zachowania łańcucha w *skończonym czasie* można porównać z wielkościami granicznymi, aby ocenić ich jakość.

11. Markowskie Monte Carlo II. Podstawowe algorytmy

11.1. Odwracalność

Najważniejsze algorytmy MCMC są oparte na idei odwracalności łańcucha Markowa.

Definicja 11.1. Łańcuch o jądrze P jest odwracalny względem rozkładu prawdopodobieństwa π , jeśli dla dowolnych $A, B \subseteq \mathcal{X}$ mamy

$$\int_A \pi(dx) P(x, B) = \int_B \pi(dy) P(y, A).$$

W skrócie,

$$\pi(dx) P(x, dy) = \pi(dy) P(y, dx).$$

Odwracalność implikuje, że rozkład π jest stacjonarny. Jest to dlatego ważne, że sprawdzanie odwracalności jest stosunkowo łatwe.

Twierdzenie 11.1. *Jeśli $\pi(dx) P(x, dy) = \pi(dy) P(y, dx)$ to $\pi P = P$.*

Dowód. $\pi P(B) = \int_{\mathcal{X}} \pi(dx) P(x, B) = \int_B \pi(dy) P(y, \mathcal{X}) = \int_B \pi(dy) = \pi(B)$. □

11.2. Algorytm Metropolisa-Hastingsa

To jest pierwszy historycznie i wciąż najważniejszy algorytm MCMC. Zakładamy, że umiemy generować łańcuch Markowa z pewnym jądrem q . Pomysł Metropolisa polega na tym, żeby zmodyfikować ten łańcuch wprowadzając specjalnie dobraną regułę akceptacji w taki sposób, żeby wymusić zbieżność do zadanego rozkładu π . W dalszym ciągu systematycznie utożsamiamy rozkłady prawdopodobieństwa z ich gęstościami, aby nie mnożyć oznaczeń. Mamy zatem:

- Rozkład docelowy: $\pi(dx) = \pi(x)dx$.
- Rozkład „propozycji”: $q(x, dy) = q(x, y)dy$.
- Reguła akceptacji:

$$a(x, y) = \frac{\pi(y)q(y, x)}{\pi(x)q(x, y)} \wedge 1.$$

Algorytm Metropolisa-Hastingsa (MH) interpretujemy jako „błądzenie losowe” zgodnie z jądrem przejścia q , zmodyfikowane poprzez odrzucanie niektórych ruchów, przy czym reguła akceptacji/odrzucania zależy w specjalny sposób od π . Pojedynczy krok algorytmu jest następujący.

Listing.

```

function KrokMH( $x$ )
  Gen  $y \sim q(x, \cdot)$ ; { propozycja }
  Gen  $U \sim U(0, 1)$ 
  if  $U > a(x, y)$  then
    begin
       $y := x$ ; { ruch odrzucony z pr-stwem  $1 - a(x, y)$  }
    end
   $KrokMH := y$ 

```

Graficznie to można przedstawić w takiej postaci.

$$X_n = x \text{ ar } [d] \text{ ar } [dl]_{a(x,y)} y \sim q(x, \cdot) \text{ ar } [dr]^{1-a(x,y)} X_{n+1} = y X_{n+1} = x$$

Łańcuch Markowa powstaje zgodnie z następującym schematem:

Listing.

```

Gen  $X_0 \sim \pi_0$ ; { start }
for  $n := 1$  to  $\infty$ 
  begin
     $X_n := KrokMH(X_{n-1})$  { krok }
  end

```

Oczywista jest analogia z podstawową metodą eliminacji. Zasadnicza różnica polega na tym, że w algorytmie MH nie „odrzucaamy” zmiennej losowej, tylko „odrzucaamy propozycje ruchu” – i stoimy w miejscu.

Jądro przejścia M-H jest następujące:

$$P(x, B) = \int_B dy q(x, y) a(x, y) + \mathbb{I}(x \in B) \int_{\mathcal{X}} dy q(x, y) [1 - a(x, y)].$$

Dla przestrzeni skończonej, jądro łańcucha MH redukuje się do macierzy prawdopodobieństw przejścia. Wzór jest w tym przypadku bardzo prosty: dla $x \neq y$,

$$P(x, y) = q(x, y) a(x, y).$$

Twierdzenie 11.2. *Jądro przejścia MH jest odwracalne względem π .*

Dowód. Ograniczmy się do przestrzeni skończonej, żeby nie komplikować oznaczeń. W ogólnym przypadku dowód jest w zasadzie taki sam, tylko napisy stają się mniej czytelne. Niech (bez straty ogólności)

$$a(x, y) = \frac{\pi(y)q(y, x)}{\pi(x)q(x, y)} \leq 1, \quad a(y, x) = 1.$$

Wtedy

$$\begin{aligned}
 \pi(x)P(x, y) &= \pi(x)q(x, y)a(x, y) \\
 &= \pi(x)q(x, y) \frac{\pi(y)q(y, x)}{\pi(x)q(x, y)} \\
 &= \pi(y)q(y, x) \\
 &= \pi(y)P(y, x) \quad \text{bo} \quad a(y, x) = 1.
 \end{aligned}$$

□

Uwagi historyczne:

- Metropolis w 1953 zaproponował algorytm, w którym zakłada się symetrię rozkładu propozycji, $q(x, y) = q(y, x)$. Warto zauważyć, że wtedy łańcuch odpowiadający q (błądzenie bez eliminacji ruchów) ma rozkład stacjonarny *jednostajny*. Reguła akceptacji przybiera postać

$$a(x, y) = \frac{\pi(y)}{\pi(x)} \wedge 1.$$

- Hastings w 1970 uogólnił rozważania na przypadek niesymetrycznego q .

Uwaga ważna:

- Algorytm MH wymaga znajomości gęstości π tylko z dokładnością do proporcjonalności, **bez stałej normującej**.

11.3. Próbnik Gibbsa

Drugim podstawowym algorytmem MCMC jest próbnik Gibbsa (PG) (*Gibbs Sampler*, GS). Załóżmy, że przestrzeń na której żyje docelowy rozkład π ma strukturę produktową: $\mathcal{X} = \prod_{i=1}^d \mathcal{X}_i$. Przyjmijmy następujące oznaczenia:

- Jeśli $\mathcal{X} \ni x = (x_i)_{i=1}^d$ to $x_{-i} = (x_j)_{j \neq i}$: wektor z pominiętą i -tą współrzędną.
- Rozkład docelowy (gęstość): $\pi(dx) = \pi(x)dx$.
- Pełne rozkłady warunkowe (*full conditionals*):

$$\pi(x_i | x_{-i}) = \frac{\pi(x)}{\pi(x_{-i})}$$

Mały krok PG jest zmianą i -tej współrzędnej (wylosowaniem nowej wartości z rozkładu warunkowego):

$$\begin{array}{c} x = (x_1, \dots, x_i, \dots, x_d) \\ \downarrow \\ \text{Gen } y_i \sim \pi(\cdot | x_{-i}) \\ \downarrow \\ Y = (x_1, \dots, y_i, \dots, x_d). \end{array}$$

Prawdopodobieństwo przejścia małego kroku PG (w przypadku przestrzeni skończonej) jest takie:

$$P_i(x, y) = \pi(y_i | x_{-i}) \mathbb{I}(x_{-i} = y_{-i}).$$

Twierdzenie 11.3. *Mały krok PG jest π -odwracalny.*

Dowód. Niech $x_{-i} = y_{-i}$. Wtedy

$$\begin{aligned} \pi(x) P_i(x, y) &= \pi(x) \pi(y_i | x_{-i}) \\ &= \pi(x_{-i}) \pi(x_i | x_{-i}) \pi(y_i | x_{-i}) \\ &= \pi(y_{-i}) \pi(x_i | x_{-i}) \pi(y_i | y_{-i}) \\ &= \pi(y) P_i(y, x). \end{aligned}$$

(skorzystaliśmy z symetrii).

□

Oczywiście, trzeba jeszcze zadbać o to, żeby łańcuch generowany przez PG był nieprzywiel-
dny. Trzeba zmieniać wszystkie współrzędne, nie tylko jedną. Istnieją dwie zasadnicze odmiany
próbnika Gibbsa, różniące się sposobem wyboru współrzędnych do zmiany.

— Losowy wybór współrzędnych, „LosPG”.

— Systematyczny wybór współrzędnych, „SystemPG”.

Losowy PG. Wybieramy współrzędną i -tą z prawdopodobieństwem $c(i)$. .

Listing.

```
function LosPG(x)
  Gen  $i \sim c(\cdot)$ ;
  Gen  $y_i := \pi(\cdot | x_{-i})$ ; { zmieniamy i-tą współrzędną }
   $y_{-i} := x_{-i}$ ; { wszystkie inne współrzędne pozostawiamy bez zmian }
  LosPG :=  $y$ 
```

Jądro przejścia w „dużym” kroku losowego PG jest takie:

$$P = \sum_{i=1}^d c(i) P_i,$$

LosPG jest **odwracalny**.

Systematyczny PG. Współrzędne są zmieniane w porządku cyklicznym.

Listing.

```
function SystemPG(x)
  begin
    Gen  $y_1 \sim \pi(\cdot | x_2, \dots, x_d)$ ;
    Gen  $y_2 \sim \pi(\cdot | y_1, x_3, \dots, x_d)$ ;
    ...
    Gen  $y_d \sim \pi(\cdot | y_1, \dots, y_{d-1})$ ;
    SystemPG :=  $y$ 
  end
```

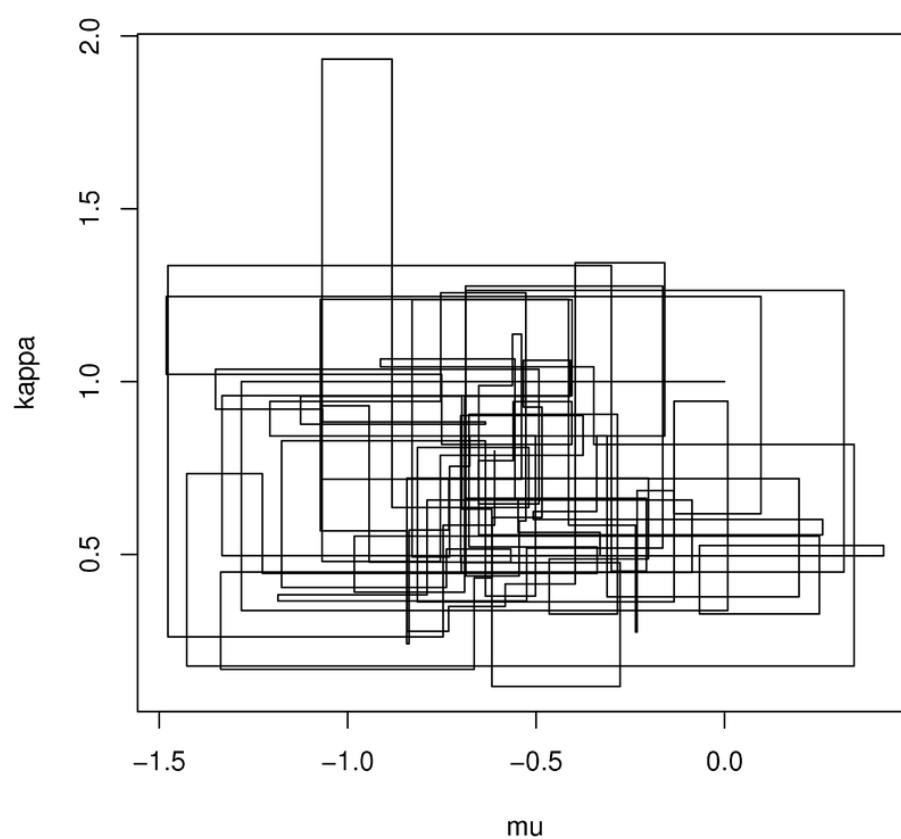
Jądro przejścia w „dużym” kroku systematycznego PG jest następujące.

$$P = P_1 P_2 \cdots P_d,$$

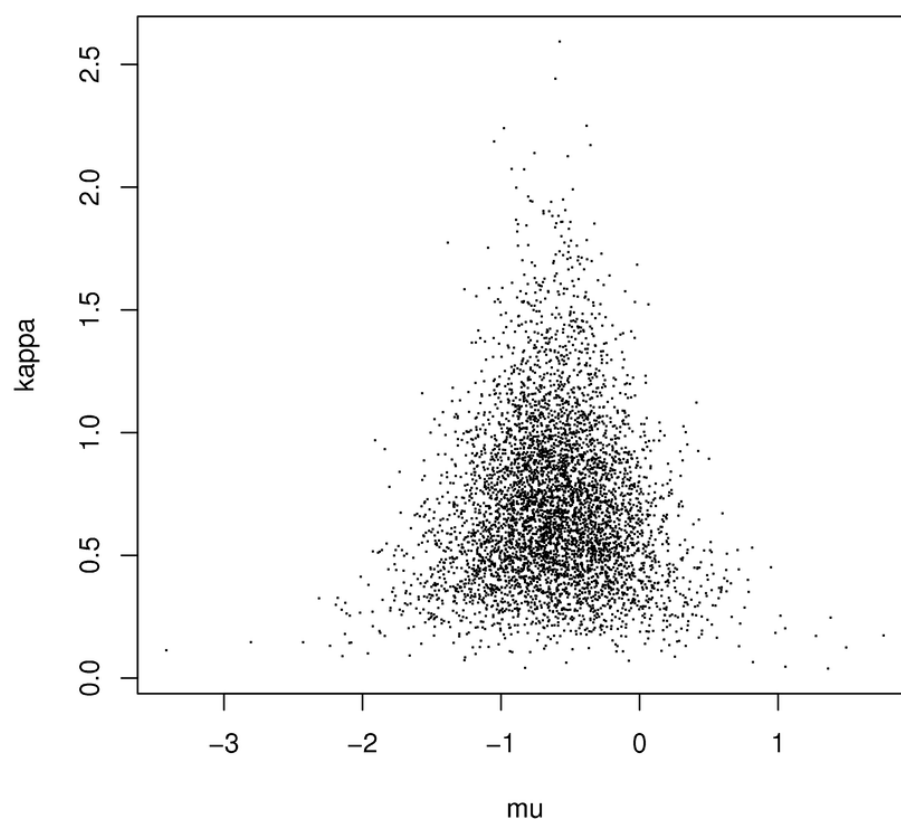
Systematyczny PG **nie jest odwracalny**. Ale jest π -stacjonarny, bo $\pi P_1 P_2 \cdots P_d = \pi$.

Przy projektowaniu konkretnych realizacji PG pojawia się szereg problemów, ważnych za-
równo z praktycznego jak i teoretycznego punktu widzenia. Jak wybrać rozkład $c(\cdot)$ w losowym
PG? Jest raczej jasne, że niektóre współrzędne powinny być zmieniane częściej, a inne rza-
dziej. Jak dobrać kolejność współrzędnych w systematycznym PG? Czy ta kolejność ma wpływ
na tempo zbieżności łańcucha. Wreszcie, w wielu przykładach można zmieniać całe „bloki”
współrzędnych na raz.

Systematyczny PG jest uważany za bardziej efektywny i częściej stosowany w praktyce. Z
drugiej strony jest trudniejszy do analizy teoretycznej, niż losowy PG.



Rysunek 11.1. Trajektoria próbnika Gibbsa w przestrzeni dwuwymiarowej.



Rysunek 11.2. Chmurka punktów wygenerowanych przez próbnik Gibbsa.

12. Markowskie Monte Carlo III. Przykłady zastosowań

12.1. Statystyka bayesowska

Algorytmy MCMC zrewolucjonizowały statystykę bayesowską. Stworzyły możliwość obliczania (w przybliżeniu) rozkładów *a posteriori* w sytuacji, gdy dokładne, analityczne wyrażenia są niedostępne. W ten sposób statystycy uwolnili się od konieczności używania nadmiernie uproszczonych modeli. Zaczęli śmiało budować modele coraz bardziej realistyczne, zwykle o strukturze hierarchicznej. Przedstawię to na jednym dość typowym przykładzie, opartym na pracy [21]. Inne przykłady i doskonały wstęp do tematyki zastosowań MCMC można znaleźć w pracy Geyera [7].

12.1.1. Hierarchiczny model klasyfikacji

Przykład 12.1 (Statystyka małych obszarów). Zaczniemy od opisu problemu tak zwanych „małych obszarów”, który jest dość ważny w dziedzinie badań reprezentacyjnych, czyli w tak zwanej „statystyce oficjalnej”. Wyobraźmy sobie, że w celu zbadania kondycji przedsiębiorstw losuje się próbkę, która liczy (powiedzmy) 3500 przedsiębiorstw z całego kraju. Na podstawie tej losowej próbki można dość wiarygodnie szacować (estymować) pewne parametry opisujące populację przedsiębiorstw w kraju. Czy jednak można z rozsądną dokładnością oszacować sprzedaż w powiecie garwolińskim? W Polsce mamy ponad 350 powiatów. Na jeden powiat przypada średnio 10 przedsiębiorstw wybranych do próbki. *Małe obszary* to pod-populacje w których rozmiar próbki nie jest wystarczający, aby zastosować „zwykłe” estymatory (średnie z próbki). Podejście bayesowskie pozwala „pożyczać informację” z innych obszarów. Zakłada się, że z każdym małym obszarem związany jest nieznany parametr, który staramy się estymować. Obserwacje pochodzące z określonego obszaru mają rozkład prawdopodobieństwa zależny od odpowiadającego temu obszarowi parametru. Parametry, zgodnie z filozofią bayesowską, traktuje się jak zmienne losowe. W najprostszej wersji taki model jest zbudowany w następujący sposób.

Model bayesowski

- $y_{ij} \sim N(\theta_i, \sigma^2)$ – badana cecha dla j -tej wylosowanej jednostki i -tego obszaru, ($j = 1, \dots, n_i$), ($i = 1, \dots, k$),
- $\theta_i \sim N(\mu, v^2)$ – interesująca nas średnia w i -tym obszarze,
- μ – średnia w całej populacji.

Ciekawe, że ten sam model pojawia się w różnych innych zastosowaniach, na przykład w matematyce ubezpieczeniowej. Przytoczymy klasyczny rezultat dotyczący tego modelu, aby wyjaśnić na czym polega wspomniane „pożyczanie informacji”.

Estymator bayesowski

W modelu przedstawionym powyżej, łatwo obliczyć estymator bayesowski (przy kwadratowej funkcji straty), czyli wartość oczekiwaną *a posteriori*. Następujący wzór jest bardzo dobrze znany specjalistom od małych obszarów i aktuariuszom.

$$\hat{\theta}_i = \mathbb{E}(\theta_i|y) = z_i \bar{y}_i + (1 - z_i)\mu, \quad z_i = \frac{n_i v^2}{n_i v^2 + \sigma^2}.$$

Estymator bayesowski dla i -tego obszaru jest średnią ważoną \bar{y} (estymatora opartego na danych z tego obszaru) i wielkości μ , która opisuje całą populację, a nie tylko i -ty obszar. Niestety, proste estymator napisany powyżej zależy od parametrów μ , σ i v , które w praktyce są nieznane i które trzeba estymować. Konsekwentnie bayesowskie podejście polega na traktowaniu również tych parametrów jako zmiennych losowych, czyli nałożeniu na nie rozkładów *a priori*. Powstaje w ten sposób model hierarchiczny.

Hierarchiczny model bayesowski

Uzupełnijmy rozpatrywany powyżej model, dobudowując „wyższe piętra” hierarchii. potraktujemy mianowicie parametry rozkładów *a priori*: μ , σ i v jako zmienne losowe i wyspecyfikujemy ich rozkłady *a priori*.

- $y_{ij} \sim N(\theta_i, \sigma^2)$,
- $\theta_i \sim N(\mu, v^2)$,
- $\mu \sim N(m, \tau^2)$,
- $\sigma^{-2} \sim \text{Gamma}(p, \lambda)$,
- $v^{-2} \sim \text{Gamma}(q, \kappa)$.

Zakładamy przy tym, że μ , σ i v są *a priori* niezależne (niestety, są one zależne *a posteriori*). Na szczycie hierarchii mamy „hiperparametry” m , τ , p , λ , q , κ , o których musimy założyć, że są znanymi liczbami.

Łączny rozkład prawdopodobieństwa wszystkich zmiennych losowych w modelu ma postać

$$p(y, \theta, \mu, \sigma^{-2}, v^{-2}) = p(y|\theta, \sigma^{-2})p(\theta|\mu, v^{-2})p(\mu)p(\sigma^{-2})p(v^{-2}).$$

We wzorze powyżej i w dalej traktujemy (trochę nieformalnie) σ^{-2} i v^{-2} jako pojedyncze symbole nowych zmiennych, żeby nie mnożyć oznaczeń. Rozkład prawdopodobieństwa *a posteriori* jest więc taki:

$$p(\theta, \mu, \sigma^{-2}, v^{-2}|y) = \frac{p(y, \theta, \mu, \sigma^{-2}, v^{-2})}{p(y)}.$$

To jest rozkład „docelowy” π , na przestrzeni $\mathcal{X} = \mathbb{R}^{k+3}$, ze nieznaną stałą normującą $1/p(y)$. Choć wygląda na papierze dość prosto, ale obliczenie rozkładów brzegowych, wartości oczekiwanych i innych charakterystyk jest, łagodnie mówiąc, trudne.

12.1.2. Próbnik Gibbsa w modelu hierarchicznym

Rozkłady warunkowe poszczególnych współrzędnych są proste i łatwe do generowania. Można te rozkłady „odczytać” uważnie patrząc na rozkład łączny:

$$\begin{aligned}
 p(\theta, \mu, v^{-2}, \sigma^{-2} | y) &\propto (\sigma^{-2})^{n/2} \exp \left\{ -\frac{\sigma^{-2}}{2} \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \theta_i)^2 \right\} \\
 &\cdot (v^{-2})^{k/2} \exp \left\{ -\frac{v^{-2}}{2} \sum_{i=1}^k (\theta_i - \mu)^2 \right\} \\
 &\cdot \exp \left\{ -\frac{\tau^{-2}}{2} (\mu - m)^2 \right\} \\
 &\cdot (\sigma^{-2})^{q-1} \exp\{-\kappa \sigma^{-2}\} \\
 &\cdot (v^{-2})^{p-1} \exp\{-\lambda v^{-2}\}.
 \end{aligned}$$

Dla ustalenia uwagi zajmijmy się rozkładem warunkowym zmiennej v^{-2} . Kolorem **niebieskim** oznaczyliśmy te czynniki łącznej gęstości, które zawierają v^{-2} . Pozostałe, czarne czynniki traktujemy jako stałe. Stąd widać, jak wygląda rozkład warunkowy v^{-2} , „przynajmniej z dokładnością do proporcjonalności:

$$\begin{aligned}
 p(v^{-2} | y, \theta, \mu, \sigma^{-2}) &\propto (v^{-2})^{k/2+p-1} \\
 &\cdot \exp \left\{ -\left(\frac{1}{2} \sum_{i=1}^k (\theta_i - \mu)^2 + \lambda \right) v^{-2} \right\}.
 \end{aligned}$$

Jest to zatem rozkład $\text{Gamma}(k/2 + p, \sum_{i=1}^k (\theta_i - \mu)^2 / 2 + \lambda)$. Zupełnie podobnie rozpoznajemy inne (pełne) rozkłady warunkowe:

$$\begin{aligned}
 v^{-2} | y, \theta, \mu, \sigma^{-2} &\sim \text{Gamma} \left(\frac{k}{2} + p, \frac{1}{2} \sum_{i=1}^k (\theta_i - \mu)^2 + \lambda \right), \\
 \sigma^{-2} | y, \theta, \mu, v^{-2} &\sim \text{Gamma} \left(\frac{n}{2} + q, \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \theta_i)^2 + \kappa \right), \\
 \mu | y, \theta, \sigma^{-2}, v^{-2} &\sim \text{N} \left(\frac{k\tau^2}{k\tau^2 + v^2} \bar{\theta} + \frac{v^2}{k\tau^2 + v^2} m, \frac{\tau^2 v^2}{k\tau^2 + v^2} \right), \\
 \theta_i | y, \theta_{-i}, \mu, \sigma^{-2}, v^{-2} &\sim \text{N} \left(\frac{nv^2}{nv^2 + \sigma^2} \bar{y}_i + \frac{\sigma^2}{nv^2 + \sigma^2} \mu, \frac{v^2 \sigma^2}{nv^2 + \sigma^2} \right),
 \end{aligned}$$

gdzie, rzecz jasna, $n = \sum n_i$, $\bar{\theta} = \sum_i \theta_i / k$ i $\theta_{-i} = (\theta_k)_{k \neq i}$. Zwróćmy uwagę, że współrzędne wektora θ są warunkowo niezależne (pełny rozkład warunkowy θ_i nie zależy od θ_{-i}). Dzięki temu możemy w próbniku Gibbsa potraktować θ jako cały „blok” współrzędnych i zmieniać „na raz”.

Próbnik Gibbsa ma w tym modelu przestrzeń stanów \mathcal{X} składającą się z punktów $x = (\theta, \mu, \sigma^{-2}, v^{-2}) \in \mathbb{R}^{k+3}$. Reguła przejścia próbnika w wersji systematycznej (duży krok „SystemPG”),

$$\underbrace{(\theta, \mu, \sigma^{-2}, v^{-2})}_{X_t} \mapsto \underbrace{(\theta, \mu, \sigma^{-2}, v^{-2})}_{X_{t+1}},$$

jest złożona z następujących „małych kroków”:

- Wylosuj $v^{-2} \sim p(v^{-2}|y, \theta, \mu, \sigma^{-2}) = \text{Gamma}(\dots)$,
- Wylosuj $\sigma^{-2} \sim p(\sigma^{-2}|y, \theta, \mu, v^{-2}) = \text{Gamma}(\dots)$,
- Wylosuj $\mu \sim p(\mu|y, \theta, \sigma^{-2}, v^{-2}) = \text{N}(\dots)$,
- Wylosuj $\theta \sim p(\theta|y, \mu, \sigma^{-2}, v^{-2}) = \text{N}(\dots)$.

Łańcuch Markowa jest zbieżny do rozkładu *a posteriori*:

$$X_t \rightarrow \pi(\cdot) = p(\theta, \mu, \sigma^{-2}, v^{-2}|y).$$

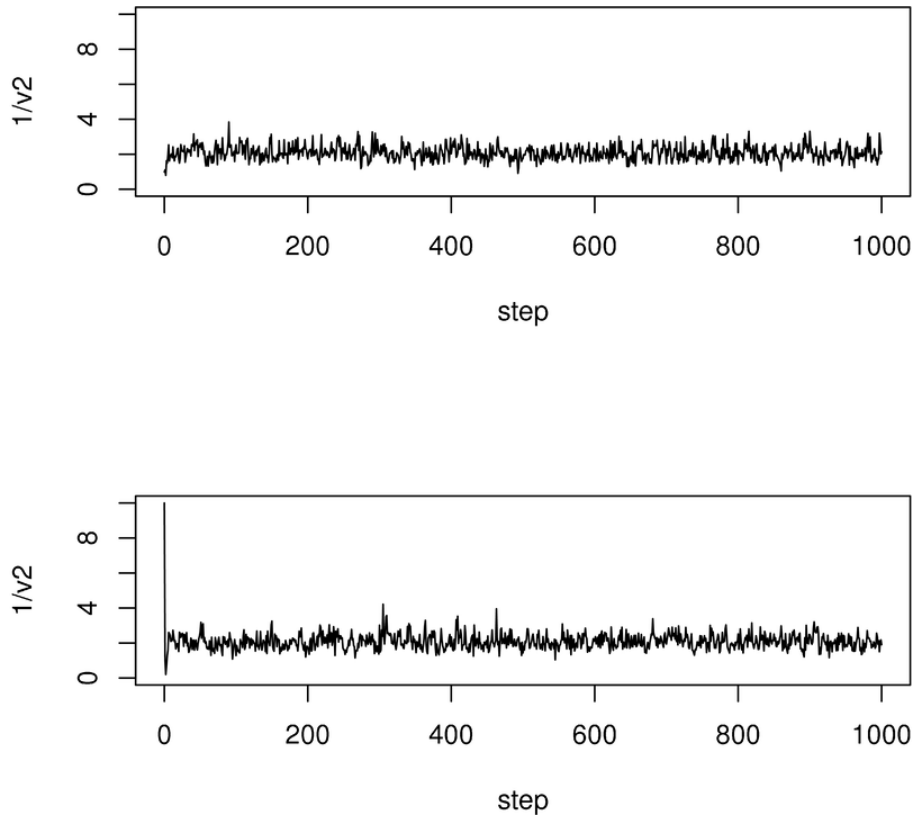
Najbardziej interesujące są w tym hierarchicznym modelu zmienne θ_i (pozostałe zmienne można uznać za „parametry zklócające”). Dla ustalenia uwagi zajmijmy się zmienną θ_1 (powiedzmy, wartością średnią w pierwszym małym obszarze). **Estymator bayesowski** jest to wartość oczekiwana *a posteriori* tej zmiennej:

$$\mathbb{E}(\theta_1|y) = \int \dots \int \theta_1 p(\theta, \mu, \sigma^{-2}, v^{-2}|y) d\theta_2 \dots d\theta_k d\mu d\sigma^{-2} dv^{-2}.$$

Aproksymacją MCMC interesującej nas wielkości są średnie wzdłuż trajektorii łańcucha:

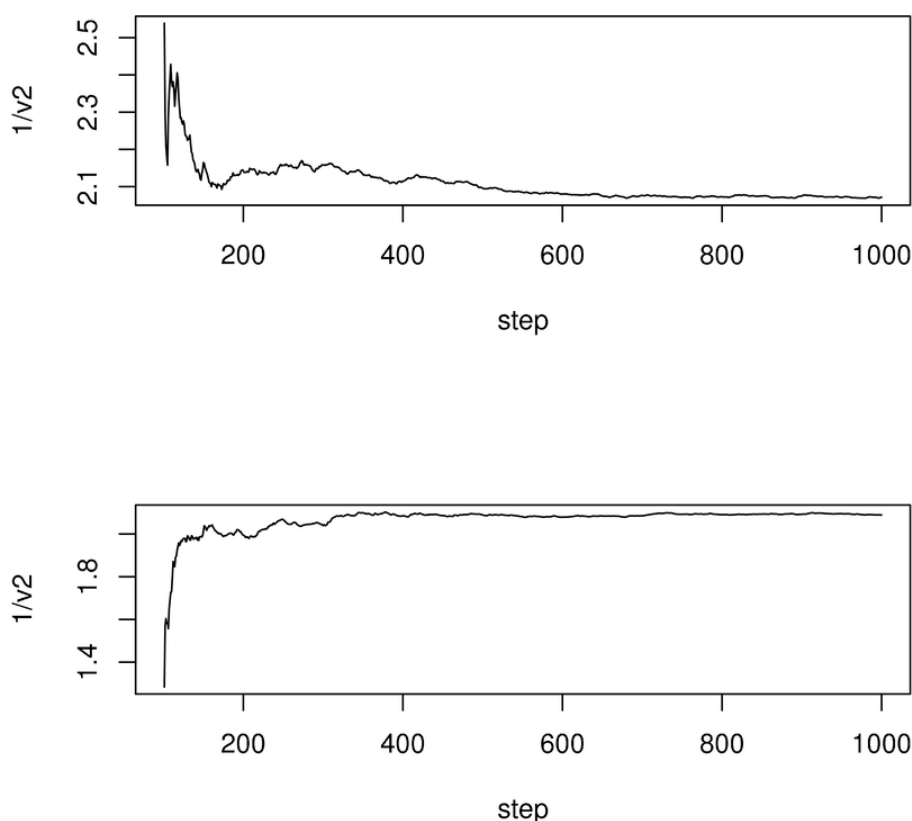
$$\theta_1(X_0), \theta_1(X_1), \dots, \theta_1(X_t), \dots$$

gdzie $\theta_1(x) = \theta_1$ dla $x = (\theta_1, \dots, \theta_k, \mu, \sigma^{-2}, v^{-2})$.



Rysunek 12.1. Trajektorie zmiennej v^{-2} dla dwóch punktów startowych.

Na Rysunku 12.1 pokazane są dwie przykładowe trajektorie współrzędnej v^{-2} dla PG poruszającego się po przestrzeni $k + 3 = 1003$ wymiarowej (model uwzględniający 1000 małych obszarów). Dwie trajektorie odpowiadają dwóm różnym punktom startowym. Dla innych zmiennych rysunki wyglądają bardzo podobnie. Uderzające jest to, jak szybko trajektoria zdaje się „osiągać” rozkład stacjonarny, przynajmniej wizualnie. Na Rysunku 12.2 pokazane są kolejne „skumulowane” średnie dla tych samych dwóch trajektorii zmiennej v^{-2} .



Rysunek 12.2. Skumulowane średnie zmiennej v^{-2} dla dwóch punktów startowych.

12.2. Estymatory największej wiarygodności

Metody MCMC w połączeniu z ideą losowania istotnego (Rozdział 8) znajdują zastosowanie również w statystyce „częstościowej”, czyli nie-bayesowskiej. Pokażę przykład, w którym oblicza się metodami Monte Carlo estymator największej wiarygodności.

12.2.1. Model auto-logistyczny

Niech $x = (x_1, \dots, x_d)$ będzie wektorem (konfiguracją) binarnych zmiennych losowych na przestrzeni $\mathcal{X} = \{0, 1\}^d$. Rozważmy następujący rozkład Gibbsa:

$$p_{\theta}(x) = \frac{1}{Z(\theta)} \exp \left\{ \sum_{i,j=1}^d \theta_{ij} x_i x_j \right\}.$$

Rolę parametru gra macierz (θ_{ij}) . Zakład się, że jest to macierz symetryczna. Stała normująca $Z(\theta) = \sum_{x \in \mathcal{X}} \exp \left\{ \sum_{i,j=1}^d \theta_{ij} x_i x_j \right\}$ jest typowo (dla dużych d) *niemożliwa do obliczenia*. Stanowi to, jak dalej zobaczymy, poważny problem dla statystyków.

Przykład 12.2 (Statystyka przestrzenna). W zastosowaniach „przestrzennych” indeks $i \in \{1, \dots, d\}$ interpretuje się jako „miejsce”. Zbiór miejsc wyposażony jest w strukturę grafu. Krawędzie łączą miejsca „sąsiadujące”. Piszemy $i \sim j$. Tego typu modele mogą opisywać na przykład rozprzestrzenianie się chorób lub występowanie pewnych gatunków. Wartość $x_i = 1$ oznacza obecność gatunku lub występowanie choroby w miejscu i . Najprostszy model zakła-

da, że każda zmienna x_i zależy tylko od swoich „sąsiadów” i to w podobny sposób w całym rozpatrywanym obszarze. W takim modelu mamy tylko dwa parametry $\theta = (\theta_0, \theta_1)$:

$$\theta_{ij} = \begin{cases} 0 & i \not\sim j, i \neq j; \\ \theta_1 & i \sim j; \\ \theta_0 & i = j. \end{cases}$$

Parametr θ_0 opisuje „skłonność” pojedynczej zmiennej do przyjmowania wartości 1, zaś parametr θ_1 odpowiada za zależność od zmiennych sąsiadujących (zakaźność choroby, powiedzmy). W typowej dla statystyki przestrzennej sytuacji, rozpatruje się nawet dziesiątki tysięcy „miejsc”. Stała $Z(\theta)$ jest wtedy sumą niewyobrażalnie wielu (dokładnie 2^d) składników.

Próbnik Gibbsa w modelu auto-logistycznym

„Pełne” rozkłady warunkowe (*full conditionals*) są w modelu auto-logistycznym bardzo proste:

$$p_\theta(x_i = 1 \mid x_{-i}) = \frac{\exp\left(\theta_{ii} + \sum_{\substack{j=1 \\ j \neq i}}^d \theta_{ij}x_j\right)}{1 + \exp\left(\theta_{ii} + \sum_{\substack{j=1 \\ j \neq i}}^d \theta_{ij}x_j\right)},$$

gdzie $x_{-i} = (x_j, j \neq i)$. Zatem:

— Symulowanie $x \sim p_\theta$ jest łatwe za pomocą próbnika Gibbsa (PG):

$$x_1 \sim p_\theta(x_1 \mid x_{-1}),$$

$$x_2 \sim p_\theta(x_2 \mid x_{-2}), \dots$$

Aproksymacja wiarygodności

Estymator największej wiarygodności obliczany metodą Monte Carlo został zaproponowany w pracy Geyer and Thopmpson (1992, *JRSS*). Rozważmy bardziej ogólną wykładniczą rodzinę rozkładów prawdopodobieństwa:

$$p_\theta(x) = \frac{1}{Z(\theta)} \exp\left[\theta^\top T(x)\right],$$

gdzie $T(x)$ jest wektorem statystyk dostatecznych. Rozkłady autologistyczne tworzą rodzinę wykładniczą. Wystarczy θ ustawić w wektor, statystykami $T(x)$ są x_i i $x_i x_j$. Logarytm wiarygodności $L(\theta) = \log p_\theta(x)$ jest dany pozornie prostym wzorem:

$$L(\theta) = \theta^\top T(x) - \log Z(\theta).$$

Kłopot jest ze stałą normującą $Z(\theta)$, która w bardziej skomplikowanych modelach jest trudna lub wręcz niemożliwa do obliczenia. Wspomnieliśmy, że tak właśnie jest dla typowych modeli autologistycznych. Istnieje jednak prosty sposób aproksymacji stałej $Z(\theta)$ metodą MC. Istotnie, wybierzmy (w zasadzie dowolny) punkt θ_* i zauważmy, że

$$Z(\theta) = \sum_x \exp\left[\theta^\top T(x)\right] = \sum_x \exp\left[(\theta - \theta_*)^\top T(x)\right] \exp\left[\theta_*^\top T(x)\right].$$

Pozwala to wyrazić wielkość $Z(\theta)/Z(\theta_*)$ jako wartość oczekiwaną względem rozkładu p_{θ_*} :

$$\begin{aligned} Z(\theta)/Z(\theta_*) &= \sum_x \exp\left[(\theta - \theta_*)^\top T(x)\right] p_{\theta_*}(x) \\ &= \mathbb{E}_{\theta_*} \exp\left[(\theta - \theta_*)^\top T(x)\right]. \end{aligned}$$

Powyższe rozważania nie są niczym nowym: to jest po prostu przedstawiony w Rozdziale 8 schemat *losowania istotnego*, w specjalnym przypadku rodziny wykładniczej. Jeśli celem jest maksymalizacja wiarygodności $Z(\theta)$, to nieznana stała $Z(\theta_*)$ zupełnie nie przeszkadza (bo interesującą nas funkcję wiarygodności wystarczy znać z dokładnością do stałej).

Jeśli umiemy generować zmienne losowe o rozkładzie p_{θ_*} , to sposób przybliżonego obliczania $\hat{Z}(\theta)$ jest w zasadzie oczywisty. Generujemy próbkę MC: $x_*(k) \sim p_{\theta_*}$, $k = 1, \dots, n_*$ gdzie θ_* jest w zasadzie dowolne, zaś n_* jest możliwie największe. Obliczamy

$$\hat{Z}(\theta)/Z(\theta_*) = \frac{1}{n_*} \sum_{k=1}^{n_*} \exp \left[(\theta - \theta_*)^\top T(x_*(k)) \right].$$

Aproksymacja logarytmu wiarygodności polega na wstawieniu $\hat{Z}(\theta)$ w miejsce $Z(\theta)$ we wzorze na $L(\theta) = \log p_\theta(x)$:

$$\hat{L}_{\text{MC}}(\theta) = \theta^\top T(x) - \log \underbrace{\sum_{k=1}^{n_*} \exp[(\theta - \theta_*)^\top T(x_*(k))]}_{\text{przybliżenie MC } Z(\theta)} + \text{const.}$$

Pozostaje wiele szczegółów do dopracowania. Umiemy obliczać wiarygodność jako funkcję θ , ale trzeba jeszcze znaleźć maksimum tej funkcji. Dobór θ_* nie wpływa na poprawność rozumowania ale może mieć zasadniczy wpływ na efektywność algorytmu. Nie będziemy się w to zagłębiać. Wspomnimy tylko, jak omawiana metoda jest związana z *markowskimi* metodami Monte Carlo, MCMC. Wróćmy do rozkładu auto-logistycznego. Jak losować próbkę z tego rozkładu? Tu przychodzi z pomocą próbnik Gibbsa. W rzeczy samej, mamy do czynienia z kolejną aproksymacją: PG jest tylko *przybliżonym* algorytmem generowania z rozkładu p_{θ_*} , które to losowanie pozwala obliczać wiarygodność *w przybliżeniu*. Momo wszystko, metody MCMC oferują pewne wyjście lepsze niż bezradne rozłożenie rąk.

13. Markowskie Monte Carlo IV. Pola losowe

13.1. Definicje

Niech $(\mathcal{S}, \mathcal{E})$ będzie nieskierowanym grafem. Wyobraźmy sobie, że elementy $s \in \mathcal{S}$ reprezentują „miejsca” w przestrzeni lub na płaszczyźnie, zaś krawędzie grafu łączą miejsca „sąsiadujące” ze sobą. Taka interpretacja jest związana z zastosowaniami do statystyki „przestrzennej” i przetwarzania obrazów. Model, który przedstawimy ma również zupełnie inne interpretacje, ale pozostaniemy przy sugestywnej terminologii „przestrzennej”:

- \mathcal{S} — zbiór miejsc,
- $\{s, t\} \in \mathcal{E}$ — miejsca s i t sąsiadują — będziemy wtedy pisać $s \sim t$,
- $\partial t = \{s : s \sim t\}$ — zbiór sąsiadów miejsca t .

Niech $\Lambda = \{1, \dots, l\}$ będzie skończonym zbiorem. Powiedzmy, że elementy $a \in \Lambda$ są „kolorami” które mogą być przypisane elementom zbioru \mathcal{S} . Konfiguracją nazywamy dowolną funkcję $x : \mathcal{S} \rightarrow \Lambda$. Będziemy mówić że $x_s = x(s)$ jest s -tą współrzędną konfiguracji x i stosować oznaczenia podobne jak dla wektorów:

$$x = (x_s) = (x_s : s \in \mathcal{S}).$$

W zadaniach przetwarzania obrazów, miejsca są pikslami na ekranie i konfigurację utożsamiamy z ich pokolorowaniem, a więc z cyfrową reprezentacją obrazu. Zbiór Λ gra rolę „palety kolorów”. Niekiedy założenie o skończoności zbioru Λ staje się niewygodne. Dla czarno-szaro-białych obrazów „pomalowanych” różnymi odcieniami szarości, wygodnie przyjąć, że $\Lambda = [0, 1]$ lub $\lambda = [0, \infty[$. Tego typu modyfikacje są dość oczywiste i nie będę się nad tym zatrzymywał. Dla ustalenia uwagi, wzory w tym podrozdziale dotyczą przypadku skończonego zbioru „kolorów”. Przestrzenią konfiguracji jest zbiór $\mathcal{X} = \Lambda^{\mathcal{S}}$. Dla konfiguracji x i miejsca t , niech

- $x_{-t} = (x_s : s \neq t) = (x_s : s \in \mathcal{S} \setminus \{t\})$ — konfiguracja z pominiętą t -tą współrzędną,
- $x_{\partial t} = (x_s : s \in \partial t)$ — konfiguracja ograniczona do sąsiadów miejsca t .

Jeśli $H : \mathcal{X} \rightarrow \mathbb{R}$ i $\beta \geq 0$ to **rozkładem Gibbsa** nazywamy rozkład prawdopodobieństwa na przestrzeni konfiguracji dany wzorem

$$\pi_\beta(x) = \frac{1}{Z(\beta)} \exp[-\beta H(x)].$$

Ze względu na inspiracje pochodzące z fizyki statystycznej, funkcję H nazywamy energią, β jest (z dokładnością do stałej) odwrotnością temperatury. Stała normująca wyraża się wzorem

$$Z(\beta) = \sum_{x \in \mathcal{X}} \exp[-\beta H(x)].$$

i jest typowo *niemożliwa do obliczenia*.

Oczywiście, każdy rozkład prawdopodobieństwa π na \mathcal{X} dale się zapisać jako rozkład Gibbsa, jeśli położyć $H(x) = -\log \pi(x)$, $\beta = 1$ i umownie przyjąć, że $-\log 0 = \infty$ (czyli konfiguracje niemożliwe mają nieskończoną energię). Nie o to jednak chodzi. Ciekawe są rozkłady Gibbsa, dla których funkcja energii ma specjalną postać związaną z topologią grafu „sąsiedztw”. Ograniczymy

się do ważnej podklasy markowskich pól losowych (MPL), mianowicie do sytuacji gdy energia jest sumą „oddziaływań” lub „interakcji” między parami miejsc sąsiadujących i składników zależnych od pojedynczych miejsc. Dokładniej, założymy że

$$H(x) = \sum_{s \sim t} V(x_s, x_t) + \sum_s U_s(x_s), \quad (13.1)$$

dla pewnych funkcji $V : \Lambda \times \Lambda \rightarrow \mathbb{R}$ i $U_s : \Lambda \rightarrow \mathbb{R}$. Funkcja $V(x_s, x_t)$ opisuje „potencjał interakcji pomiędzy s i t ”, zaś $U_s(a)$ jest wielkością związaną z „tendencją miejsca s do przybrania koloru a ”. Zwróćmy uwagę, że potencjał V jest jednorodny ($V(a, b)$ zależy tylko od „kolorów” $a, b \in \Lambda$ ale nie od miejsc), zaś $U_s(a)$ może zależeć zarówno od $a \in \Lambda$ jak i od $s \in \mathcal{S}$. W modelach fizyki statystycznej zazwyczaj $U_s(a) = U(a)$ jest jednorodnym „oddziaływaniem zewnętrznym” ale w modelach rekonstrukcji obrazów nie można tego zakładać.

Przykład 13.1 (Model Potts). Niech Λ będzie zbiorem skończonym i

$$H(x) = J \sum_{s \sim t} \mathbb{I}(x_s \neq x_t).$$

Ta funkcja opisuje „tendencję sąsiednich miejsc do przybierania tego samego koloru”. Jeśli $J > 0$ to preferowane są konfiguracje złożone z dużych, jednobarwnych plam.

13.2. Generowanie markowskich pól losowych

Użyteczność MPL w różnorodnych zastosowaniach związana jest z istnieniem efektywnych algorytmów symulacyjnych. Wszystko opiera się na próbniku Gibbsa i następującym prostym fakcie.

Twierdzenie 13.1 (Pełne rozkłady warunkowe dla MPL). *Jeżeli π_β jest rozkładem Gibbsa z energią daną wzorem (13.1), to*

$$\pi_\beta(x_s | x_{-s}) = \pi_\beta(x_s | x_{\partial s}) = \frac{1}{Z_s(\beta)} \exp[-\beta H_s(x)],$$

gdzie

$$H_s(x) = \sum_{t \in \partial s} V(x_s, x_t) + U(x_s),$$

$Z_s(\beta) = \sum_{a \in \Lambda} \exp[H_s(x_{a \rightsquigarrow s})]$. Symbol $x_{a \rightsquigarrow s}$ oznacza konfigurację powstałą z x przez wpisanie koloru a w miejscu s .

Dowód. Skorzystamy z elementarnej definicji prawdopodobieństwa warunkowego (poniżej piszemy $\pi_\beta(\cdot) = \pi(\cdot)$, bo parametr β jest ustalony):

$$\begin{aligned}
\pi(x_s | x_{-s}) &= \frac{\pi(x)}{\pi(x_{-s})} = \frac{\pi(x)}{\sum_a \pi(x_{a \rightsquigarrow s})} \\
&= \frac{\exp -\beta H(x)}{\sum_a \exp -\beta H(x_{a \rightsquigarrow s})} \\
&= \frac{\exp -\beta \left(\sum_{t: t \rightsquigarrow s} V(x_s, x_t) + \sum_{t \rightsquigarrow w, t \neq s, w \neq s} V(x_t, x_w) + U_s(x_s) + \sum_{t \neq s} U_t(x_t) \right)}{\sum_a \exp -\beta \left(\sum_{t: t \rightsquigarrow s} V(a, x_t) + \sum_{t \rightsquigarrow w, t \neq s, w \neq s} V(x_t, x_w) + U_s(a) + \sum_{t \neq s} U_t(x_t) \right)} \\
&= \frac{\exp -\beta \left(\sum_{t: t \rightsquigarrow s} V(x_s, x_t) + U_s(x_s) \right)}{\sum_a \exp -\beta \left(\sum_{t: t \rightsquigarrow s} V(a, x_t) + U_s(a) \right)} \\
&= \frac{\exp -\beta H_s(x)}{Z_s(\beta)}.
\end{aligned}$$

Ponieważ otrzymany wynik zależy tylko od x_s i $x_{\partial s}$, więc $\pi(x_s | x_{-s}) = \pi(x_s | x_{\partial s})$. Ten wniosek jest pewną formą własności Markowa. \square

Zauważmy, że obliczenie $H_s(x)$ jest łatwe, bo suma $\sum_{t \in \partial s} \dots$ zawiera tylko tyle składników, ile jest sąsiadów miejsca s . Obliczenie $Z_s(\beta)$ też jest łatwe, bo suma $\sum_{a \in \Lambda} \dots$ zawiera tylko $l = |\Lambda|$ składników. Ale nawet nie musimy obliczać stałej normującej $Z_s(\beta)$ żeby generować z rozkładu

$$\pi(x_s = a | x_{\partial s}) \propto \exp -\beta \left(\sum_{t: t \rightsquigarrow s} V(a, x_t) + U_s(a) \right). \quad (13.2)$$

Na tym opiera się implementacja próbnika Gibbsa. Wersję PG z „systemstycznym przeglądem miejsc” można zapisać tak:

Listing.

```

for  $s \in \mathcal{S}$  do
  begin
    Gen  $a \sim \pi(x_s = \cdot | x_{\partial s})$ ;
     $x := x_{a \rightsquigarrow s}$ 
  end

```

Faktycznie już ten algorytm spotkaliśmy w Podrozdziale 12.2, dla szczególnego przypadku modelu auto-logistycznego.

13.3. Rekonstrukcja obrazów

Bayesowski model rekonstrukcji obrazów został zaproponowany w pracy Gemana i Geman w 1987 roku. Potem zdobył dużą popularność i odniósł wiele sukcesów. Model łączy idee zaczerpnięte ze statystyki bayesowskiej i fizyki statystycznej. Cyfrową reprezentację obrazu utożsamiamy z konfiguracją kolorów na wierzchołkach grafu, czyli z elementem przestrzeni $\mathcal{X} = \Lambda^{\mathcal{S}}$, zdefiniowanej w Podrozdziale 13.1. Przyjmijmy, że „idealny obraz”, czyli to co chcielibyśmy zrekonstruować jest konfiguracją $x \in \mathcal{X}$. Niestety, obraz jest „zakłócony” lub „zaszumiony”.

Możemy tylko obserwować konfigurację y reprezentującą zakłócony obraz. Zbiór kolorów w obrazie y nie musi być identyczny jak w obrazie x . Powiedzmy, że $y \in \mathcal{Y} = \Gamma^S$. Ważne jest to, że zniekształcenie modelujemy probabilistycznie przy pomocy rodziny rozkładów warunkowych $f(y|x)$. Dodatkowo zakładamy, że obraz x pojawia się losowo, zgodnie z rozkładem prawdopodobieństwa $\pi(x)$. Innymi słowy, „idealny” obraz x oraz „zniekształcony” obraz y traktujemy jako realizacje zmiennych losowych $X : \Omega \rightarrow \mathcal{X}$ i $Y : \Omega \rightarrow \mathcal{Y}$,

$$\pi(x) = \mathbb{P}(X = x), \quad f(y|x) = \mathbb{P}(Y = y|X = x).$$

W ten sposób buduje się statystyczny model bayesowski, w którym

- Y jest obserwowaną zmienną losową,
- x jest nieznanym parametrem traktowanym jako zmienna losowa X .

Oczywiście, π gra rolę rozkładu *a priori*, zaś f jest wiarogodnością. Być może użycie literki x na oznaczenie parametru jest niezgodne z tradycyjnymi oznaczeniami statystycznymi, ale z drugiej strony jest wygodne. Wzór Bayesa mówi, że rozkład *a posteriori* jest następujący.

$$\pi_y(x) = \mathbb{P}(X = x|Y = y) \propto f(y|x)\pi(x).$$

Pomysł Gemana i Gemana polegał na tym, żeby modelować rozkład *a priori* π jako MPL. Załóżmy, że π jest rozkładem Gibbsa,

$$\pi(x) \propto \exp(-H(x)), \quad (13.3)$$

gdzie

$$H(x) = J \sum_{s \sim t} V(x_s, x_t). \quad (13.4)$$

Energia „*a priori*” zawiera tu tylko składniki reprezentujące oddziaływania między parami miejsc sąsiednich. Funkcja $V(a, b)$ zazwyczaj ma najmniejszą wartość dla $a = b$ i rośnie wraz z „odległością” między a i b (jakkolwiek tę odległość zdefiniujemy). W ten sposób „nagradza” konfiguracje w których sąsiednie miejsca są podobnie pokolorowane. Im większy parametr $J > 0$, tym bardziej prawdopodobne są obrazy zawierające jednolite plamy kolorów.

Trzeba jeszcze założyć coś o „wiarogodności” f . Dla uproszczenia opiszę tylko najprostszy model, w którym kolor y_s na obserwowanym obrazie zależy tylko od koloru x_s na obrazie idealnym. Intuicyjnie znaczy to, że „zaszumienie” ma ściśle lokalny charakter. Matematycznie znaczy to, że

$$f(y|x) = \prod_s f(y_s|x_s)$$

(pozwolę sobie na odrobinę nieścisłości aby uniknąć nowego symbolu na oznaczenie $f(y_s|x_s)$). Zapiszemy teraz „wiarogodność” f w postaci zlogarytmowanej. Jeśli położymy $-\log f(y_s|x_s) = U_s(x_s)$ pamiętając, że y jest w świecie bayesowskim ustalone, to otrzymujemy następujący wzór:

$$\pi_y(x) \propto \exp(-H_y(x)), \quad (13.5)$$

gdzie

$$H_y(x) = J \sum_{s \sim t} V(x_s, x_t) - \sum_s \log f(y_s|x_s). \quad (13.6)$$

Okazuje się zatem, że rozkład *a posteriori* ma podobną postać do rozkładu *a priori*. Też jest rozkładem Gibbsa, a różnica polega tylko na dodaniu składników reprezentujących oddziaływania zewnętrzne $U_s(x_s) = -\log f(y_s|x_s)$. Pamiętajmy przy tym, że y jest w świecie bayesowskim ustalone. W modelu rekonstrukcji obrazów „oddziaływania zewnętrzne” zależą od y i „wymuszają podobieństwo” rekonstruowanego obrazu do obserwacji. Z kolei „oddziaływania między parami” są odpowiedzialne za wygładzenie obrazu. Lepiej to wyjaśnimy na przykładzie.

Przykład 13.2 (Losowe „przekłamanie koloru” i wygładzanie Potts’a). Załóżmy, że $\Lambda = \{1, \dots, l\}$ jest naprawdę paletą kolorów, na przykład

$$\Lambda = \{\text{Czerwony}, \text{Niebieski}, \text{Pomarańczowy}, \text{Zielony}\}.$$

Przypuśćmy, że mechanizm losowego „przekłamania” polega na tym, że w każdym pikslu, kolor obecny w idealnym obrazie x jest z prawdopodobieństwem $1 - \varepsilon$ niezmieniony, a z prawdopodobieństwem ε zmienia się na losowo wybrany inny kolor. Tak więc zarówno x jak i y należą do tej samej przestrzeni Λ^S ,

$$f(y_s|x_s) = \begin{cases} 1 - \varepsilon & \text{dla } y_s = x_s; \\ \varepsilon/(l-1) & \text{dla } y_s \neq x_s. \end{cases}$$

Można za rozkład *a priori* przyjąć rozkład Potts’a z Przykładu 13.1. Rozkład *a posteriori* ma funkcję energii daną następującym wzorem (z $J > 0$):

$$H_y(x) = J \sum_{s \sim t} \mathbb{I}(x_s \neq x_t) - \sum_s [\log(1 - \varepsilon) \mathbb{I}(x_s = y_s) + \log(\varepsilon/(l-1)) \mathbb{I}(x_s \neq y_s)].$$

Pierwszy składnik w tym wzorze pochodzi od rozkładu *a priori* (z modelu Potts’a) i „nagradza” konfiguracje w których dużo sąsiednich punktów jest pomalowanych na ten sam kolor. Powoduje to, że obrazy x składające się z jednolitych dużych „plam” są preferowane. Drugi składnik pochodzi od obserwowanej konfiguracji y i jest najmniejszy dla $x = y$. Powoduje to, że obrazy x mało się różniące od y są bardziej prawdopodobne. Rozkład *a posteriori* jest pewnym kompromisem pomiędzy tymi dwoma konkurującymi składnikami. Parametr J jest „wagą” pierwszego składnika i dlatego odgrywa rolę „parametru wygładzającego”. Im większe J tym odtwarzany obraz będzie bardziej regularny (a tym mniej będzie starał się upodobnić do y). I odwrotnie, małe J powoduje ściślejsze dopasowanie x do y ale mniejszą „regularność” x .

Jeszcze lepiej to samo widać na przykładzie tak zwanego „szumu gaussowskiego”.

Przykład 13.3 (Addytywny szum gaussowski). Załóżmy, że x jest konfiguracją „poziomów szarości” czyli, powiedzmy, $\Lambda \subseteq [0, \infty[$. Mechanizm losowego „zaszumienia” polega na tym, że zamiast poziomu szarości x_s obserwujemy $y_s \sim N(x_s, \sigma^2)$. Innymi słowy,

$$f(y_s|x_s) \propto \exp \left[-\frac{1}{2\sigma^2} (y_s - x_s)^2 \right].$$

Przestrzenią obserwowanych konfiguracji y jest tutaj (formalnie) \mathbb{R}^S (faktycznie, raczej $[0, \infty[^S$). Rozkład *a posteriori* ma funkcję energii daną następującym wzorem:

$$H_y(x) = J \sum_{s \sim t} V(x_s \neq x_t) + \frac{1}{2\sigma^2} \sum_s (y_s - x_s)^2.$$

Jeśli rozpatrujemy model ze skończoną liczbą poziomów szarości dla konfiguracji x to można pierwszy składnik określić tak jak w poprzednim przykładzie, czyli zapożyczyć z modelu Potts’a. Bardziej naturalne jest określenie $V(a, b)$ w taki sposób, aby większe różnice pomiędzy poziomami a i b były silniej karane. parametr J jest, jak poprzednio, odpowiedzialny za stopień wygładzenia.

14. Markowskie Monte Carlo V. Elementarna teoria łańcuchów Markowa

14.1. Podstawowe określenia i oznaczenia

W tym rozdziale rozważamy jednorodny łańcuch Markowa $X_0, X_1, \dots, X_n, \dots$ na skończonej przestrzeni stanów $\mathcal{X} = \{1, \dots, d\}$. Będziemy posługiwać się wygodną i zwięzłą notacją wektorowo-macierzową. Macierz przejścia o wymiarach $d \times d$ oznaczamy $P = (P(x, y))_{x, y \in \mathcal{X}}$. Rozkład początkowy utożsamiamy z wektorem wierszowym $\xi^\top = (\xi(1), \dots, \xi(x), \dots, \xi(d))$. W dalszym ciągu, mówiąc o łańcuchu Markowa, będziemy mieli na myśli ustaloną macierz przejścia P i dowolnie wybrany rozkład początkowy ξ . Przyjmujemy oznaczenie $\mathbb{P}_\xi(\cdot)$. W szczególności, $\mathbb{P}_x(\cdot) = \mathbb{P}(\cdot | X_0 = x)$, dla $x \in \mathcal{X}$. Analogicznie będziemy oznaczali wartość oczekiwaną: \mathbb{E}_ξ lub \mathbb{E}_x . Zauważmy, że $\mathbb{P}(X_{n+2} = y | X_n = x) = \sum_z P(x, z)P(z, y) = P^2(x, y)$. Ogólniej, macierz przejścia w m krokach jest m -tą potęgą macierzy P :

$$\mathbb{P}(X_{n+m} = y | X_n = x) = P^m(x, y).$$

Rozkład *brzegowy* zmiennej losowej X_n jest wektorem $\xi^\top P^n$:

$$\mathbb{P}(X_n = x) = (\xi^\top P^n)(x).$$

Z macierzą stochastyczną P związany jest graf skierowany opisujący możliwe przejścia łańcucha. Jest to graf $(\mathcal{X}, \mathcal{E})$, gdzie zbiorem wierzchołków jest przestrzeń stanów, zaś $\mathcal{E} \subseteq \mathcal{X} \times \mathcal{X}$ jest zbiorem takich par (x, y) , że $P(x, y) > 0$. Pojęcia, o których teraz będziemy mówić, są związane tylko ze strukturą grafu, czyli z położeniem niezerowych elementów macierzy przejścia.

Mówimy, że stan y jest *osiągalny* ze stanu x , jeśli $x = y$ lub istnieje liczba naturalna $n \geq 1$ i ciąg stanów $x = x_0, x_1, x_2, \dots, x_{n-1}, x_n = y$ taki, że $P(x_{i-1}, x_i) > 0$ dla $i = 1, \dots, n$. Równoważnie, $P^n(x, y) > 0$ dla pewnego $n \geq 0$. Będziemy stosowali oznaczenie „ $x \rightarrow y$ ”.

Stany x i y *komunikują* się, co zapiszemy „ $x \leftrightarrow y$ ”, jeśli $x \rightarrow y$ i $y \rightarrow x$.

Stan x jest *istotny*, jeśli dla każdego y takiego, że $x \rightarrow y$ mamy $y \rightarrow x$. W przeciwnym razie stan x nazywamy *nieistotnym*. Stan x jest więc istotny, jeśli istnieje stan y do którego można przejść z x , ale nie można wrócić do x . Zbiór stanów istotnych oznaczmy przez \mathcal{C} , zaś zbiór stanów nieistotnych przez \mathcal{T} . Zbiór \mathcal{C} jest (dla łańcucha o skończonej przestrzeni stanów) zawsze niepusty. Zbiór \mathcal{T} może być pusty.

Jeśli $x \leftrightarrow y$ dla dowolnych $x, y \in \mathcal{X}$, czyli wszystkie stany komunikują się, to łańcuch nazywamy *nieprzywiedlnym*. Oczywiście, wszystkie stany łańcucha nieprzywiedlnego są istotne, $\mathcal{X} = \mathcal{C}$. Łańcuch nieprzywiedlny nazywamy *nieokresowym*, jeśli dla dowolnych $x, y \in \mathcal{X}$ istnieje n_0 takie, że dla każdego $n \geq n_0$ mamy $P^n(x, y) > 0$. Wspomnijmy, że przyjęta przez nas definicja nieokresowości bywa formułowana w nieco inny, ale równoważny sposób. Dla naszych potrzeb wystarczy następujące proste spostrzeżenie: jeśli łańcuch jest nieprzywiedlny i $P(x, x) > 0$ dla przynajmniej jednego stanu x , to łańcuch jest nieokresowy. Latwo zauważyć, że dla łańcucha nieprzywiedlnego i nieokresowego, dla dostatecznie dużych n macierz P^n ma wszystkie elementy niezerowe. Wynika to z faktu, że rozważamy łańcuchy ze skończoną przestrzenią stanów.

Jeśli $x \leftrightarrow y$ dla dowolnych $x, y \in \mathcal{C}$, czyli wszystkie stany *istotne* komunikują się, to łańcuch nazywamy *jednoklasowym*. Łańcuch jednoklasowy może mieć niepusty zbiór stanów nieistotnych \mathcal{T} . Łańcuch jednoklasowy możemy „przerobić” na nieprzywiedlny jeśli ograniczymy przestrzeń do stanów istotnych. Macierz $P|_{\mathcal{C}} = (P(x, y)_{x, y \in \mathcal{C}})$ jest, jak łatwo widzieć, stochastyczna.

Interesują nas głównie łańcuchy, które „zmierzają w kierunku położenia równowagi”. Aby uściślić co to znaczy „równowaga”, przypomnijmy pojęcie stacjonarności. Rozkład π jest stacjonarny jeśli dla każdego stanu y ,

$$\pi(y) = \sum_x \pi(x)P(x, y).$$

W notacji macierzowej: $\pi^\top = \pi^\top P$. Stąd oczywiście wynika, że $\pi^\top = \pi^\top P^n$.

Poniższy prosty fakt można uzasadnić na wiele sposobów. W następnym podrozdziale przytoczymy, wraz z dowodem, piękne twierdzenie Kaca (Twierdzenie 14.2), które implikuje Twierdzenie 14.1.

Twierdzenie 14.1. *Jeśli łańcuch Markowa jest nieprzywiedlny, to istnieje dokładnie jeden rozkład stacjonarny π , przy tym $\pi(x) > 0$ dla każdego $x \in \mathcal{X}$.*

Z Twierdzenia 14.1 łatwo wynika następujący wniosek.

Wniosek 14.1. *Jeśli łańcuch Markowa jest jednoklasowy, to istnieje dokładnie jeden rozkład stacjonarny π , przy tym $\pi(x) > 0$ dla każdego $x \in \mathcal{C}$, czyli dla wszystkich stanów istotnych oraz $\pi(x) = 0$ dla każdego $x \in \mathcal{T}$, czyli dla stanów nieistotnych.*

Uwaga 14.1. Podkreślmy stale obowiązujące w tym rozdziale założenie, że *przestrzeń stanów jest skończona*. To założenie jest istotne w Twierdzeniu 14.1 i to samo dotyczy dalszych rozważań. Istnieją co prawda odpowiedniki sformułowanych tu twierdzeń dla przypadku ogólnej przestrzeni stanów (nieskończonej, a nawet „ciągłej” takiej jak \mathbb{R}^d) ale wymagają one dodatkowych, niełatwych do sprawdzenia założeń. Przystępny i bardzo elegancki wykład teorii łańcuchów Markowa na ogólnej przestrzeni stanów można znaleźć w pracy Nummelina [17]. Przeglądowy artykuł Robertsa i Rosenthala [20] zawiera dużo dodatkowych informacji na ten temat. Obie cytowane prace koncentrują się na tych własnościach łańcuchów, które są istotne z punktu widzenia algorytmów Monte Carlo. Z kolei piękna książka Brémaud [4] ogranicza się do przestrzeni dyskretnych (skończonych lub przeliczalnych).

14.2. Regeneracja

Przedstawimy w tym podrozdziale konstrukcję, która prowadzi do łatwych i eleganckich dowodów twierdzeń granicznych. Podstawowa idea jest następująca. Wyróżnia się jeden ustalony stan, powiedzmy $z \in \mathcal{X}$. W każdym momencie wpadnięcia w z , następuje „odnowienie” i dalsza ewolucja łańcucha jest niezależna od przeszłości.

Niech, dla $z \in \mathcal{X}$,

$$T^z = \min\{n > 0 : X_n = z\}. \quad (14.1)$$

Przyjmujemy przy tym naturalną konwencję: $T^z = \infty$, jeśli $X_n \neq z$ dla każdego $n \geq 1$. Zmienna losowa T^z jest więc czasem pierwszego dojścia do stanu z . Jeśli założymy, że łańcuch startuje z punktu z , to T^z jest czasem pierwszego powrotu.

Lemat 14.1. *Jeżeli łańcuch jest nieprzywiedlny, to istnieją stałe c i $\gamma < 1$ takie, że dla dowolnego rozkładu początkowego ξ ,*

$$\mathbb{P}_\xi(T^z > n) \leq c\gamma^n.$$

Dowód. Dla uproszczenia przyjmijmy dodatkowe założenie, że łańcuch jest nieokresowy. Wtedy dla dostatecznie dużych k wszystkie elementy macierzy P^k są niezerowe. Ustalmy k i znajdziemy liczbę $\delta > 0$ taką, że $P^k(x, z) \geq \delta$ dla wszystkich x (jest to możliwe, bo łańcuch ma skończoną liczbę stanów). Dla dowolnego n , dobierzmy takie m , że $mk \leq n < (m+1)k$. Mamy wówczas

$$\begin{aligned} \mathbb{P}_\xi(T^y > n) &\leq \mathbb{P}_\xi(T^y > mk) \\ &\leq \mathbb{P}_\xi(X_0 \neq z, X_k \neq z, \dots, X_{mk} \neq z) \\ &= \sum_{x_0 \neq z, x_1 \neq z, \dots, x_m \neq z} \xi(x_0) P^k(x_0, x_1) \cdots P^k(x_{m-1}, x_m) \\ &\leq (1 - \delta)^m \leq c\gamma^n, \end{aligned}$$

dla $\gamma = (1 - \delta)^{1/k}$ i $c = (1 - \delta)^{-1}$.

W przypadku łańcucha okresowego dowód nieco się komplikuje i, choć nie jest trudny, zostanie pominięty. \square

Wniosek 14.2. *Dla łańcucha nieprzywiedlnego mamy $\mathbb{P}_\xi(T^z < \infty) = 1$, co więcej $\mathbb{E}_\xi T^z < \infty$, co więcej $\mathbb{E}_\xi (T^z)^k < \infty$ dla dowolnego k , a nawet $\mathbb{E}_\xi \exp(\lambda T^z) < \infty$ przynajmniej dla pewnych dostatecznie małych wartości $\lambda > 0$ (w istocie dla $\lambda < -\log \gamma$).*

Podamy teraz bardzo ciekawą interpretację rozkładu stacjonarnego, wykazując przy okazji jego istnienie (Twierdzenie 14.1). Ustalmy dowolnie wybrany stan z . Udowodnimy, że średni czas, spędzony przez łańcuch w stanie y pomiędzy wyjściem z z i pierwszym powrotem do z jest proporcjonalny do $\pi(y)$, prawdopodobieństwa stacjonarnego.

Twierdzenie 14.2 (Kaca). *Założmy, że łańcuch jest nieprzywiedlny. Ustalmy $z \in \mathcal{X}$ i zdefiniujmy miarę α wzorem*

$$\alpha(y) = \mathbb{E}_z \sum_{i=0}^{T^z-1} \mathbb{I}(X_i = y) = \mathbb{E}_z \sum_{i=1}^{T^z} \mathbb{I}(X_i = y).$$

Wtedy:

- (i) *Miara α jest stacjonarna, czyli $\alpha^\top P = \alpha$.*
- (ii) *Miara α jest skończona, $\alpha(\mathcal{X}) = \mathbb{E}_z(T^z) = m < \infty$.*
- (iii) *Unormowana miara $\alpha/m = \pi$ jest jedynym rozkładem stacjonarnym.*

Dowód. Dla uproszczenia będziemy pisali $T^z = T$, $\mathbb{P}_z = \mathbb{P}$ i $\mathbb{E}_z = \mathbb{E}$. Zauważmy, że

$$\begin{aligned} \alpha(x) &= \mathbb{E} \sum_{i=0}^{T-1} \mathbb{I}(X_i = x) = \mathbb{E} \sum_{i=0}^{\infty} \mathbb{I}(X_i = x, T > i) \\ &= \sum_{i=0}^{\infty} \mathbb{P}(X_i = x, T > i). \end{aligned}$$

Udowodnimy teraz (i). Jeśli $y \neq z$, to

$$\begin{aligned} \sum_x \alpha(x)P(x, y) &= \sum_x \sum_{i=0}^{\infty} \mathbb{P}(X_i = x, T > i)P(x, y) \\ &= \sum_{i=0}^{\infty} \sum_x \mathbb{P}(X_i = x, T > i)P(x, y) \\ &= \sum_{i=0}^{\infty} \mathbb{P}(X_{i+1} = y, T > i+1) = \sum_{i=1}^{\infty} \mathbb{P}(X_i = y, T > i) \\ &= \alpha(y), \end{aligned}$$

ponieważ $\mathbb{P}(X_0 = y) = 0$, bo $\mathbb{P}(X_0 = z) = 1$. Dla $y = z$ mamy z kolei

$$\begin{aligned} \sum_x \alpha(x)P(x, z) &= \sum_x \sum_{i=0}^{\infty} \mathbb{P}(X_i = x, T > i)P(x, z) \\ &= \sum_{i=0}^{\infty} \sum_x \mathbb{P}(X_i = x, T > i)P(x, z) \\ &= \sum_{i=0}^{\infty} \mathbb{P}(X_{i+1} = z, T = i+1) = \sum_{i=1}^{\infty} \mathbb{P}(T = i) \\ &= 1 = \alpha(z), \end{aligned}$$

co kończy dowód (i).

Część (ii) jest łatwa. Równość

$$\alpha(\mathcal{X}) = \sum_y \alpha(y) = \mathbb{E}T$$

wynika wprost z definicji miary α . Fakt, że $m = \mathbb{E}T < \infty$ jest wnioskiem z Lematu 14.1.

Punkt (iii): istnienie rozkładu stacjonarnego

$$\pi(y) = \frac{\alpha(y)}{m}.$$

jest natychmiastowym wnioskiem z (i) i (ii). Jednoznaczność rozkładu stacjonarnego dla jest nietrudna do bezpośredniego udowodnienia. Pozostawiamy to jako ćwiczenie. W najbardziej interesującym nas przypadku łańcucha nieokresowego, jednoznaczność wyniknie też ze Słabego Twierdzenia Ergodycznego, które udowodnimy w następnym podrozdziale. \square

Odnotujmy ważny wniosek wynikający z powyższego twierdzenia:

$$\pi(z) = \frac{1}{\mathbb{E}_z(T^z)}.$$

Zjawisko odnowienia, czyli regeneracji pozwala sprowadzić badanie łańcuchów Markowa do rozpatrywania niezależnych zmiennych losowych, a więc do bardzo prostej i dobrze znanej sytuacji. Aby wyjaśnić to bliżej, zauważmy następującą oczywistą równość. Na mocy własności Markowa i jednorodności,

$$\begin{aligned} &\mathbb{P}(X_{n+1} = x_1, \dots, X_{n+k} = x_k | T^z = n) \\ &\mathbb{P}(X_{n+1} = x_1, \dots, X_{n+k} = x_k | X_n = z) \\ &= \mathbb{P}_z(X_1 = x_1, \dots, X_k = x_k). \end{aligned}$$

Zatem warunkowo, dla $T^z = n$, łańcuch „regeneruje się w momencie n ” i zaczyna się zachowywać dokładnie tak, jak łańcuch który wystartował z punktu z w chwili 0. Niezależnie od przeszłości! Ponieważ z jest ustalone, będziemy odtąd pomijali górny indeks przy $T = T^z$. Zdefiniujemy kolejne momenty odnowienia, czyli czasy odwiedzin stanu z :

$$T = T_1 = \min\{n > 0 : X_n = z\},$$

$$T_k = \min\{n > T_{k-1} : X_n = z\}.$$

Momenty $0 < T_1 < \dots < T_k < \dots$ dzielą trajektorię łańcucha na następujące „losowe wycieczki”, czyli losowej długości ciągi zmiennych losowych:

$$\underbrace{X_0, \dots, X_{T_1-1}}_{T_1}, \quad \underbrace{X_{T_1}, \dots, X_{T_2-1}}_{T_2-T_1}, \quad \underbrace{X_{T_2}, \dots, X_{T_3-1}}_{T_3-T_2}, \quad \dots$$

$$\quad \quad \quad \uparrow \quad \quad \quad \uparrow \quad \quad \quad \dots$$

$$\quad \quad \quad X_{T_1} = z \quad \quad \quad X_{T_2} = z \quad \quad \quad \dots$$

Wycieczka zaczyna się w punkcie z i kończy tuż przed powrotem do z . Oznaczmy k -tą wycieczkę symbolem Ξ_k :

$$\Xi = \Xi_1 = (X_0, \dots, X_{T-1}, T),$$

$$\Xi_k = (X_{T_{k-1}}, \dots, X_{T_k-1}, T_k - T_{k-1})$$

Z tego, co powiedzieliśmy wcześniej wynika, że wszystkie „wycieczki” są niezależne. Co więcej wycieczki Ξ_k mają ten sam rozkład, z wyjątkiem być może początkowej, czyli Ξ_1 . Jeśli rozkład początkowy jest skupiony w punkcie z , to również wycieczka Ξ_1 ma ten sam rozkład (0 jest wtedy momentem odnowienia).

Podejście regeneracyjne, czyli rozbięcie łańcucha na niezależne wycieczki prowadzi do ładnych i łatwych dowodów PWL i CTG dla łańcuchów Markowa. Sformułujemy najpierw pewną wersję *Mocnego Prawa Wielkich Liczb*. Rozważmy funkcję f o wartościach rzeczywistych, określoną na przestrzeni stanów. Przypomnijmy, że $\mathbb{E}_\pi f = \sum_{x \in \mathcal{X}} \pi(x) f(x)$.

Twierdzenie 14.3 (Mocne Twierdzenie Ergodyczne). *Jeśli X_n jest nieprzywiedlnym łańcuchem Markowa, to dla dowolnego rozkładu początkowego ξ i każdej funkcji $f : \mathcal{X} \rightarrow \mathbb{R}$,*

$$\frac{1}{n} \sum_{i=0}^{n-1} f(X_i) \longrightarrow \mathbb{E}_\pi f \quad (n \rightarrow \infty)$$

z prawdopodobieństwem 1.

Dowód. Zdefiniujmy sumy blokowe:

$$\Xi_0(f) = \sum_{i=0}^{T-1} f(X_i),$$

$$\Xi_k(f) = \sum_{i=T_k}^{T_{k+1}-1} f(X_i).$$

Niech $N(n) = \max\{k : T_k \leq n\}$, czyli $T_{N(n)}$ jest ostatnią regeneracją przed momentem n :

$$0, \dots, T_1 - 1, \quad T_1, \dots, T_{N(n)}, \dots, n, \dots, T_{N(n)+1} - 1, \quad T_{N(n)+1}, \dots$$

$$\quad \quad \quad \uparrow \quad \quad \quad \uparrow \quad \quad \quad \uparrow \quad \quad \quad \uparrow$$

$$\quad \quad \quad X = z \quad \quad \quad X = z \quad \quad \quad \bullet \quad \quad \quad X = z$$

Oczywiście,

$$T_{N(n)} \leq n < T_{N(n)+1}. \quad (14.2)$$

Wiemy, że $\mathbb{E}_z T = m < \infty$. Ponieważ T_k jest sumą k niezależnych zmiennych losowych (długości wycieczek) to wnioskujemy, że $T_k/k \rightarrow m$ z prawdopodobieństwem 1, na mocy *zwykłego Prawa Wielkich Liczb*. Rzecz jasna, tak samo $T_{k+1}/k \rightarrow m$. Podzielmy nierówność (14.2) stronami przez $N(n)$ i przejdźmy do granicy (korzystając z tego, że $N(n) \rightarrow \infty$ prawie na pewno). Twierdzenie o trzech ciągach pozwala wywnioskować, że

$$\frac{N(n)}{n} \rightarrow \frac{1}{m} \text{ p.n.}$$

Założmy teraz, że $f \geq 0$ i powtórzmy bardzo podobne rozumowanie dla sum

$$\sum_{j=1}^{N(n)} \Xi_j(f) \leq S_n(f) = \sum_{i=0}^{n-1} f(X_i) \leq \sum_{j=1}^{N(n)+1} \Xi_j(f). \quad (14.3)$$

Po lewej i po prawej stronie mamy sumy niezależnych składników $\Xi_j(f)$. Korzystamy z PWL dla niezależnych zmiennych, dzielimy (14.3) stronami przez $N(n)$ i przechodzimy do granicy. Otrzymujemy

$$\frac{S_n(f)}{N(n)} \rightarrow \mathbb{E}_z \Xi(f) \text{ p.n.}$$

a więc

$$\frac{S_n(f)}{n} \rightarrow \frac{\mathbb{E}_z \Xi(f)}{m} = \frac{1}{m} \sum_x \alpha(x) f(x) = \sum_x \pi(x) f(x) \text{ p.n.}$$

Ostatnia równość wynika z Twierdzenia Kaca. Przypomnijmy, że $\alpha(x)$ jest „średnim czasem spędzonym w stanie x ” podczas pojedynczej wycieczki.

Jeśli funkcja f nie jest nieujemna, to możemy zastosować rozkład $f = f^+ - f^-$ i wykorzystać już udowodniony wynik. \square

Na podobnej idei oparty jest „regeneracyjny” dowód Centralnego Twierdzenia Granicznego (istnieją też zupełnie inne dowody).

Twierdzenie 14.4 (Centralne Twierdzenie Graniczne). *Jeśli X_n jest łańcuchem nieprzywiedlnym, to dla dowolnego rozkładu początkowego ξ i każdej funkcji $f : \mathcal{X} \rightarrow \mathbb{R}$,*

$$\frac{1}{\sqrt{n}} \left(\sum_{i=0}^{n-1} [f(X_i) - \mathbb{E}_\pi f] \right) \rightarrow_d N(0, \sigma_{\text{as}}^2(f)), \quad (n \rightarrow \infty).$$

Ponadto, dla dowolnego rozkładu początkowego ξ zachodzi wzór (10.4), czyli $(1/n)\text{Var}_\xi \sum_{i=0}^{n-1} f(X_i) \rightarrow \sigma_{\text{as}}^2(f)$ przy $n \rightarrow \infty$.

Szkic dowodu. Trochę więcej jest tu technicznych zawiłości niż w dowodzie PWL, wobec tego zdecydowałem się pominąć szczegóły. W istocie, przedstawię tylko bardzo pobieżnie główną ideę. Bez straty ogólności założmy, że $\pi^\top f = 0$. Tak jak w dowodzie PWL, sumę $S_n(f) = \sum_{i=0}^{n-1} f(X_i)$ przybliżamy sumą niezależnych składników, które odpowiadają całkowitym wycieczkom: $S_n(f) \simeq S_{T_{N(n)}}(f) = \sum_{j=1}^{N(n)} \Xi_j(f)$. Ze zwykłego CTG dla niezależnych zmiennych o jednakowym rozkładzie otrzymujemy

$$\frac{1}{\sqrt{k}} \sum_{j=1}^k \Xi_j(f) \rightarrow_d N(0, \text{Var}_z \Xi(f)).$$

Jeśli „podstawimy” w miejsce k zmienną losową $N(n)$ i wykorzystamy fakt, że $N(n) \simeq n/m$ (PWL gwarantuje, że $N(n)/n \rightarrow 1/m$), to nie powinien dziwić następujący wniosek:

$$\frac{1}{\sqrt{n}} S_n(f) \rightarrow_d N(0, \text{Var}_z \Xi(f)/m).$$

W ten sposób „otrzymujemy” tezę. □

Przytoczony powyżej rozumowanie daje ciekawe przedstawienie asymptotycznej wariancji:

$$\sigma_{\text{as}}^2(f) = \text{Var}_z \Xi(f) / \mathbb{E}_z T. \quad (14.4)$$

Istnieje wiele wyrażeń na asymptotyczną wariancję, przy tym różne wzory wymagają różnych założeń. Najbardziej znany jest wzór (10.5). Sformułujemy w jawny sposób potrzebne założenia, i przepisemy ten wzór w postaci macierzowej. Najpierw musimy wprowadzić jeszcze parę nowych oznaczeń. Wygodnie będzie utożsamić funkcję f z wektorem kolumnowym

$$f = \begin{pmatrix} f(1) \\ \vdots \\ f(d) \end{pmatrix}.$$

Niech

$$\Pi = \text{diag}(\pi(1), \dots, \pi(d)),$$

gdzie π oznacza, jak zwykle, rozkład stacjonarny. Możemy teraz napisać

$$\mathbb{E}_\pi f(X_0) = \sum_x \pi(x) f(x) = \pi^\top f$$

oraz

$$\begin{aligned} \text{Var}_\pi f(X_0) &= \mathbb{E}_\pi f^2(X_0) - [\mathbb{E}_\pi f(X_0)]^2 \\ &= f^\top \Pi f - f^\top \pi \pi^\top f \\ &= f^\top \Pi (I - 1\pi^\top) f. \end{aligned}$$

Podobnie,

$$\begin{aligned} \text{Cov}_\pi [f(X_0), f(X_n)] &= \mathbb{E}_\pi [f(X_0) f(X_n)] - [\mathbb{E}_\pi f(X_0)]^2 \\ &= f^\top \Pi P^n f - f^\top \pi \pi^\top f \\ &= f^\top \Pi (P^n - 1\pi^\top) f. \end{aligned}$$

W powyższym wzorze, I jest macierzą identycznościową, 1 oznacza kolumnę jedynek. Łatwo sprawdzić, że $P^n - 1\pi^\top = (P - 1\pi^\top)^n$ dla $n > 0$ (ale nie dla $n = 0$). Jeśli P jest macierzą przejścia nieprzywiedlnego i nieokresowego łańcucha Markowa, to $P^n - 1\pi^\top \rightarrow 0$ przy $n \rightarrow \infty$ na mocy Słabego Twierdzenia Ergodycznego. Stąd wynika zbieżność następujących szeregów:

$$\begin{aligned} A &= \sum_{n=0}^{\infty} (P^n - 1\pi^\top), \\ Z &= \sum_{n=0}^{\infty} (P - 1\pi^\top)^n = (I - P + 1\pi^\top)^{-1}. \end{aligned}$$

Macierz Z nazywamy *macierzą fundamentalną*. Ponieważ $P^0 - 1\pi^\top = I - 1\pi^\top$ zaś $(P - 1\pi^\top)^0 = I$ więc $Z = A + 1\pi^\top$.

Stwierdzenie 14.1 (Asymptotyczna wariancja). *Jeśli łańcuch jest nieprzystawny i nieokresowy, to zachodzi wzór (10.5), czyli*

$$\sigma_{\text{as}}^2(f) = \sigma^2(f) \sum_{n=-\infty}^{\infty} \rho_n(f),$$

gdzie $\sigma^2(f) = \text{Var}_{\pi} f(X_0)$ i $\rho_n(f) = \text{corr}_{\pi}[f(X_0), f(X_n)]$, W postaci macierzowej asymptotyczna wariancja wyraża się następująco

$$\sigma_{\text{as}}^2(f) = f^{\top} (2\Pi Z - \pi\pi^{\top} - \Pi) f = f^{\top} (2\Pi A + \pi\pi^{\top} - \Pi) f.$$

Szkic dowodu. Załóżmy, że rozkładem początkowym jest π i skorzystamy ze stacjonarności łańcucha:

$$\begin{aligned} \frac{1}{n} \text{Var}_{\pi} \left(\sum_{i=0}^{n-1} f(X_i) \right) &= \frac{1}{n} \sum_{i=0}^{n-1} \text{Var}_{\pi}(X_i) + \frac{2}{n} \sum_{i=0}^{n-1} \sum_{j=i}^{n-1} \text{Cov}_{\pi}(f(X_i), f(X_j)) \\ &= \text{Var}_{\pi} f(X_0) + 2 \sum_{k=1}^{n-1} \frac{n-k}{n} \text{Cov}_{\pi}(f(X_0), f(X_k)) \\ &\rightarrow \text{Var}_{\pi} f(X_0) + 2 \sum_{k=1}^{\infty} \text{Cov}_{\pi}(f(X_0), f(X_k)), \quad (n \rightarrow \infty). \end{aligned}$$

Poprawność przejścia do granicy w ostatniej linijce wynika z elementarnego faktu, że dla dowolnego ciągu liczbowego a_n mamy $\lim_{n \rightarrow \infty} \sum_{k=1}^{n-1} \frac{n-k}{n} a_k = \sum_{k=1}^{\infty} a_k$, o ile szereg po prawej stronie równości jest zbieżny. To zaś, przy założeniu nieprzystawności i nieokresowości, wynika ze STE (skorzystaliśmy już z tego faktu uzasadniając poprawność definicji Z i A).

Pokazaliśmy w ten sposób, że $(1/n) \text{Var}_{\pi} \sum_{i=0}^{n-1} f(X_i)$ zmierza do granicy, która jest równa prawej stronie wzoru (10.5). Pominiemy uzasadnienie, że można zastąpić rozkład stacjonarny π przez dowolny rozkład początkowy ξ oraz że wzór (10.5) definiuje tę samą wielkość co (14.4).

Macierzowe wyrażenia $\sigma_{\text{as}}^2(f)$ wynikają ze wzorów na kowariancje oraz z określenia macierzy A i Z . \square

Zauważmy, że ani CTG ani PWL nie wymagały założenia o nieokresowości ale w 14.1 to założenie jest potrzebne.

Jak widać, asymptotyczna wariancja wyraża się w postaci formy kwadratowej $\sigma_{\text{as}}^2(f) = f^{\top} C f$ o współczynnikach

$$\begin{aligned} C(x, y) &= \pi(x)A(x, y) + A(y, x)\pi(y) + \pi(x)\pi(y) - \pi(x)\mathbb{I}(x = y), \\ &= \pi(x)Z(x, y) + Z(y, x)\pi(y) - \pi(x)\pi(y) - \pi(x)\mathbb{I}(x = y). \end{aligned}$$

Okazuje się, że „asymptotyczne obciążenie” estymatora $\hat{\theta}_n = (1/n) \sum_{i=0}^{n-1} f(X_i)$ wartości oczekiwanej $\theta = \mathbb{E}_{\pi} f$ można napisać w postaci formy dwuliniowej $\xi^{\top} A f$, gdzie ξ jest rozkładem początkowym. Podsumowując,

$$\begin{aligned} \text{Var}_{\xi}(\hat{\theta}_n) &= \frac{1}{n} f^{\top} C f + o\left(\frac{1}{n}\right), \\ \mathbb{E}_{\xi} \hat{\theta}_n - \theta f &= \frac{1}{n} \xi^{\top} A f + o\left(\frac{1}{n}\right). \end{aligned}$$

Uzasadnienie drugiej części powyższego wzoru pozostawiam jako ćwiczenie. Stąd z łatwością otrzymujemy ważny wzór (10.6) z Rozdziału 10 (wyrażenie na błąd średniokwadratowy).

14.3. Łańcuchy sprzężone i zbieżność rozkładów

Dowód Słabego Twierdzenia Ergodycznego, który przedstawimy, opiera się na tak zwanym sprzęganiu (ang. *coupling*), czyli metodzie „dwóch cząstek”. Ta metoda, jak się okaże, nie tylko pozwala udowodnić zbieżność rozkładów prawdopodobieństwa, ale daje w wielu przypadkach bardzo dobre oszacowania *szybkości zbieżności*.

14.3.1. Odległość pełnego wahan

Najpierw zajmiemy się określeniem odległości między rozkładami. Dla naszych celów najbardziej przydatna będzie następująca metryka. Niech ν i λ będą dwoma rozkładami prawdopodobieństwa na skończonej przestrzeni \mathcal{X} . Odległość *pełnego wahan* pomiędzy ν i λ określamy wzorem

$$\|\nu - \lambda\|_{\text{tv}} = \max_{\mathcal{A} \subseteq \mathcal{X}} |\nu(\mathcal{A}) - \lambda(\mathcal{A})|. \quad (14.5)$$

Jak zwykle, możemy utożsamić rozkład prawdopodobieństwa na \mathcal{X} z funkcją, przypisującą prawdopodobieństwa pojedynczym punktom $x \in \mathcal{X}$. Zauważmy, że

$$\|\nu - \lambda\|_{\text{tv}} = \frac{1}{2} \sum_{x \in \mathcal{X}} |\nu(x) - \lambda(x)|. \quad (14.6)$$

Istotnie, ponieważ rozpatrujemy dwie miary probabilistyczne, dla których $\nu(\mathcal{X}) = \lambda(\mathcal{X}) = 1$, więc $\|\nu - \lambda\| = \nu(\mathcal{B}) - \lambda(\mathcal{B})$ dla $\mathcal{B} = \{x : \nu(x) > \lambda(x)\}$. Ale $\sum_{x \in \mathcal{B}} (\nu(x) - \lambda(x)) = \sum_{x \in \mathcal{X} \setminus \mathcal{B}} (\lambda(x) - \nu(x)) = \frac{1}{2} \sum_{x \in \mathcal{X}} |\nu(x) - \lambda(x)|$.

Dla zmiennej losowej $X : \Omega \rightarrow \mathcal{X}$ napis $X \sim \nu$ będzie oznaczał fakt, że X ma rozkład prawdopodobieństwa ν , czyli $\mathbb{P}(X = x) = \nu(x)$,

Lemat 14.2. *Jeżeli $X, Y : \Omega \rightarrow \mathcal{X}$ są dwiema zmiennymi losowymi określonymi na tej samej przestrzeni probabilistycznej i $X \sim \nu$ i $Y \sim \lambda$, to*

$$\|\nu - \lambda\|_{\text{tv}} \leq \mathbb{P}(X \neq Y).$$

Dowód. Niech $d = \mathbb{P}(X \neq Y)$. Dla dowolnego $\mathcal{A} \subseteq \mathcal{X}$ mamy

$$\nu(\mathcal{A}) = \mathbb{P}(X \in \mathcal{A}) \leq \mathbb{P}(Y \in \mathcal{A}) + \mathbb{P}(X \neq Y) = \lambda(\mathcal{A}) + d.$$

Symetrycznie, $\nu(\mathcal{A}) \leq \lambda(\mathcal{A}) + d$. Zatem $\|\nu - \lambda\|_{\text{tv}} \leq d$. □

Interesujące jest, że Lemat 14.2 daje się, w pewnym sensie, odwrócić. Co prawda, to nie będzie potrzebne w dowodzie Słabego Twierdzenia Ergodycznego, ale później okaże się bardzo pomocne.

Lemat 14.3. *Jeżeli ν i λ są rozkładami prawdopodobieństwa na \mathcal{X} , to istnieją zmienne losowe X i Y określone na tej samej przestrzeni probabilistycznej, takie, że $X \sim \nu$ i $Y \sim \lambda$ i*

$$\|\nu - \lambda\|_{\text{tv}} = \mathbb{P}(X \neq Y).$$

Dowód. Niech $\|\nu - \lambda\|_{\text{tv}} = d$. Bez straty ogólności możemy przyjąć, że X i Y są zmiennymi losowymi określonymi na przestrzeni probabilistycznej $\Omega = \mathcal{X} \times \mathcal{X}$. Należy podać *łączny* rozkład zmiennych losowych X i Y , czyli miarę probabilistyczną χ na $\mathcal{X} \times \mathcal{X}$ taką, że $\sum_y \chi(x, y) = \nu(x)$, $\sum_x \chi(x, y) = \lambda(y)$ i $\sum_x \chi(x, x) = 1 - d$.

Niech

$$\chi(x, x) = \min(\nu(x), \lambda(x)) = \begin{cases} \nu(x) & \text{dla } x \in \mathcal{A}; \\ \lambda(x) & \text{dla } x \in \mathcal{B}, \end{cases}$$

gdzie $\mathcal{A} = \{x : \nu(x) \leq \lambda(x)\}$ i $\mathcal{B} = \{x : \nu(x) > \lambda(x)\}$.

Mamy oczywiście $d = 1 - \sum_x \chi(x, x)$ i jest jasne, że tabelka łącznego rozkładu $\chi(x, y) = \mathbb{P}(X = x, Y = y)$ musi być postaci macierzy blokowej

$$\begin{array}{c} x \in \mathcal{A} \left\{ \begin{array}{cc} D_{\mathcal{A}} & 0 \end{array} \right. \\ x \in \mathcal{B} \left\{ \begin{array}{cc} G & D_{\mathcal{B}} \end{array} \right. , \\ \underbrace{\hspace{1cm}}_{y \in \mathcal{A}} \quad \underbrace{\hspace{1cm}}_{y \in \mathcal{B}} \end{array}$$

gdzie $D_{\mathcal{A}}$ i $D_{\mathcal{B}}$ są macierzami diagonalnymi. Pozostaje tylko odpowiednio „rozmieścić pozostałą masę prawdopodobieństwa” d w macierzy G . Możemy na przykład przyjąć, dla $x \in \mathcal{B}$ i $y \in \mathcal{A}$,

$$\chi(x, y) = \frac{1}{d} (\nu(x) - \lambda(x)) (\lambda(y) - \nu(y)).$$

Mamy wtedy $\sum_{y \in \mathcal{A}} \chi(x, y) = \nu(x) - \lambda(x)$, więc $\sum_y \chi(x, y) = \nu(x)$ dla $x \in \mathcal{B}$ i podobnie $\sum_x \chi(x, y) = \lambda(y)$ dla $y \in \mathcal{A}$. Określony przez nas rozkład łączny χ ma więc masę $1 - d$ na przekątnej i żądane rozkłady brzegowe. \square

14.3.2. Sprzęganie

Rozważmy „podwójny” łańcuch Markowa (X_n, X'_n) na przestrzeni stanów $\mathcal{X} \times \mathcal{X}$. Przypuśćmy, że każda z dwóch „współrzędnych”, oddzielnie rozpatrywana, jest łańcuchem o macierzy przejścia P . Mówiąc dokładniej, zakładamy, że

$$\begin{aligned} \mathbb{P}(X_{n+1} = y, X'_{n+1} = y' | X_n = x, X'_n = x', X_{n-1}, X'_{n-1}, \dots, X_0, X'_0) \\ = \bar{P}((x, x'), (y, y')), \end{aligned}$$

gdzie macierz przejścia \bar{P} podwójnego łańcucha spełnia następujące warunki:

$$\begin{aligned} \sum_{y'} \bar{P}((x, x'), (y, y')) &= P(x, y) \quad \text{dla każdego } x' \\ \sum_y \bar{P}((x, x'), (y, y')) &= P(x', y') \quad \text{dla każdego } x. \end{aligned} \tag{14.7}$$

Widać, że $X_0, X_1, \dots, X_n, \dots$ jest łańcuchem Markowa z prawdopodobieństwami przejścia P i to samo można powiedzieć o $X'_0, X'_1, \dots, X'_n, \dots$. Załóżmy ponadto, że od momentu, gdy oba łańcuchy się spotkają, dalej „poruszają się” już razem. Innymi słowy,

$$\bar{P}((x, x), (y, y')) = \begin{cases} P(x, y) & \text{jeśli } y = y', \\ 0 & \text{jeśli } y \neq y'. \end{cases}$$

Nazwiemy konstrukcję takiej pary *sprzęganiem* łańcuchów (bardziej znany jest angielski termin *coupling*).

Oznaczmy przez T moment spotkania się łańcuchów:

$$T = \min\{n > 0 : X_n = X'_n\}. \tag{14.8}$$

Podstawową rolę odgrywa następujące spostrzeżenie:

$$\|\mathbb{P}(X_n \in \cdot) - \mathbb{P}(X'_n \in \cdot)\|_{\text{tv}} \leq \mathbb{P}(X_n \neq X'_n) = \mathbb{P}(T > n).$$

Jeśli teraz łańcuch X'_n „wystartuje” z rozkładu stacjonarnego, czyli $X_0 \sim \pi$ to $X_n \sim \pi$ dla każdego n i otrzymujemy

$$\|\mathbb{P}(X_n \in \cdot) - \pi(\cdot)\|_{\text{tv}} \leq \mathbb{P}(T > n). \quad (14.9)$$

Aby udowodnić zbieżność $\mathbb{P}(X_n \in \cdot) \rightarrow \pi(\cdot)$ wystarczy skonstruować parę łańcuchów sprzężonych, które się spotkają z prawdopodobieństwem 1: $\mathbb{P}(T < \infty) = 1$. Możemy teraz udowodnić (10.1), przynajmniej dla łańcuchów na skończonej przestrzeni stanów.

Twierdzenie 14.5 (Słabe Twierdzenie Ergodyczne). *Jeśli łańcuch Markowa na skończonej przestrzeni stanów jest nieprzywiedlny i nieokresowy, to*

$$\|\mathbb{P}(X_n \in \cdot) - \pi(\cdot)\|_{\text{tv}} \rightarrow 0.$$

Dowód. Rozważmy parę łańcuchów sprzężonych, które poruszają się *niezależnie* aż do momentu spotkania.

$$\tilde{P}((x, x'), (y, y')) = \begin{cases} P(x, y)P(x', y') & \text{jeśli } x \neq x', \\ P(x, y) & \text{jeśli } x = x' \text{ i } y = y', \\ 0 & \text{jeśli } x = x' \text{ i } y \neq y'. \end{cases} \quad (14.10)$$

Żeby pokazać, że $\mathbb{P}(T < \infty) = 1$ wystarczy zauważyć, że do przed momentem spotkania, łańcuch podwójny ewoluuje zgodnie z prawdopodobieństwami przejścia

$$\tilde{P}((x, x'), (y, y')) = P(x, y)P(x', y').$$

Łańcuch odpowiadający \tilde{P} jest nieprzywiedlny. Istotnie, możemy znaleźć takie n_0 , że dla $n \geq n_0$ wszystkie elementy macierzy P^n są niezerowe. Stąd $P^n((x, x'), (y, y')) = P^n(x, y)P^n(x', y') > 0$ dla dowolnych x, x', y, y' . Wystarczy teraz powołać się na Wniosek 14.2: podwójny łańcuch z prawdopodobieństwem 1 prędzej czy później dojdzie do każdego punktu przestrzeni $\mathcal{X} \times \mathcal{X}$, a zatem musi dojść do „przekątnej” $\{(x, x) : x \in \mathcal{X}\}$. \square

Uwaga 14.2. W dowodzie Twierdzenia 14.5 wykorzystaliśmy w istotny sposób nieokresowość macierzy przejścia P (dla pojedynczego łańcucha), choć to mogło nie być wyraźnie widoczne. Jeśli P jest nieprzywiedlna ale okresowa, wtedy \tilde{P} jest nieprzywiedlna. Na przykład, niech $\mathcal{X} = \{0, 1\}$ i

$$P = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}.$$

Wtedy, oczywiście, $\tilde{P}^n((0, 0), (0, 1)) = 0$ bo $P^n(0, 0) = 0$ dla nieparzystych n zaś $P^n(0, 1) = 0$ dla parzystych n .

Ten sam trywialny przykład pokazuje, że dla łańcuchów okresowych teza Słabego Twierdzenia Ergodycznego nie jest prawdziwa.

W istocie, przytoczony przez nas dowód Twierdzenia 14.5 daje nieco więcej, niż tylko zbieżność rozkładów. Z Wniosku 14.2 wynika, że

$$\|\mathbb{P}(X_n \in \cdot) - \pi(\cdot)\|_{\text{tv}} \leq c\gamma^n$$

dla pewnych stałych $c < \infty$ i $\gamma < 1$. Dla naszych celów takie ogólnikowe stwierdzenie nie jest wystarczające. Skupimy się na przykładach łańcuchów używanych w algorytmach MCMC,

dla których znajdziemy jawne oszacowania, z konkretnymi stałymi. Zobaczymy, że użycie niezależnych kopii łańcucha w dowodzie Twierdzenia 14.5, wzór (14.10) jest konstrukcją dalece nieoptymalną. W wielu przykładach istnieją łańcuchy sprzężone znacznie szybciej „zmierzające do spotkania”.

Przykład 14.1 (Błądzenie po kostce). Niech $\mathcal{X} = \{0, 1\}^n$ i $\pi = U(\mathcal{X})$, czyli $\pi(x) = 1/2^n$ dla każdego x . Rozważmy łańcuch Markowa X_n , którego krok polega na wylosowaniu jednej, losowo wybranej współrzędnej z rozkładu $(1/2, 1/2)$ na zbiorze $\{0, 1\}$ i pozostawieniu pozostałych współrzędnych bez zmian. Formalnie,

$$P(x, y) = \frac{1}{2d} \sum_{i=1}^d \mathbb{I}(x_{-i} = y_{-i}).$$

Jest to zatem „losowe błądzenie” wzdłuż krawędzi n -wymiarowej kostki lub inaczej próbnik Gibbsa. Rzecz jasna, dokładne genrowanie z rozkładu jednostajnego na kostce jest łatwe i nie potrzebujemy do tego łańcuchów Markowa, ale nie o to teraz chodzi. Chcemy zilustrować jak metoda sprzęgania pozwala oszacować szybkość zbieżności łańcucha na możliwie prostym przykładzie. Skonstruujmy parę łańcuchów sprzężonych w taki sposób: wybieramy współrzędną i oraz losujemy jej nową wartość z rozkładu $(1/2, 1/2)$ po czym zmieniamy w ten sam sposób obie kopie. Formalnie,

$$\bar{P}((x, x'), (y, y')) = \frac{1}{2d} \sum_{i=1}^d \mathbb{I}(x_{-i} = y_{-i}, x'_i = y'_i, y_i = y'_i).$$

Jest jasne, że to jest poprawny *coupling* (sprzęganie), to znaczy spełnione są równania (14.10). Spotkanie obu kopii nastąpi *najpóźniej* w momencie gdy każda ze współrzędnych zostanie wybrana przynajmniej raz. Zatem

$$\mathbb{P}(T > n) \leq \left(1 - \frac{1}{d}\right)^n.$$

Nie trudno wyobrazić sobie, że dla *niezależnego* couplingu określonego wzorem (14.10), czasu oczekiwania na spotkanie obu kopii jest na ogół dużo, dużo dłuższy.

15. Markowskie Monte Carlo VI. Oszacowania dokładności

W tym rozdziale zajmiemy się problemem oszacowania błędu markowskich algorytmów Monte Carlo. W odróżnieniu od wstępnych rozważań w Rozdziale 10, będą nas interesowały ściśle nierówności, a nie oceny oparte na twierdzeniach granicznych. Skupimy się na rozkładzie spektralnym macierzy przejścia. Jest to najbardziej znana i zapewne najskuteczniejsza metoda otrzymywania dobrych oszacowań, przynajmniej dla łańcuchów odwracalnych na przestrzeni skończonej.

15.1. Reprezentacja spektralna macierzy odwracalnej

Rozważmy łańcuch nieprzywiedlny i odwracalny z macierzą przejścia P i rozkładem stacjonarnym π^\top . Zakładamy więc, że dla dowolnych $x, y \in \mathcal{X}$ spełniona jest zależność

$$\pi(x)P(x, y) = \pi(y)P(y, x).$$

Niech

$$\Pi = \text{diag}(\pi) = \begin{pmatrix} \pi(1) & 0 & \cdots & 0 \\ 0 & \pi(2) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \pi(d) \end{pmatrix}.$$

Warunek odwracalności możemy zapisać w postaci $\Pi P = P^\top \Pi$. Wyposaźmy przestrzeń \mathbb{R}^d w iloczyn skalarny

$$\langle f, g \rangle_\pi = f^\top \Pi g = \sum_x f(x)g(x)\pi(x),$$

gdzie $f, g \in \mathbb{R}$ traktujemy jak wektory kolumnowe. Oczywiście, norma funkcji f jest zdefiniowana wzorem $\|f\|_\pi^2 = \langle f, f \rangle_\pi$. Macierz P traktujemy jako operator działający z lewej strony na funkcje, z prawej strony na rozkłady prawdopodobieństwa. Zgodnie z regułami mnożenia macierzy i wektorów, wzory na Pf i $\xi^\top P$ są następujące:

$$Pf(x) = \sum_y P(x, y)f(y), \quad \xi^\top P(y) = \sum_x \xi(x)P(x, y).$$

Odwracalność P implikuje

$$\langle f, Pg \rangle_\pi = f^\top \Pi Pg = f^\top P^\top \Pi g = g^\top \Pi Pf = \langle Pf, g \rangle_\pi.$$

Znaczy to, że P jest macierzą operatora samosprężonego względem iloczynu skalarnego $\langle \cdot, \cdot \rangle_\pi$. W skrócie powiemy, że P jest π -samosprężona. Wartości własne P są rzeczywiste i zawarte w przedziale $[-1, 1]$. Uporządkujmy je w kolejności malejącej:

$$1 = \lambda_0 > \lambda_1 \geq \cdots \lambda_{d-1} \geq -1.$$

Wiadomo, że $\lambda_0 = 1$ jest pojedynczą wartością własną. Jeśli łańcuch jest nieokresowy, to $\lambda_{d-1} > -1$. Niech

$$\Lambda = \text{diag}(1, \lambda_1, \dots, \lambda_{d-1}).$$

Reprezentacja spektralna macierzy P jest następująca:

$$P = V\Lambda V^\top \Pi,$$

gdzie

$$V^\top \Pi V = I$$

lub, równoważnie, $VV^\top \Pi = I$. Zauważmy, że kolumny macierzy

$$V = (1, v_1, \dots, v_{d-1})$$

są prawostronnymi wektorami własnymi macierzy P tworzącymi bazę π -ortonormalną. Mamy więc $Pv_i = \lambda v_i$ (oczywiście, $v_0 = 1$ jest tu wektorem jedynek) oraz

$$\langle v_i, v_j \rangle_\pi = v_i^\top \Pi v_j = \mathbb{I}(i = j).$$

Lewostronne wektory własne P (zapisane wierszowo) są postaci $v_i^\top \Pi$: mamy $v_i^\top \Pi P = \lambda_i v_i^\top \Pi$. W szczególności, $v_0^\top \Pi = \pi^\top$. Zapiszmy reprezentację spektralną P w bardziej jawnej formie:

$$P = \sum_{i \geq 0} \lambda_i v_i v_i^\top \Pi,$$

przy tym

$$\sum_{i \geq 0} v_i v_i^\top \Pi = I.$$

Zauważmy jeszcze, że pierwszy (lub raczej - zerowy) składnik jest macierzą stabilną (o jednokowych wierszach):

$$v_0 v_0^\top \Pi = 1\pi^\top.$$

Reprezentacja spektralna prowadzi do zgrabnych wyrażeń na potęgę macierzy. Latwo zauważyć, że $P^n = V\Lambda^n V^\top \Pi$, czyli

$$P^n = \sum_{i \geq 0} \lambda_i^n v_i v_i^\top \Pi = 1\pi^\top + \sum_{i \geq 1} \lambda_i^n v_i v_i^\top \Pi. \quad (15.1)$$

Wiemy, że wszystkie wartości własne λ_i z wyjątkiem zerowej ($\lambda_0 = 1$) oraz, być może, ostatniej ($\lambda_{d-1} \geq -1$) są co do modułu mniejsze niż 1. Jeżeli więc $\lambda_{d-1} > -1$, to wszystkie składniki sumy we wzorze (15.1) z wyjątkiem początkowego zmierzają do zera i w rezultacie

$$P^n \rightarrow 1\pi^\top \quad (n \rightarrow \infty).$$

Jest to nic innego jak teza Słabego Twierdzenia Ergodycznego (Twierdzenie 14.5), otrzymana zupełnie inną metodą, przy założeniu odwracalności. Można pokazać, że warunek $\lambda_{d-1} > -1$ jest równoważny nieokresowości, i jest konieczny (jeśli $\lambda_{d-1} = -1$ to łańcuch ma okres 2).

Następujący lemat będzie podstawą dalszych rozważań i umożliwi „przerobienie” STE na jawne wyniki. Zdefiniujemy

$$\lambda = \max(\lambda_1, |\lambda_{d-1}|).$$

Założmy, że łańcuch jest nieokresowy, więc $\lambda < 1$. Pokażemy, że operator P ograniczony do podprzestrzeni $\{f : f \perp 1\} \subset \mathbb{R}^d$ ortogonalnej do funkcji stałych jest zwężający, ze stałą $\lambda < 1$.

Lemat 15.1. *Jeżeli $\langle f, 1 \rangle_\pi = 0$ to $\|Pf\|_\pi \leq \lambda \|f\|_\pi$.*

Dowód. Wystarczy zauważyć, że

$$\begin{aligned}
 \langle Pf, Pf \rangle_\pi &= \langle f, P^2 f \rangle_\pi \\
 &= f^\top \Pi 1 \pi^\top f + f^\top \sum_{i \geq 1} \lambda_i^{2n} v_i v_i^\top \Pi f \\
 &= \sum_{i \geq 1} \lambda_i^{2n} \langle v_i, f \rangle_\pi^2 \\
 &\leq \lambda^{2n} \sum_{i \geq 1} \langle v_i, f \rangle_\pi^2 \\
 &= \lambda^{2n} \langle f, f \rangle.
 \end{aligned}$$

Korzystamy tu z faktu, że P jest samosprężony, ze wzoru (15.1) i z tego, że $f^\top \Pi 1 = 0$. \square

15.1.1. Oszacowanie szybkości zbieżności

Przejdźmy teraz do jawnych oszacowań szybkości zbieżności w STE. Wyniki zawarte z tym podrozdziałem pochodzą z pracy Diaconisa i Strooka [5]. Dla rozkładu prawdopodobieństwa ξ^\top definiujemy „odległość” χ^2 od rozkładu stacjonarnego wzorem

$$\chi^2(\pi, \xi) = \sum_x \frac{(\xi(x) - \pi(x))^2}{\pi(x)} = (\xi^\top - \pi^\top) \Pi^{-1} (\xi - \pi).$$

Ta „odległość” nie ma własności symetrii, więc nie jest metryką, ale to nie przeszkadza. Istotna jest interpretacja χ^2 jako „odstępstwa od stacjonarności”. Wykorzystamy podejście spektralne, w szczególności Lemat 15.1. Niech $\chi = \sqrt{\chi^2}$.

Stwierdzenie 15.1 (Diaconis i Strook). *Jeśli łańcuch odwracalny, nieprzywiedlny i nieokresowy ma rozkład początkowy ξ , to*

$$\chi(\pi, P^n \xi) \leq \lambda^n \chi(\pi, \xi).$$

Dowód. Zastosujmy Lemat 15.1 do wektora $\Pi^{-1}(\xi - \pi)$ który, jak łatwo zauważyć, jest prostopadły do 1. Otrzymujemy

$$\chi(\pi, P^n \xi) = \|P^n \Pi^{-1}(\xi - \pi)\|_\pi \leq \lambda^n \|\Pi^{-1}(\xi - \pi)\|_\pi = \lambda^n \chi(\pi, \xi).$$

\square

Wniosek 15.1. *Dla łańcucha o rozkładzie początkowym skupionym w punkcie x ,*

$$\chi(\pi, P^n(x, \cdot)) \leq \lambda^n \sqrt{\frac{1 - \pi(x)}{\pi(x)}} \leq \frac{\lambda^n}{\sqrt{\pi(x)}}.$$

Istotnie, mamy jeszcze jedno wyrażenie na „odległość” χ^2 : dla dowolnego rozkładu ξ ,

$$\chi^2(\pi, \xi) = \xi^\top \Pi^{-1} \xi - 1.$$

Wystarczy teraz podstawić $\xi(y) = \mathbb{I}(y = x)$, aby otrzymać $\chi^2 = 1/\pi(x) - 1$.

15.1.2. Oszacowanie normy pełnego wahania

Istnieje prosta nierówność pomiędzy normą pełnego wahania i „odległością” χ^2 :

$$\|\xi - \pi\|_{\text{tv}} = \frac{1}{2} \sum_x |\xi(x) - \pi(x)| \leq \frac{1}{2} \chi(\pi, \xi).$$

Wynika to z następującego rachunku:

$$\begin{aligned} 4 \|\xi - \pi\|_{\text{tv}}^2 &= \left(\sum_x \frac{|\xi(x) - \pi(x)|}{\sqrt{\pi(x)}} \sqrt{\pi(x)} \right)^2 \\ &\leq \sum_x \frac{(\xi(x) - \pi(x))^2}{\pi(x)} \sum_x \pi(x) \quad [\text{Cauchy-Schwarz}] \\ &= \chi^2(\pi, \xi). \end{aligned}$$

Stąd natychmiast otrzymujemy wniosek

$$\|P_n(x, \cdot) - \pi\|_{\text{tv}} \leq \frac{1}{2} \lambda^n \sqrt{\frac{1 - \pi(x)}{\pi(x)}} \leq \frac{\lambda^n}{2\sqrt{\pi(x)}}.$$

15.1.3. Oszacowanie obciążenia estymatora

Rozważmy zadanie obliczania wartości oczekiwanej $\theta = \mathbb{E}_\pi f = \pi^\top f$ dla pewnej funkcji f . Niech $\bar{f} = f - 1\pi^\top f$ oznacza „scentrowaną” funkcję f . Natychmiast widać, że $\bar{f} \perp 1$ i możemy zastosować Lemat 15.1. Stąd już tylko mały krok do oszacowania różnicy między wartością oczekiwaną $\mathbb{E}_\xi f(X_n) = \xi^\top P^n f$ i wartością stacjonarną θ .

Stwierdzenie 15.2. *Jeśli łańcuch jest odwracalny, nieprzywiedlny, nieokresowy i ma rozkład początkowy ξ , to*

$$|\mathbb{E}_\xi f(X_n) - \theta| \leq \lambda^n \chi(\pi, \xi) \sigma(f),$$

gdzie $\sigma^2(f) = \|\bar{f}\|_\pi^2$ jest wariancją stacjonarną funkcji f .

Dowód. Mamy

$$\begin{aligned} |\mathbb{E}_\xi f(X_n) - \theta| &= |(\xi^\top - \pi^\top) P^n f| = |(\xi^\top - \pi^\top) P^n \bar{f}| \\ &= |\langle \Pi^{-1}(\xi - \pi), P^n \bar{f} \rangle_\pi| \\ &\leq \|\langle \Pi^{-1}(\xi - \pi) \rangle_\pi\|_\pi \|P^n \bar{f}\|_\pi \quad [\text{Cauchy-Schwarz}] \\ &\leq \chi(\pi, \xi) \lambda^n \|\bar{f}\|_\pi \quad [\text{Lemat 15.1}]. \end{aligned}$$

□

Rozważmy teraz naturalny estymator

$$\hat{\theta}_{t,n} = \frac{1}{n} \sum_{i=t}^{t+n-1} f(X_i).$$

Jest to średnia wzdłuż trajektorii łańcucha, długości n i „opóźniona” o t . Idea jest jasna: ignorujemy początkowy odcinek trajektorii długości t (tak zwany okres *burn-in*) aby dać łańcuchowi czas na zbliżenie od rozkładu stacjonarnego. Później obliczymy średnią. W ten sposób redukujemy obciążenie. Precyzuje to następujący wniosek.

Wniosek 15.2. *Dla dowolnego n mamy*

$$|\mathbb{E}_\xi \hat{\theta}_{t,n} - \theta| \leq \frac{\lambda^t}{1 - \lambda} \chi(\pi, \xi) \sigma(f).$$

Wynika to z nierówności trójkąta i wzoru na sumę szeregu geometrycznego:

$$|\mathbb{E}_\xi \hat{\theta}_{t,n} - \theta| \leq \sum_{i=t}^{t+n-1} |\mathbb{E}_\xi f(X_i) - \theta| \leq \sum_{i=t}^{\infty} \lambda^i \chi(\pi, \xi) \sigma(f).$$

Zwróćmy uwagę, że obciążenie maleje w tempie geometrycznym przy $t \rightarrow \infty$ ale zachowuje się zaledwie jak $O(1/n)$ przy ustalonym t i $n \rightarrow \infty$ (ponieważ początkowe wyrazy sumy mają na obciążenie wpływ dominujący).

15.2. Oszacowanie błędu średniokwadratowego estymatora

Oszacowanie błędu średniokwadratowego (BŚK) estymatora MCMC jest znacznie trudniejsze i subtelniejsze, niż obciążenia.

15.2.1. Asymptotyczna wariancja

Zacznijmy od wyprowadzenia kolejnego wzoru na asymptotyczną wariancję. Przypomnijmy, że zgodnie ze Stwierdzeniem 14.1,

$$\sigma_{\text{as}}^2(f) = f^\top C f,$$

gdzie

$$\begin{aligned} C &= \Pi(2Z - I - 1\pi^\top) = 2\Pi Z - \Pi - \pi\pi^\top \\ &= \Pi(2A - I + 1\pi^\top) = 2\Pi A - \Pi + \pi\pi^\top. \end{aligned}$$

Skorzystajmy z reprezentacji spektralnej macierzy P :

$$P - 1\pi^\top = \sum_{i \geq 1} \lambda_i v_i v_i^\top \Pi = V \left(\begin{array}{c|ccc} 0 & & \cdots & 0 \\ \hline & \lambda_1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_{d-1} \end{array} \right) V^\top \Pi$$

Z definicji macierzy A i wzoru (15.1), ponieważ $\sum_{n=0}^{\infty} \lambda_i^n = 1/(1 - \lambda_i)$, więc

$$A = \sum_{n=0}^{\infty} (P^n - 1\pi^\top) = V \left(\begin{array}{c|ccc} 0 & & \cdots & 0 \\ \hline & \ddots & \cdots & \\ \vdots & \vdots & 1/(1 - \lambda_i) & \vdots \\ 0 & & \cdots & \ddots \end{array} \right) V^\top \Pi.$$

Wreszcie, ponieważ $2/(1 - \lambda_i) - 1 = (1 + \lambda_i)/(1 - \lambda_i)$, więc

$$\Pi(2A - I + 1\pi^\top) = \Pi V \left(\begin{array}{c|ccc} 0 & & \cdots & 0 \\ \hline & \ddots & \cdots & \\ \vdots & \vdots & (1 + \lambda_i)/(1 - \lambda_i) & \vdots \\ 0 & & \cdots & \ddots \end{array} \right) V^\top \Pi.$$

Zauważmy teraz, że wektor $\alpha = V^\top \Pi f$ zawiera współrzędne wektora f w bazie ON złożonej z prawych wektorów własnych: $\alpha_i = v_i^\top \Pi f = \langle v_i, f \rangle_\pi$. Udowodniliśmy w ten sposób następujący fakt.

Stwierdzenie 15.3. *Dla łańcucha odwracalnego, wzór na asymptotyczną wariancję przybiera postać*

$$\sigma_{\text{as}}^2(f) = \sum_{i \geq 1} \alpha_i^2 \frac{1 + \lambda_i}{1 - \lambda_i},$$

gdzie $\alpha_i = v_i^\top \Pi f = \langle v_i, f \rangle_\pi$.

Wynika stąd ważna nierówność.

Wniosek 15.3. *Dla łańcucha odwracalnego mamy następujące oszacowanie asymptotycznej wariancji:*

$$\sigma_{\text{as}}^2(f) \leq \frac{1 + \lambda_1}{1 - \lambda_1} \sigma^2(f)$$

gdzie λ_1 jest największą wartością własną mniejszą od 1, a $\sigma^2(f)$ jest wariancją stacjonarną.

Istotnie,

$$\sigma_{\text{as}}^2(f) \leq \sum_{i \geq 1} \alpha_i^2 \frac{1 + \lambda_i}{1 - \lambda_i} \leq \frac{1 + \lambda_1}{1 - \lambda_1} \sum_{i \geq 1} \alpha_i^2,$$

a łatwo widzieć, że $\sum_{i \geq 1} \alpha_i^2 = \|\bar{f}\|_\pi^2 = \sigma^2(f)$. Zwróćmy uwagę, że we Wniosku 15.3 występuje λ_1 , a nie $\lambda = \max(\lambda_1, |\lambda_{d-1}|)$ (największa co do modułu wartość własna mniejsza od 1).

Na zakończenie przytoczę jeszcze kilka sugestywnych wzorów.

Uwaga 15.1. Macierz $L = I - P$ jest nazywana laplasjanem. Zauważmy, że

$$L = \sum_{i \geq 1} (1 - \lambda_i) v_i v_i^\top \Pi.$$

Ponieważ

$$A = \sum_{i \geq 1} \frac{1}{1 - \lambda_i} v_i v_i^\top \Pi,$$

uzasadnia to interpretację macierzy A jako „uogólnionej odwrotności” laplasjanu.

Dorzućmy przy okazji jeszcze jedno wyrażenie na asymptotyczną wariancję:

$$\sigma_{\text{as}}^2(f) = \langle f, (2A - \mathbb{I} + \mathbf{1}\pi^\top) f \rangle_\pi = \langle f, Af \rangle_\pi + \langle Af, f \rangle_\pi - \sigma^2(f).$$

Jeśli $\pi^\top f = 0$, to $\sigma^2(f) = \text{Var}_\pi f = \langle f, f \rangle_\pi$ i ostatni wzór możemy przepisać w postaci

$$\sigma_{\text{as}}^2(f) = \sigma^2(f) \left(2 \frac{\langle f, Af \rangle_\pi}{\langle f, f \rangle_\pi} - 1 \right)$$

15.2.2. Oszacowanie BŚK

Wyniki w tym podrozdziale zostały otrzymane przez Aldousa [1]. Pomysł polega na tym, żeby najpierw otrzymać nierówność dla łańcucha stacjonarnego, a potem postarać się o uogólnienie dla łańcucha o dowolnym rozkładzie początkowym. Przypomnijmy oznaczenia $\theta = \mathbb{E}_\pi f$, $\hat{\theta}_n = \sum_{i=0}^{n-1} f(X_i)$.

Stwierdzenie 15.4 (Aldous, 1987). *Dla łańcucha nieprzywiedlnego, odwracalnego i stacjonarnego, MSE_π można oszacować w następujący sposób:*

$$\mathbb{E}_\pi(\hat{\theta}_n - \theta)^2 \leq \frac{1 + \tilde{\lambda}}{1 - \tilde{\lambda}} \sigma^2(f)$$

gdzie $\tilde{\lambda} = \max(\lambda_1, 0)$, zaś λ_1 jest największą wartością własną mniejszą od 1.

Dowód. Korzystamy ze stacjonarności i z rozkładu spektralnego macierzy P .

$$\begin{aligned}
\mathbb{E}_\pi(\hat{\theta}_n - \theta)^2 &= \mathbb{E}_\pi \left(\frac{1}{n} \sum_{i=0}^{n-1} \bar{f}(X_i) \right)^2 \\
&= \frac{1}{n^2} \sum_{i=0}^{n-1} \mathbb{E}_\pi \bar{f}(X_i)^2 + \frac{2}{n^2} \sum_{i=0}^{n-1} \sum_{j=i+1}^{n-1} \mathbb{E}_\pi \bar{f}(X_i) \bar{f}(X_j) \\
&= \frac{1}{n} \mathbb{E}_\pi \bar{f}(X_0)^2 + \frac{2}{n^2} \sum_{k=1}^{n-1} (n-k) \mathbb{E}_\pi \bar{f}(X_0) \bar{f}(X_k) \quad [k = j - i] \\
&= \frac{1}{n} \langle \bar{f}, \bar{f} \rangle_\pi + \frac{2}{n^2} \sum_{k=1}^{n-1} (n-k) \langle \bar{f}, P^k \bar{f} \rangle_\pi \\
&= \frac{1}{n} \langle \bar{f}, \bar{f} \rangle_\pi + \frac{2}{n^2} \sum_{k=1}^{n-1} (n-k) \sum_{i \geq 1} \tilde{\lambda}_i^k \langle v_i, \bar{f} \rangle_\pi^2 \\
&= \frac{1}{n} \langle \bar{f}, \bar{f} \rangle_\pi + \frac{2}{n^2} \sum_{i \geq 1} \langle v_i, \bar{f} \rangle_\pi^2 \sum_{k=1}^{n-1} (n-k) \lambda_i^k \\
&\leq \frac{1}{n} \sigma^2(f) + \frac{2}{n^2} \sum_{\substack{i \geq 1 \\ \lambda_i > 0}} \langle v_i, \bar{f} \rangle_\pi^2 \sum_{k=1}^{n-1} (n-k) \lambda_i^k \quad [\text{pomijamy składniki dla } \lambda_i < 0] \\
&\leq \frac{1}{n} \sigma^2(f) + \frac{2}{n^2} \sigma^2(f) \sum_{k=1}^{n-1} (n-k) \tilde{\lambda}^k \\
&\leq \frac{1}{n} \sigma^2(f) + \frac{2}{n} \sigma^2(f) \sum_{k=1}^{n-1} \tilde{\lambda}^k \\
&\leq \frac{2}{n} \sigma^2(f) \left(1 + \frac{2\tilde{\lambda}}{1 - \tilde{\lambda}} \right) = \frac{1}{n} \sigma^2(f) \frac{1 + \tilde{\lambda}}{1 - \tilde{\lambda}}.
\end{aligned}$$

Zwróćmy uwagę na miejsce, w którym pomijamy składniki odpowiadające ujemnym wartościom własnym. Uzasadnienie jest takie, że dla $\lambda_i < 0$, składniki sumy $\sum_{k=1}^{n-1} (n-k) \lambda_i^k$ są naprzemiennie ujemne i dodatnie, o malejących wartościach bezwzględnych. Stąd wynika, że $\sum_{k=1}^{n-1} (n-k) \lambda_i^k < 0$ i można tę sumę w nierówności opuścić. Jest to ciekawe zjawisko: *ujemne wartości własne pomagają, zmniejszając błąd!* Działają podobnie jak zmienne antytetyczne. \square

Niech teraz $E(x)$ oznacza błąd średniokwadratowy dla łańcucha startującego z punktu $x \in \mathcal{X}$,

$$e_n(x) = \mathbb{E}_x(\hat{\theta}_n - \theta)^2.$$

Rozpatrzmy łańcuch o rozkładzie początkowym ξ . Niech

$$\hat{\theta}_{t,n} = \frac{1}{n} \sum_{i=t}^{t+n-1} f(X_i).$$

będzie średnią długości n obliczany po odrzuceniu t początkowych zmiennych.

Stwierdzenie 15.5. *Dla dla łańcucha nieprzywiedlnego, odwracalnego startującego z dowolnego rozkładu ξ , MSE_ξ można oszacować w następujący sposób:*

$$\mathbb{E}_\xi(\hat{\theta}_{t,n} - \theta)^2 \leq \frac{1 + \tilde{\lambda}}{1 - \tilde{\lambda}} \sigma^2(f) + \lambda^t \chi(\pi, \xi) \max_x |\bar{f}(x)|^2$$

gdzie $\tilde{\lambda} = \max(\lambda_1, 0)$ i $\lambda = \max(\lambda_1, |\lambda_{d-1}|)$.

Dowód. Na mocy Stwierżeń 15.4 i 15.2 mamy

$$\begin{aligned}\mathbb{E}_\xi(\hat{\theta}_{t,n} - \theta)^2 &= \mathbb{E}_\xi e_n(X_t) \leq \mathbb{E}_\pi e_n + |\mathbb{E}_\xi e_n(X_t) - \mathbb{E}_\pi e_n| \\ &\leq \frac{1}{n} \sigma^2(f) \frac{1 + \tilde{\lambda}}{1 - \tilde{\lambda}} + \lambda^t \chi(\pi, \xi) \sigma(e_n) \\ &\leq \frac{1}{n} \sigma^2(f) \frac{1 + \tilde{\lambda}}{1 - \tilde{\lambda}} + \lambda^t \chi(\pi, \xi) \max_x |\bar{f}(x)|^2,\end{aligned}$$

bo $e_n(y) \leq \max_x |\bar{f}(x)|^2$, a więc $\sigma(e_n) \leq \max_x |\bar{f}(x)|^2$. \square

Nierówność podane przez Aldousa w cytowanej pracy była nieco inna.

Stwierdzenie 15.6 (Aldous, 1987). *Dla dla łańcucha nieprzywiedlnego, odwracalnego startującego z dowolnego rozkładu ξ mamy następujące oszacowanie błędu MSE_ξ :*

$$\mathbb{E}_\xi(\hat{\theta}_{t,n} - \theta)^2 \leq \left(1 + \frac{\lambda^t}{\pi_*}\right) \frac{1 + \tilde{\lambda}}{1 - \tilde{\lambda}} \sigma^2(f),$$

gdzie $\pi_* = \min_x \pi(x)$.

Dowód. Istotnie, założmy, że łańcuch startuje z deterministycznie wybranego punktu $x \in \mathcal{X}$, czyli $\xi = P(x, \cdot)$. Jeśli otrzymamy oszacowanie niezależne od x , dowód będzie zakończony.

$$\begin{aligned}\mathbb{E}_x(\hat{\theta}_{t,n} - \theta)^2 &= \sum_y P^t(x, y) e_n(y) \\ &= \sum_y P^t(x, y) e_n(y) = \sum_y \frac{P^t(x, y)}{\pi(y)} \pi(y) e_n(y) \\ &\leq \left(1 + \frac{\lambda^t}{\pi_*}\right) \sum_y \pi(y) e_n(y) = \left(1 + \frac{\lambda^t}{\pi_*}\right) MSE_\pi.\end{aligned}$$

i zastosujmy Stwierdzenie 15.4. Wykorzystaliśmy tu nierówność $P^t(x, y)/\pi(y) \leq 1 + \lambda^t/\pi_*$, która wynika z rozkładu spektralnego:

$$\begin{aligned}\frac{P^t(x, y)}{\pi(y)} &= P^t \Pi^{-1}(x, y) = 1_x^\top P^t \Pi^{-1} 1_y \\ &= 1_x^\top \sum_{i=0}^{d-1} \lambda_i^t v_i v_i^\top 1_y = 1 + 1_x^\top \sum_{i=1}^{d-1} \lambda_i^t v_i v_i^\top 1_y \\ &\leq 1 + \lambda^t \sum_{i=1}^{d-1} |1_x^\top v_i v_i^\top 1_y| \\ &\leq 1 + \lambda^t \sqrt{\sum_{i=1}^{d-1} (1_x^\top v_i)^2} \sqrt{\sum_{i=1}^{d-1} (v_i 1_y^\top)^2} \\ &\leq 1 + \frac{\lambda^t}{\sqrt{\pi(x)} \sqrt{\pi(y)}} \leq 1 + \frac{\lambda^t}{\pi_*}.\end{aligned}$$

W powyższym wzorze symbol 1_x^\top oznacza wektor $(0, \dots, 0, 1, 0, \dots, 0)$, gdzie jedynka stoi na x -tym miejscu. Skorzystaliśmy z nierówności Schwarza. \square

Literatura

- [1] D. Aldous: On the Markov Chain Simulation Method for Uniform Combinatorial Distributions and Simulated Annealing, *Probability in the Engineering and Informational Science*, pp. 33–45, 1987.
- [2] S. Asmussen and P.W. Glynn: *Stochastic Simulation, Algorithms and Analysis*, Springer, 2007.
- [3] S. Asmussen: *Ruin Probabilities*, World Scientific, 2000, 2001.
- [4] P. Bremaud: *Markov Chains: Gibbs Fields, Monte Carlo Simulation and Queues*, Springer Verlag, 1999.
- [5] P. Diaconis, D. Strook (1991): *Geometric bounds for eigenvalues of Markov chains*, *Annals of Applied Probability* 1 (1), 36–61.
- [6] S. Geman and D. Geman (1984): Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE-PAMI*, 6, 721–741.
- [7] C.J. Geyer (1992): Practical Markov Chain Monte Carlo. *Statistical Science* 7 (4), 473–511.
- [8] C.J. Geyer (1995, 2005): Markov chain Monte Carlo Lecture Notes. Dostępne na www.stat.umn.edu/geyer.
- [9] C.J. Geyer, E.A. Thompson (1992). Constrained Monte Carlo maximum likelihood for dependent data. *J. R. Statistical Society B*, 54, 3, 657–699.
- [10] W.K. Hastings (1970): Monte Carlo sampling methods using Markov chains and their applications, *Biometrika* 57, 97–109.
- [11] M. Jerrum, A. Sinclair (1996): The Markov chain Monte Carlo method: an approach to approximate counting and integration, In *Approximation Algorithms for NP-hard Problems*, (Dorit Hochbaum, ed.), PWS, 1996.
- [12] : M. Jerrum (1998): Mathematical foundations of the Markov chain Monte Carlo method, In *Probabilistic Methods for Algorithmic Discrete Mathematics*, Springer 1998.
- [13] G. Jones (2004): On the Markov chain Central Limit Theorem, *Probability Surveys* 1, 299–320.
- [14] F.K.C. Kingman: *Procesy Poissona*, PWN 2002.
- [15] J.S. Liu: *Monte Carlo Strategies in Scientific Computing*, Springer 2004.
- [16] N. Metropolis, A.W. Rosenbluth, M.N. Rosenbluth, A.H. Teller, E. Teller (1953): Equation of state calculation by fast computing machines, *Journal of Chemical Physics*, 21 (6), 1087–1092.
- [17] E. Nummelin (2002): MC's for MCMC'ists. *International Statistical Review*, 70, 215–240.
- [18] B.D. Ripley: *Stochastic Simulation*, Wiley & Sons, 1987.
- [19] C.P. Robert, G. Casella: *Monte Carlo Statistical Methods*, Springer 2004.
- [20] G.O. Roberts, J.S. Rosenthal (2004): General state space Markov chains and MCMC algorithms. *Probability Surveys* 1, 20–71.
- [21] J.S. Rosenthal (1995): Rates of convergence for Gibbs sampling for variance component models, *Annals of Statistics* 23, 740–761.
- [22] M. Rybiński: *Krótkie wprowadzenie do R dla programistów, z elementami statystyki opisowej*, WMIM UW 2009.
- [23] R. Zieliński, R. Wieczorkowski: *Komputerowe generatory liczb losowych*, WNT, Warszawa, 1997.