

Algorytmy Statystyki Praktycznej

Błażej Miasojedow

Instytut Matematyczny PAN

7 kwietnia 2018

Podstawowe informacje.

- Błażej Miasojedow, email: bmiasojedow@impan.pl
- Konsultacje po wykładach.
- www.mimuw.edu.pl/~bmia/algorytmy

Zasady zaliczenia przedmiotu.

- Laboratorium : 20 pkt.
- Egzamin : 20 pkt.
- Ocena końcowa : **Laboratorium + Egzamin**

Egzamin odbędzie się na ostatnim wykładzie.

O czym będzie wykład?

Większość metod statystycznych w praktyce wymaga dodatkowych metod obliczeniowych, pozwalających na efektywne wyznaczanie estymatorów lub innych wielkości w realnych problemach.

- 1 W statystyce częstościowej, najpopularniejszym o dobrych własnościach estymatorem jest ENW postaci:

$$\hat{\theta} = \arg \max_{\theta \in \Theta} L(\theta).$$

- 2 W statystyce bayesowskiej estymatory są postaci:

$$\hat{\theta} = \mathbb{E}(\theta|X),$$

gdzie typowo nie można wyznaczyć analitycznie rozkładu $\theta|X$.

O czym będzie wykład?

Przedstawimy metody statystyki obliczeniowej ich własności i ograniczenia. Wykład będzie się składał z 3 części:

- 1 Optymalizacja (wypukła).
- 2 Wnioskowanie dla ukrytych modeli Markowa.
- 3 Metody MCMC.

Spis treści

1 Optymalizacja

- Wstęp
- Metody gradientowe
- Algorytm Newtona-Raphsona
- Iterowane ważone najmniejsze kwadraty
- Proximal gradient
- ADMM
- Expectation Maximization

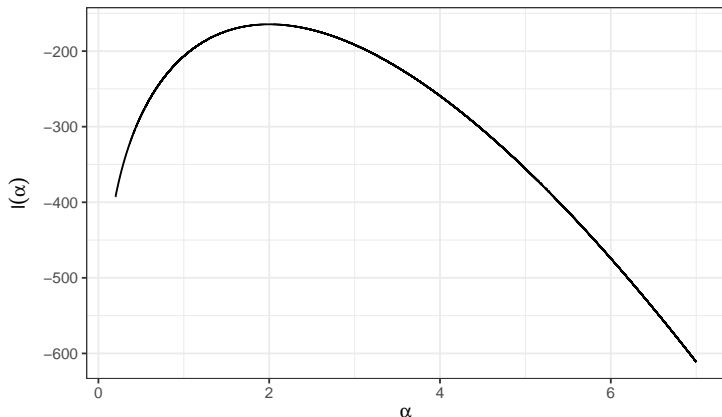
1 Optymalizacja

- Wstęp
- Metody gradientowe
- Algorytm Newtona-Raphsona
- Iterowane ważone najmniejsze kwadraty
- Proximal gradient
- ADMM
- Expectation Maximization

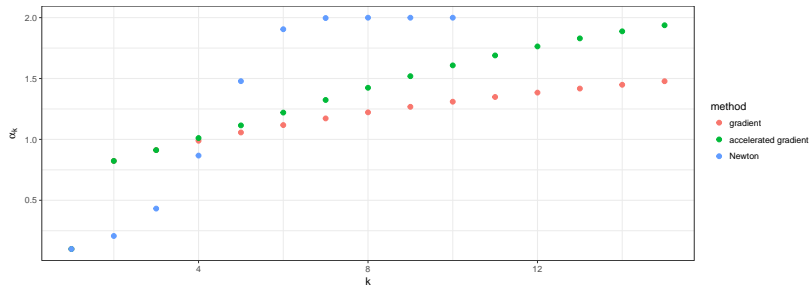
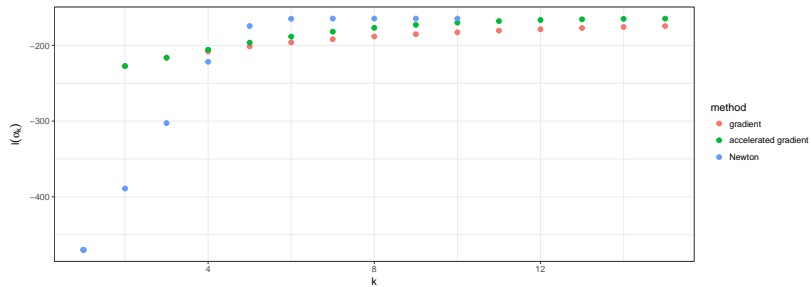
Problem 1.

Niech X_1, X_2, \dots, X_n próbka z rozkładu $\Gamma(\alpha, 1)$. Jak wyznaczyć estymator największej wiarygodności α ?

$$\hat{\alpha} = \arg \max_{\alpha} L(\alpha) = \Gamma(\alpha)^n \prod_{i=1}^n x_i^{\alpha-1}.$$



Problem 1. cd.



1 Optymalizacja

- Wstęp
- Metody gradientowe
- Algorytm Newtona-Raphsona
- Iterowane ważone najmniejsze kwadraty
- Proximal gradient
- ADMM
- Expectation Maximization

Spadek po gradiencie

Problem:

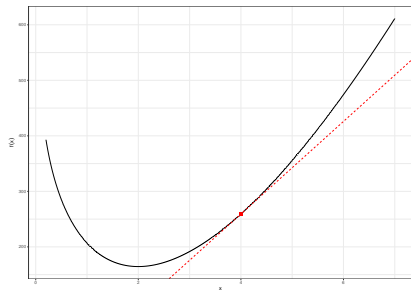
$$\arg \min_{\mathbf{x}} f(\mathbf{x}),$$

f -wypukła, klasy C^1 .

Algorytm:

$$\mathbf{x}_k = \mathbf{x}_{k-1} - \gamma_k \nabla f(\mathbf{x}_{k-1})$$

- Jakie długości kroków wybrać? (stałe, malejące)
- Czy i jak szybko zbiega ten algorytm?



Funkcje wypukłe

Definicja

Funkcja $f: X \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$ jest wypukła wtedy i tylko wtedy gdy dla dowolnego $t \in [0, 1]$ oraz dowolnych $x, y \in X$ zachodzi:

$$f(tx + (1 - t)y) \leq tf(x) + (1 - t)f(y).$$

Stwierdzenie

Różniczkowalna funkcja f jest wypukła na zbiorze wypukłym X wtedy i tylko wtedy gdy dla dowolnych $x, y \in X$ zachodzi:

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle.$$

Dowód

Niech f wypukła stąd

$$f(ty + (1 - t)x) \leq tf(y) + (1 - t)f(x).$$

Z definicji gradientu mamy

$$\begin{aligned} \langle \nabla f(x), y - x \rangle &= \lim_{t \rightarrow 0^+} \frac{f(x + t(y - x)) - f(x)}{t} \\ &\leq \lim_{t \rightarrow 0^+} \frac{tf(y) - tf(x)}{t} = f(y) - f(x) \end{aligned}$$

Dowód c.d.

Niech $z = tx + (1 - t)y$

$$f(x) \geq f(z) + \langle \nabla f(z), x - z \rangle$$

$$f(y) \geq f(z) + \langle \nabla f(z), y - z \rangle$$

Mnożymy nierówności przez t i $1 - t$ odpowiednio i dodajemy stronami. Otrzymujemy

$$tf(x) + (1 - t)f(y) \geq f(z) + \langle \nabla f(z), tx + (1 - t)y - z \rangle = f(z).$$

Silna wypukłość

Definicja

Funkcję f klasy C^1 nazywamy m -silnie wypukłą, dla $m \geq 0$ wtedy i tylko wtedy gdy

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{m}{2} \|x - y\|^2.$$

Uwagi:

- Równoważnie silną wypukłość można zdefiniować przez

$$f(tx + (1 - t)y) \leq tf(x) + (1 - t)f(y) - \frac{1}{2}mt(1 - t)\|x - y\|^2$$

- f jest m silnie wypukła wtedy i tylko wtedy gdy $f(x) - \frac{m}{2}\|x\|^2$ jest wypukła.
- Gdy $m = 0$ silna wypukłość oznacza zwyczajną wypukłość.

L-gładkość

Definicja

Funkcje f nazywamy L -gładką jeżeli jej gradient ∇f jest Lipschitzowski ze stałą L :

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|.$$

Stwierdzenie

Jeżeli f jest L -gładka to dla dowolnego \mathbf{x}, \mathbf{y} zachodzi

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2}\|\mathbf{y} - \mathbf{x}\|^2. \quad (1)$$

Dowód

Zdefiniujmy $g(t) = f(x + t(y - x))$. Wówczas $g'(t) = \langle \nabla f(x + t(y - x)), y - x \rangle$ oraz

$$f(y) - f(x) = g(1) - g(0) = \int_0^1 g'(t) dt = \int_0^1 \langle \nabla f(x + t(y - x)), y - x \rangle dt.$$

Stąd

$$\begin{aligned} f(y) - f(x) - \langle \nabla f(x), y - x \rangle &= \int_0^1 \langle \nabla f(x + t(y - x)), y - x \rangle - \langle \nabla f(x), y - x \rangle dt \\ &= \int_0^1 \langle \nabla f(x + t(y - x)) - \nabla f(x), y - x \rangle dt \\ &\leq \|y - x\| \int_0^1 \|\nabla f(x + t(y - x)) - \nabla f(x)\| dt \\ &\leq L \|y - x\|^2 \int_0^1 t dt = \frac{L}{2} \|y - x\|^2 \end{aligned}$$

Monotoniczność spadku po gradiencie

$$\mathbf{x}_k = \mathbf{x}_{k-1} - \gamma_k \nabla f(\mathbf{x}_{k-1})$$

Lemat (1)

Jeśli nierówność (1) jest spełniona dla $\mathbf{y} = \mathbf{x}_k$ $\mathbf{x} = \mathbf{x}_{k-1}$ oraz

$L = \frac{1}{\gamma_k}$ to

$$f(\mathbf{x}_k) \leq f(\mathbf{x}_{k-1})$$

Dowód

Zdefiniujmy $Q_\gamma(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2\gamma} \|\mathbf{y} - \mathbf{x}\|^2$.

- $Q_\gamma(\mathbf{x}, \mathbf{x}) = f(\mathbf{x})$
- $\arg \min_{\mathbf{y}} Q_{\gamma_k}(\mathbf{x}_{k-1}, \mathbf{y}) = \mathbf{x}_{k-1} - \gamma_k \nabla f(\mathbf{x}_{k-1}) = \mathbf{x}_k$
- Jeśli $\gamma < L^{-1}$ to

$$Q_\gamma(\mathbf{x}, \mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2 \geq f(\mathbf{y})$$

Czyli

$$f(\mathbf{x}_{k-1}) = Q(\mathbf{x}_{k-1}, \mathbf{x}_{k-1}) \geq Q(\mathbf{x}_{k-1}, \mathbf{x}_k) \geq f(\mathbf{x}_k).$$

Zbieżność spadku po gradientie

TWIERDZENIE

Jeśli f jest wypukła, L -gładka i osiąga minimum w x^* oraz dla pewnego $a > 0$ zachodzi $aL^{-1} \leq \gamma_k$ oraz spełnione są założenia Lematu (1) to dla dowolnego k

$$f(x_k) - f(x^*) \leq \frac{L \|x_0 - x^*\|^2}{2ak}.$$

Ponadto jeśli f jest m -silnie wypukła to

$$\|x_k - x^*\|^2 \leq \left(1 - \frac{ma}{L}\right)^k \|x_0 - x^*\|^2.$$

Dowód I

Lemat

Jeśli f jest m -silnie wypukła, spełnione są założenia Lematu (1) oraz dla pewnego $a > 0$ zachodzi $aL^{-1} \leq \gamma_k$ to dla dowolnego k i dowolnego y

$$2\gamma_k(f(x_k) - f(y)) \leq \left(1 - \frac{ma}{L}\right) \|x_{k-1} - y\|^2 - \|x_k - y\|^2$$

Dowód twierdzenia

Niech x^* punkt w którym f osiąga minimum wtedy $f(x_k) - f(x^*) > 0$ i z lematu otrzymujemy

$$0 \leq \left(1 - \frac{ma}{L}\right) \|x_{k-1} - y\|^2 - \|x_k - y\|^2$$

Czyli

$$\|x_k - x^*\|^2 \leq \left(1 - \frac{ma}{L}\right) \|x_{k-1} - x^*\|^2$$

Dowód II

Iterując tą nierówność otrzymujemy 2 część twierdzenia.

Z lematu otrzymujemy

$$2\gamma_i(f(x_i) - f(x^*)) \leq \|x_{i-1} - x^*\|^2 - \|x_i - x^*\|^2$$

Sumujemy nierówności od $i = 1$ do k

$$\sum_{i=1}^k 2\gamma_i(f(x_i) - f(x^*)) \leq \|x_0 - x^*\|^2 - \|x_k - x^*\|^2$$

Z monotoniczności $f(x_i)$ oraz z $\gamma_i \geq aL^{-1}$ otrzymujemy

$$\frac{2ak}{L}(f(x_k) - f(x^*)) \leq \|x_0 - x^*\|^2 - \|x_k - x^*\|^2$$

Dowód lematu I

Z nierówności (1) mamy dla $L = \frac{1}{\gamma_k}$

$$\begin{aligned}f(x_k) &\leq f(x_{k-1}) + \langle \nabla f(x_{k-1}), x_k - x_{k-1} \rangle + \frac{L}{2} \|x_k - x_{k-1}\|^2 \\&= f(x_{k-1}) - \frac{1}{\gamma_k} \|x_k - x_{k+1}\|^2 + \frac{L}{2} \|x_k - x_{k-1}\|^2 \\&= f(x_{k-1}) - \frac{1}{\gamma_k} \left(1 - \frac{L\gamma_k}{2}\right) \|x_k - x_{k-1}\|^2 \\&= f(x_{k-1}) - \frac{1}{2\gamma_k} \|x_k - x_{k-1}\|^2\end{aligned}$$

Z wypukłości

$$\begin{aligned}f(y) &\geq f(x_{k-1}) + \langle \nabla f(x_k), y - x_{k-1} \rangle + \frac{m}{2} \|x_{k-1} - y\|^2 \\&= f(x_{k-1}) - \frac{1}{\gamma_k} \langle x_k - x_{k-1}, y - x_{k-1} \rangle + \frac{m}{2} \|x_{k-1} - y\|^2\end{aligned}$$

Dowód lematu II

Mnożąc nierówności przez $2\gamma_k$ i odejmując stronami otrzymujemy.

$$\begin{aligned}2\gamma_k(f(x_k) - f(y)) &\leq 2\langle x_k - x_{k-1}, y - x_{k-1} \rangle - \|x_k - x_{k-1}\|^2 - \gamma_k m \|x_{k-1} - y\|^2 \\&= (1 - m\gamma_k)\|x_{k-1} - y\|^2 - \|x_{k-1} - y\|^2 + 2\langle x_k - x_{k-1}, y - x_{k-1} \rangle - \|x_k - x_{k-1}\|^2 \\&= (1 - m\gamma_k)\|x_{k-1} - y\|^2 - \|x_{k-1} - y - (x_{k-1} - x_k)\|^2 \\&= (1 - m\gamma_k)\|x_{k-1} - y\|^2 - \|x_k - y\|^2\end{aligned}$$

Jak wybrać γ_k

- Optymalne jest wybrać $\gamma_k = L$, jednak obliczenie L (dobre oszacowanie) może, być trudne lub niemożliwe.
- Z dowodu wynika, że jeśli $\gamma_k \rightarrow 0$ oraz $\sum_{k=0}^{\infty} \gamma_k = \infty$ to też otrzymamy zbieżność ale w tempie zależnym od tempa zbieżności γ_k .
- Innym rozwiązaniem jest backtracking. Ustalamy γ_0 i wybieramy $0 < \eta < 1$ I wykonujemy algorytm

- 1 $\gamma_k = \gamma_{k-1}$
- 2 $\mathbf{x}_k = \mathbf{x}_{k-1} - \gamma_k \nabla f(\mathbf{x}_{k-1})$
- 3 Jeśli

$$f(\mathbf{x}_k) \leq f(\mathbf{x}_{k-1}) + \langle \nabla f(\mathbf{x}_{k-1}), \mathbf{x}_k - \mathbf{x}_{k-1} \rangle + \frac{1}{2\gamma_k} \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2$$

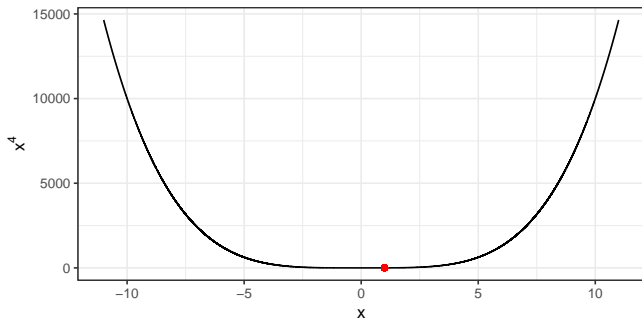
to $k \rightarrow k + 1$. W przeciwnym przypadku $\gamma_k = \eta\gamma_k$ i wracamy do kroku 2.

Istotność założenia o L-gładkości

Jeśli funkcja nie jest L gładka to algorytm gradientowy ze stałym krokiem, a nawet z γ_k dążącym do zera może uciekać do nieskończoności.

Przykład:

$$\arg \min_x x^4, \quad x_0 = 1 \quad \gamma_k = \frac{1}{k}$$

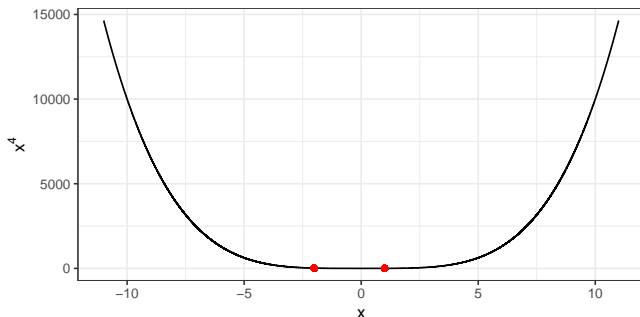


Istotność założenia o L-gładkości

Jeśli funkcja nie jest L gładka to algorytm gradientowy ze stałym krokiem, a nawet z γ_k dążącym do zera może uciekać do nieskończoności.

Przykład:

$$\arg \min_x x^4, \quad x_0 = 1 \quad \gamma_k = \frac{1}{k}$$

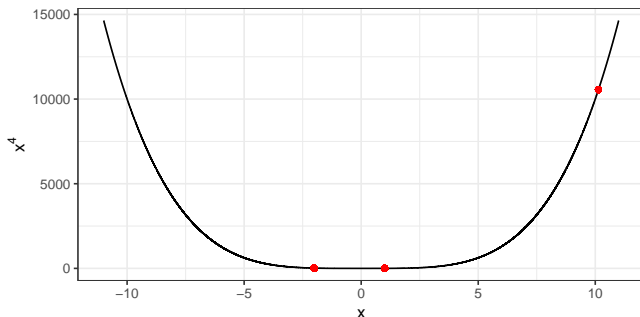


Istotność założenia o L-gładkości

Jeśli funkcja nie jest L-gładka to algorytm gradientowy ze stałym krokiem, a nawet z γ_k dążącym do zera może uciekać do nieskończoności.

Przykład:

$$\arg \min_x x^4, \quad x_0 = 1 \quad \gamma_k = \frac{1}{k}$$



Przyspieszenie Nesterova

Można pokazać, że optymalne tempo zbieżności algorytmów pierwszego rzędu (opartych na gradiencie) jest $\mathcal{O}(1/k^2)$. Spadek po gradiencie zbiega w tempie $\mathcal{O}(1/k)$. Czy można go poprawić?

Nesterov 1983

Niech f wypukła

$$\arg \min_x f(x)$$

- 1 $y_0 = x_0, t_0 = 0$
- 2 $x_k = y_{k-1} - \gamma_k \nabla f(y_{k-1})$
- 3 $t_k = \frac{1 + \sqrt{1 + 4t_{k-1}^2}}{2}$
- 4 $y_k = x_k + \frac{t_{k-1} - 1}{t_k} (x_k - x_{k-1})$

Przyspieszenie Nesterowa

TWIERDZENIE

Niech f wypukła, L -gładka oraz osiąga minimum w x^* . Jeśli $\gamma_k = \gamma \leq L^{-1}$ to

$$f(x_k) - f(x^*) \leq \frac{2\|x_0 - x^*\|^2}{\gamma k^2}.$$

Dowód I

Tak jak w poprzednim dowodzie z wypukłości i L-gładkości

$$\begin{aligned}f(\mathbf{x}_k) &\leq f(\mathbf{y}_{k-1}) - \frac{1}{2\gamma} \|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2 \\f(z) &\geq f(\mathbf{y}_{k-1}) - \frac{1}{\gamma} \langle \mathbf{x}_k - \mathbf{y}_{k-1}, z - \mathbf{y}_{k-1} \rangle\end{aligned}$$

Odejmując stronami uzyskujemy

$$f(\mathbf{x}_k) - f(z) \leq -\frac{1}{2\gamma} \|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2 + \frac{1}{\gamma} \langle \mathbf{x}_k - \mathbf{y}_{k-1}, z - \mathbf{y}_{k-1} \rangle$$

Używamy tej nierówności z $z = \mathbf{x}_{k-1}$ i $z = \mathbf{x}^*$

$$\begin{aligned}f(\mathbf{x}_k) - f(\mathbf{x}_{k-1}) &\leq -\frac{1}{2\gamma} \|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2 + \frac{1}{\gamma} \langle \mathbf{x}_k - \mathbf{y}_{k-1}, \mathbf{x}_{k-1} - \mathbf{y}_{k-1} \rangle \\f(\mathbf{x}_k) - f(\mathbf{x}^*) &\leq -\frac{1}{2\gamma} \|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2 + \frac{1}{\gamma} \langle \mathbf{x}_k - \mathbf{y}_{k-1}, \mathbf{x}^* - \mathbf{y}_{k-1} \rangle\end{aligned}$$

Dowód II

Mnożymy pierwszą nierówność przez $t_{k-1} - 1$ i dodajemy do drugiej

$$t_{k-1}(f(x_k) - f(x^*)) - (t_{k-1} - 1)(f(x_{k-1}) - f(x^*)) \leq \\ - \frac{t_{k-1}}{2\gamma} \|x_k - y_{k-1}\|^2 + \frac{1}{\gamma} \langle x_k - y_{k-1}, (t_{k-1} - 1)x_{k-1} - t_{k-1}y_{k-1} + x^* \rangle$$

Oznaczmy przez $\Delta_k = f(x_k) - f(x^*)$ i pomnóżmy nierówność przez t_{k-1} ($t_{k-2}^2 = t_{k-1}^2 - t_{k-1}$)

$$t_{k-1}^2 \Delta_k - t_{k-2}^2 \Delta_{k-1} \leq \frac{-1}{2\gamma} \left[\|(x_k - y_{k-1})t_{k-1}\|^2 \right. \\ \left. + 2t_{k-1} \langle x_k - y_{k-1}, t_{k-1}y_{k-1} - (t_{k-1} - 1)x_{k-1} - x^* \rangle \right]$$

Dowód III

Po przekształceniu

$$t_{k-1}^2 \Delta_k - t_{k-2}^2 \Delta_{k-1} \leq \frac{-1}{2\gamma}$$

$$\left[\|x_k t_{k-1} - (t_{k-1} - 1)x_{k-1} - x^*\|^2 - \|t_{k-1}y_{k-1} - (t_{k-1} - 1)x_{k-1} - x^*\|^2 \right]$$

Dodatkowo

$$y_k = x_k + \frac{t_{k-1} - 1}{t_k} (x_k - x_{k-1})$$

$$t_k y_k = t_k x_k + (t_{k-1} - 1)(x_k - x_{k-1})$$

$$t_k y_k - (t_k - 1)x_k = t_{k-1}x_k - (t_{k-1} - 1)x_{k-1}$$

Oznaczając przez $u_k = t_{k-1}x_k - (t_{k-1} - 1)x_{k-1} - x^*$

Stąd

$$t_{k-1}^2 \Delta_k - t_{k-2}^2 \Delta_{k-1} \leq \frac{1}{2\gamma} \left(\|u_{k-1}\|^2 - \|u_k\|^2 \right)$$

Dowód IV

Sumując nierówność od $i = 2$ do k otrzymujemy

$$t_{k-1}^2 \Delta_k \leq \frac{1}{2\gamma} (\|u_1\|^2 - \|u_k\|^2) \leq \frac{1}{2\gamma} \|u_1\|^2$$

Gdzie

$$u_1 = t_0 x_k - (t_0 - 1)x_0 - x^* = x_0 - x^*$$

Oraz dla $k \geq 2$ (Indukcja)

$$t_{k-1} \geq \frac{k}{2}$$

Backtracking

- 1 $y_0 = x_0, t_0 = 0, \gamma_0$ i $0 < \eta < 1$
- 2 $\gamma_k = \gamma_{k-1}$
- 3 $x_k = y_{k-1} - \gamma_k \nabla f(y_{k-1})$
- 4 Jeśli

$$f(x_k) \leq f(y_{k-1}) + \langle \nabla f(y_{k-1}), x_k - y_{k-1} \rangle + \frac{1}{2\gamma_k} \|x_k - y_{k-1}\|^2$$

to krok 5. W przeciwnym przypadku $\gamma_k = \eta\gamma_k$ i wracamy do kroku 3.

- 5 $t_k = \frac{1 + \sqrt{1 + 4t_{k-1}^2}}{2}$
- 6 $y_k = x_k + \frac{t_{k-1} - 1}{t_k} (x_k - x_{k-1})$

TWIERDZENIE

Niech f wypukła, L -gładka oraz osiąga minimum w x^* . Jeśli γ_k jak powyżej to

$$f(x_k) - f(x^*) \leq \frac{2L \|x_0 - x^*\|^2}{\eta k^2}.$$

1 Optymalizacja

- Wstęp
- Metody gradientowe
- Algorytm Newtona-Raphsona
- Iterowane ważone najmniejsze kwadraty
- Proximal gradient
- ADMM
- Expectation Maximization

Algorytm Newtona - Raphsona

Niech f klasy C^2 , rozwiązanie problemu

$$\arg \min_x f(x)$$

srowadza się do rozwiązania równania

$$\nabla f(x) = 0 \tag{2}$$

Ze wzoru Taylora

$$\nabla f(x) \approx \nabla f(x_0) + \nabla^2 f(x_0)(x - x_0)$$

Czyli możemy przybliżać rozwiązanie (2) poprzez

$$x = x_0 - (\nabla^2 f(x_0))^{-1} \nabla f(x_0)$$

Algorytm Newtona-Raphsona

- 1 Ustalamy punkt startowy x_0 .
- 2 Kolejne przybliżenia rozwiązania uzyskujemy

$$x_k = x_{k-1} - (\nabla^2 f(x_{k-1}))^{-1} \nabla f(x_{k-1})$$

Kwadratowa zbieżność metody Newtona-Raphsona

Algorytm jest kwadratowo zbieżny do x^* jeżeli istnieje M takie, że

$$\lim_{k \rightarrow \infty} \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|^2} < M$$

TWIERDZENIE

Niech f funkcja jest klasy C^3 i osiąga minimum w x^* . Jeżeli istnieje $\epsilon > 0$, $M < \infty$, takie że

$$\sup_{\{x, y: \|x - x^*\| < \epsilon, \|y - x^*\| < \epsilon\}} \|(\nabla^2 f(x))^{-1} \nabla^3 f(y)\| < M$$

oraz $\|x_0 - x^*\| < \min(\epsilon, M^{-1})$, $\nabla^2 f$ odwracalna, to dla każdego k

$$\frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|^2} < M$$

Dowód

Ponieważ w \mathbf{x}^* jest minimum f ze wzoru Taylora dla pewnego ξ “pomiędzy” \mathbf{x}_k i \mathbf{x}^* otrzymujemy

$$0 = \nabla f(\mathbf{x}^*) = \nabla f(\mathbf{x}_k) + \nabla^2 f(\mathbf{x}_k)(\mathbf{x}^* - \mathbf{x}_k) + \nabla^3 f(\xi)(\mathbf{x}^* - \mathbf{x}_k)^{\otimes 2}$$

Mnożąc, stronami przez $(\nabla^2 f(\mathbf{x}_k))^{-1}$ otrzymujemy

$$0 = (\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k) + (\mathbf{x}^* - \mathbf{x}_k) + (\nabla^2 f(\mathbf{x}_k))^{-1} \nabla^3 f(\xi)(\mathbf{x}^* - \mathbf{x}_k)^{\otimes 2}$$

Z definicji \mathbf{x}_{k+1} otrzymujemy

$$0 = \mathbf{x}^* - \mathbf{x}_{k+1} + (\nabla^2 f(\mathbf{x}_k))^{-1} \nabla^3 f(\xi)(\mathbf{x}^* - \mathbf{x}_k)^{\otimes 2}$$

Przenosząc na drugą stronę i przykładając normę otrzymujemy

$$\|\mathbf{x}^* - \mathbf{x}_{k+1}\| \leq \|(\nabla^2 f(\mathbf{x}_k))^{-1} \nabla^3 f(\xi)\| \|\mathbf{x}^* - \mathbf{x}_k\|^2$$

Własności metody Newtona-Raphsona

- Przy dodatkowych warunkach regularności można pokazać globalną kwadratową zbieżność.
- W ogólności obszar kwadratowej zbieżności może być mały i przy złym wyborze punktu startowego algorytm może wpaść w cykl lub nawet być rozbieżny.
- W przestrzeni wyosko wymiarowej, problemem może być obliczanie macierzy drugiej pochodnej i jej odwracanie.

1 Optymalizacja

- Wstęp
- Metody gradientowe
- Algorytm Newtona-Raphsona
- Iterowane ważone najmniejsze kwadraty
- Proximal gradient
- ADMM
- Expectation Maximization

Regresja logistyczna

Y_1, Y_2, \dots, Y_n - binarne zmienne objaśniane, x_1, x_2, \dots, x_n - zmienne objaśniające z \mathbb{R}^p . Oznaczmy przez $p_i = \mathbb{P}(Y_i = 1)$

Regresja logistyczna

$$\text{logit}(p_i) = \log\left(\frac{p_i}{1 - p_i}\right) = x_i^T \beta$$

Log wiarygodność

$$\ell(\beta) = \sum_{i=1}^n Y_i x_i^T \beta - \sum_{i=1}^n \log(1 + \exp(x_i^T \beta))$$

Estymator największej wiarygodności

$$\hat{\beta} = \arg \max_{\beta} \ell(\beta) = \arg \min_{\beta} -\ell(\beta)$$

Jak obliczyć $\hat{\beta}$?

Uogólnione modele liniowe

Regresja logistyczna jest szczególnym przypadkiem uogólnionych modeli liniowych. Rodzina wykładnicza

$$f_{\theta}(y) = \exp \left[\frac{\theta y - b(\theta)}{a(\phi)} \right] h(y)$$

- θ - parametr kanoniczny
- ϕ - parametr dyspersji
- $\mathbb{E}(Y) = b'(\theta) = \mu$
- $\text{var}(Y) = a(\phi)b''(\theta)$

Uogólniony model liniowy $Y_1, \dots, Y_n \sim f_{\theta_i}$

$$g(\mu_i) = \mathbf{x}_i^T \beta = \eta_i$$

g - funkcja łącząca (jeśli $g(\mu_i) = \theta_i$ to kanoniczna funkcja łącząca)

Estymator największej wiarygodności w GLM I

$$\ell(\beta) = \sum_{i=1}^n \frac{\theta_i y_i - b(\theta_i)}{a(\phi)}$$

Chcemy wyprowadzić algorytm Newtona-Raphsona (przybliżony).

$$\frac{\partial}{\partial \beta_j} \ell =? \quad \frac{\partial^2}{\partial \beta_j \partial \beta_k} \ell =?$$

Dla uproszczenia zapisu przyjmijmy na chwilę $n = 1$ i pomińmy indeks i

$$\frac{\partial}{\partial \beta_j} \ell(\beta) = \frac{\partial \ell}{\partial \theta} \frac{\partial \theta}{\partial \mu} \frac{\partial \mu}{\partial \eta} \frac{\partial \eta}{\partial \beta_j}$$

Estymator największej wiarygodności w GLM II

$$\begin{aligned}\frac{\partial \ell}{\partial \theta} &= \frac{y - b'(\theta)}{a(\phi)} \\ \frac{\partial \theta}{\partial \mu} &= \frac{1}{b''(\theta)} = \frac{a(\phi)}{\text{var}(Y)} \\ \frac{\partial \eta}{\partial \beta_j} &= x_{ij}\end{aligned}$$

W przypadku kanonicznej funkcji łączącej $\eta = (b')^{-1}(\mu)$

$$\frac{\partial \mu}{\partial \eta} = b''(\theta)$$

Stąd

$$\frac{\partial \ell}{\partial \beta_j} = \frac{y - \mu}{\text{var}(Y)} x_{ij} \frac{\partial \mu}{\partial \eta}$$

Estymator największej wiarygodności w GLM III

i w przypadku kanonicznej funkcji łączącej

$$\frac{\partial \ell}{\partial \beta_j} = \frac{y - \mu}{a(\phi)} x_{ij}$$

Fisher scoring

W algorytmie Newtona-Raphsona chcemy przybliżyć $\frac{\partial^2 \ell}{\partial \beta_j \partial \beta_k}$ przez $\mathbb{E} \frac{\partial^2}{\partial \beta_j \partial \beta_k} \ell$. Czyli

$$\nabla^2 \ell(\beta) \approx -I(\beta)$$

$I(\beta)$ - informacja Fishera

W przypadku kanonicznej funkcji łączącej

$$\nabla^2 \ell(\beta) = -I(\beta)$$

Estymator największej wiarygodności w GLM IV

Korzystając z faktu

$$\begin{aligned}\mathbb{E} \frac{\partial^2 \ell}{\partial \beta_j \partial \beta_k} &= -\mathbb{E} \frac{\partial \ell}{\partial \beta_j} \frac{\partial \ell}{\partial \beta_k} \\ &= \mathbb{E} \frac{(y - \mu)^2}{\text{var}(Y)^2} x_{ij} x_{ik} \left(\frac{\partial \mu}{\partial \eta} \right)^2 = \frac{x_{ij} x_{ik}}{\text{var}(Y)} \left(\frac{\partial \mu}{\partial \eta} \right)^2\end{aligned}$$

Oznaczmy

$$A = \text{diag} \left(\frac{\frac{\partial \mu_i}{\partial \eta_i}}{\text{var}(Y_i)} \right) \quad W = \text{diag} \left(\frac{\left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2}{\text{var}(Y_i)} \right)$$

Otrzymujemy

$$\nabla \ell(\beta) = X^T A (y - \mu), \quad I(\beta) = X^T W X$$

Estymator największej wiarygodności w GLM V

Oraz algorytm Fisher scoring przyjmuje postać

$$\begin{aligned}\beta^{k+1} &= \beta^k + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{A}(y - \mu) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} [\mathbf{X}^T \mathbf{W} \mathbf{X} \beta^k + \mathbf{X}^T \mathbf{A}(y - \mu)]\end{aligned}$$

$$\mathbf{A} \operatorname{diag} \left(\frac{\partial \mu}{\partial \eta} \right) = \mathbf{W}, \quad \mathbf{X}^T \beta = \eta$$

Otrzymujemy iterowane ważone najmniejsze kwadraty

$$\beta^{k+1} = (\mathbf{X}^T \mathbf{W}^k \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}^k \mathbf{z}^k,$$

gdzie

$$\mathbf{z} = \eta - \operatorname{diag} \left(\frac{\partial \eta}{\partial \mu} \right) (y - \mu)$$

To jeszcze nie jest praktyczny algorytm, pozostaje pytanie jak numerycznie rozwiązać równanie liniowe

Powrót do regresji logistycznej

$$Y_i \sim \text{Bin}(1, p_i)$$

$$\text{logit}(p_i) = \mathbf{x}_i^T \boldsymbol{\beta}, \quad p_i = \frac{1}{1 + \exp(-\mathbf{x}_i^T \boldsymbol{\beta})}$$

Gęstość można zapisać w postaci wykładniczej

$$f(y; p) = \exp(y \text{logit}(p) + \log(1 - p))$$

kanoniczny parametr $\theta = \text{logit}(p)$

$$f(y; \theta) = \exp(y\theta - \log(1 + \exp(\theta)))$$

$$\mathbb{E}(Y) = p, \quad \text{var}(Y) = p(1 - p), \quad \eta = \theta, \quad \frac{\partial \eta}{\partial \mu} = \frac{1}{p(1 - p)}$$

Stąd W i z są postaci

$$W = \text{diag}(p_i(1 - p_i)), \quad z_i = \mathbf{x}_i^T \boldsymbol{\beta} + \frac{y_i - p_i}{p_i(1 - p_i)}$$

Obliczanie rozwiązania MNK

Rozkład Choleskiego Dodatnio określoną macierz A można przedstawić w postaci

$$A = LL^T$$

gdzie L dolnotrójkątna.

Równanie

$$X^T X \beta = X^T y$$

Rozwiązujemy

① $X^T X = LL^T$

②

$$Lz = X^T Y$$

$$L^T \beta = z$$

Obliczanie rozwiązania MNK II sposób

Rozkład QR Macierz X $n \times p$, $n > p$ można przedstawić w postaci

$$X = QR$$

gdzie Q - $n \times p$ o ortogonalnych kolumnach $Q^T Q = I$ i R górnotrójkątna $p \times p$. Wówczas

$$\beta = (X^T X)^{-1} X^T Y = (R^T Q^T Q R)^{-1} R^T Q^T Y = R^{-1} Q^T Y$$

I dostajemy równanie

$$R\beta = Q^T Y$$

Dodatkowo

$$\hat{Y} = X\hat{\beta} = Q Q^T Y$$

1 Optymalizacja

- Wstęp
- Metody gradientowe
- Algorytm Newtona-Raphsona
- Iterowane ważone najmniejsze kwadraty
- Proximal gradient
- ADMM
- Expectation Maximization

LASSO

Model liniowy:

$$Y = X\beta + \epsilon$$

Metoda najmniejszych kwadratów:

$$\arg \min_{\beta} \|Y - X\beta\|^2$$

LASSO:

$$\arg \min_{\beta} \{ \|Y - X\beta\|^2 + \lambda \sum_{i=1}^p |\beta_i| \}$$

Jak obliczyć rozwiązanie LASSO?

Subgradient

Jeśli f jest wypukła i gładka to

$$f(y) \geq f(x) + \nabla f(x)^T (y - x)$$

Definicja

Jeśli f jest wypukła to g nazywamy subgradientem f w punkcie x wtedy i tylko wtedy gdy dla każdego y

$$f(y) \geq f(x) + g^T (y - x)$$

Zbiór wszystkich subgradientów f w punkcie x nazywamy subróżniczką i będziemy oznaczać przez $\partial f(x)$

Własności subrózniczki

Jeśli f wypukła to:

- $\partial f(x)$ jest niepustym wypukłym zbiorem
- Jeśli $\partial f(x) = \{g\}$ wtedy i tylko wtedy gdy f różniczkowalna w x oraz $g = \nabla f(x)$
- f osiąga minimum w x^* wtedy i tylko wtedy gdy $0 \in \partial f(x^*)$
- $\partial(f_1 + f_2) = \partial f_1 + \partial f_2$

Spadek po subgradinecie

$$\mathbf{x}_k = \mathbf{x}_{k-1} - \gamma_k \mathbf{g}_{k-1}, \quad \mathbf{g}_{k-1} \in \partial f(\mathbf{x}_{k-1})$$

Własności: Jeśli f G -Lipshitzowska oraz $\|\mathbf{x}_0 - \mathbf{x}^*\| < R$ to

- Dla $\gamma_k = \gamma$

$$\lim_{k \rightarrow \infty} f(\mathbf{x}_k^N) \leq f(\mathbf{x}^*) + G^2 \gamma$$

- Jeśli $\gamma_k \rightarrow 0, \sum \gamma_k^2 < \infty$ oraz $\sum \gamma_k = \infty$ to

$$\lim_{k \rightarrow \infty} f(\mathbf{x}_k^N) = f(\mathbf{x}^*)$$

- Jeśli $\gamma_k = \frac{R}{G\sqrt{k}}$

$$f(\mathbf{x}_k^N) - f(\mathbf{x}^*) \leq \frac{RG}{\sqrt{k}}$$

Proximal operator

$$\text{prox}_{\gamma, f}(\mathbf{x}) = \arg \min_y \left\{ f(y) + \frac{1}{2\gamma} \|\mathbf{x} - y\|^2 \right\}$$

Dla dowolnej wypukłej funkcji f jeżeli

$$y = \text{prox}_{\gamma, f}(\mathbf{x})$$

to

$$y \in \mathbf{x} - \gamma \partial f(y)$$

Przykład:

$$\text{prox}_{\gamma, \lambda|\cdot|} = \text{sign}(\mathbf{x})(|\mathbf{x}| - \lambda\gamma)_+$$

Proximal gradient

Problem:

$$\arg \min_{\mathbf{x}} \{f(\mathbf{x}) + g(\mathbf{x})\}$$

- f wypukła i gładka
- g wypukła

Algorytm

$$\mathbf{x}_{k+1} = \text{prox}_{\gamma_k, g}(\mathbf{x}_k - \gamma_k \nabla f(\mathbf{x}_k))$$

Stąd

$$\mathbf{x}_{k+1} \in \mathbf{x}_k - \gamma_k \nabla f(\mathbf{x}_k) - \gamma_k \partial g(\mathbf{x}_{k+1})$$

Własności proximal gradient

Jeśli f jest L -gładka oraz $\gamma \leq L^{-1}$ to

$$f(x) + g(x) \leq f(y) + \nabla f(y)^T(x - y) + \frac{1}{2\gamma} \|y - x\|^2 + g(x) = Q_\gamma(x, y)$$

Przekształcając

$$Q_\gamma(x, y) = f(y) - \frac{\gamma}{2} \|\nabla f(y)\|^2 + g(x) + \frac{1}{2\gamma} \|x - (y - \gamma \nabla f(y))\|^2$$

Czyli

$$\text{prox}_{\gamma_k, g}(x_k - \gamma_k \nabla f(x_k)) = \arg \min_x Q_{\gamma_k}(x, x_k)$$

Zatem

$$\begin{aligned} f(x_k) + g(x_k) &= Q_{\gamma_k}(x_k, x_k) \geq \arg \min_x Q_{\gamma_k}(x, x_k) \\ &= Q_{\gamma_k}(x_{k+1}, x_k) \geq f(x_{k+1}) + g(x_{k+1}) \end{aligned}$$

Zbieżność proximal gradient

Oznaczmy przez $F = f + g$

Lemat

Jeśli f jest wypukła i L -gładka, g jest wypukła oraz $\gamma_k \leq L^{-1}$ to dla dowolnego y

$$2\gamma_k(F(x_k) - F(y)) \leq \|x_{k-1} - y\|^2 - \|x_k - y\|^2$$

Dowód:

Niech $v \in \partial g(x_k)$ takie, że

$$x_k = x_{k-1} - \gamma_k \nabla f(x_{k-1}) - \gamma_k v$$

Mamy trzy nierówności

$$F(x_k) \leq g(x_k) + f(x_{k-1}) + \nabla f(x_{k-1})^T (x_k - x_{k-1}) + \frac{1}{2\gamma_k} \|x_k - x_{k-1}\|^2$$

$$f(y) \geq f(x_{k-1}) + \nabla f(x_{k-1})^T (y - x_{k-1})$$

$$g(y) \geq g(x_k) + v^T (y - x_k)$$

Dowód lematu c.d.

Odejmując nierówności otrzymujemy

$$\begin{aligned} F(\mathbf{x}_k) - F(y) &\leq \frac{1}{2\gamma_k} \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 + \nabla f(\mathbf{x}_{k-1})^T (\mathbf{x}_k - \mathbf{x}_{k-1}) \\ &\quad - \nabla f(\mathbf{x}_{k-1})^T (y - \mathbf{x}_{k-1}) - v^T (y - \mathbf{x}_k) \end{aligned}$$

Przekształcając prawą stronę otrzymujemy

$$F(\mathbf{x}_k) - F(y) \leq \frac{1}{2\gamma_k} \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 + (\nabla f(\mathbf{x}_{k-1}) + v)^T (\mathbf{x}_k - y)$$

$$F(\mathbf{x}_k) - F(y) \leq \frac{1}{2\gamma_k} \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 - \frac{1}{\gamma_k} (\mathbf{x}_k - \mathbf{x}_{k-1})^T (\mathbf{x}_k - y)$$

$$2\gamma_k (F(\mathbf{x}_k) - F(y)) \leq \|\mathbf{x}_k - \mathbf{x}_{k-1} - \mathbf{x}_k + y\|^2 - \|\mathbf{x}_k - y\|^2$$

Zbieżność proximal gradient

TWIERDZENIE

Niech $F = f + g$, gdzie f jest wypukła i L -gładka, g jest wypukła oraz F osiąga minimum w x^* . Jeśli $\gamma_k = \gamma < L^{-1}$ to

$$F(x_k) - F(x^*) \leq \frac{L \|x_0 - x^*\|^2}{2k}.$$

Przyspieszenie Nesterova dla proximal gradient

- 1 $y_0 = x_0, t_0 = 0, \gamma_0$ i $0 < \eta < 1$
- 2 $\gamma_k = \gamma_{k-1}$
- 3 $x_k = \text{prox}_{\gamma_k g}(y_{k-1} - \gamma_k \nabla f(y_{k-1}))$
- 4 Jeśli

$$f(x_k) \leq f(y_{k-1}) + \langle \nabla f(y_{k-1}), x_k - y_{k-1} \rangle + \frac{1}{2\gamma_k} \|x_k - y_{k-1}\|^2$$

to krok 5. W przeciwnym przypadku $\gamma_k = \eta\gamma_k$ i wracamy do kroku 3.

- 5 $t_k = \frac{1 + \sqrt{1 + 4t_{k-1}^2}}{2}$
- 6 $y_k = x_k + \frac{t_{k-1} - 1}{t_k} (x_k - x_{k-1})$

TWIERDZENIE

Niech f wypukła, L -gładka, g wypukła oraz F osiąga minimum w x^* . Jeśli γ_k jak powyżej to

$$F(x_k) - F(x^*) \leq \frac{2L \|x_0 - x^*\|^2}{\eta k^2}.$$

1 Optymalizacja

- Wstęp
- Metody gradientowe
- Algorytm Newtona-Raphsona
- Iterowane ważone najmniejsze kwadraty
- Proximal gradient
- ADMM
- Expectation Maximization

Alternating Direction Method of Multipliers

Rozważmy problem

$$\min_{x: Ax=b} f(x)$$

Lagrangian dla tego problemu

$$\mathcal{L}(x, \mu) = f(x) + \mu^T (Ax - b)$$

Problem dualny

$$\max_{\mu} g(\mu), \quad g(\mu) = \inf_x \mathcal{L}(x, \mu)$$

Rozwiązanie

$$x^* = \arg \min_x \mathcal{L}(x, \mu^*)$$

Algorytm gradientowy dla problemu dualnego

$$\mu_{k+1} = \mu_k + \gamma_k \nabla g(\mu_k),$$

gdzie

$$\nabla g(\mu_k) = A\tilde{x} - b, \quad \tilde{x} = \arg \min \mathcal{L}(x, \mu_k).$$

Czyli

$$\begin{aligned} x_{k+1} &= \arg \min \mathcal{L}(x, \mu_k) \\ \mu_{k+1} &= \mu_k + \gamma_k (Ax_{k+1} - b) \end{aligned}$$

Algorytm jest zbieżny przy restrykcyjnych założeniach. Nie praktyczny

Metoda mnożników I

Warunki konieczne

$$Ax^* - b = 0, \quad \nabla \mathcal{L}(x^*, \mu^*) = \nabla f(x^*) + A^T \mu^* = 0$$

Chcielibyśmy aby

$$\nabla \mathcal{L}(x_k, \mu_k) = 0$$

Równoważnie, warunki konieczne można zapisać w postaci

$$Ax^* - b = 0, \quad \nabla f(x^*) + A^T \mu^* + \rho A^T (Ax^* - b) = 0$$

Co odpowiada

$$Ax^* - b = 0, \quad \nabla \mathcal{L}_\rho(x^*, \mu^*) = 0,$$

gdzie

$$\mathcal{L}_\rho(x, \mu) = f(x) + \mu^T (Ax - b) + \frac{\rho}{2} \|Ax - b\|^2$$

Metoda mnożników II

Algorytm:

$$\begin{aligned}x_{k+1} &= \arg \min_x \mathcal{L}_\rho(x, \mu_k) \\ \mu_{k+1} &= \mu_k + \rho(Ax_{k+1} - b)\end{aligned}$$

Mamy

$$\begin{aligned}0 &= \nabla \mathcal{L}_\rho(x_{k+1}, \mu_k) \\ &= \nabla f(x_{k+1}) + A^T \mu_k + \rho A^T (Ax_{k+1} - b) \\ &= \nabla f(x_{k+1}) + A^T (\mu_k + \rho(Ax_{k+1} - b)) \\ &= \nabla f(x_{k+1}) + A^T \mu_{k+1}\end{aligned}$$

II warunek osiągnany jest asymptotycznie

$$Ax_k - b \rightarrow 0$$

ADMM

Rozważmy problem

$$\min_{\{x,z: Ax+Bz=c\}} f(x) + g(z)$$

f, g wypukłe

$$\mathcal{L}_\rho(x, z, \mu) = f(x) + g(z) + \mu^T(Ax + Bz - c) + \frac{\rho}{2}\|Ax + Bz - c\|^2$$

Algorytm

- 1 $x_{k+1} = \arg \min_x \mathcal{L}_\rho(x, z_k, \mu_k)$
- 2 $z_{k+1} = \arg \min_z \mathcal{L}_\rho(x_{k+1}, z, \mu_k)$
- 3 $\mu_{k+1} = \mu_k + \rho(Ax_{k+1} + Bz_{k+1} - c)$

Uwagi

- Kroki 2 i 3 można zastąpić przez rozwiązanie numeryczne.
- Algorytm jest zbieżny jeżeli \mathcal{L}_0 ma punkt przegięcia.

ADMM dla LASSO

Problem

$$\min_{\beta} \frac{1}{2} \|Y - X\beta\|^2 + \lambda \|\beta\|_1$$

można zapisać jako

$$\min_{\{\beta, z: \beta=z\}} \frac{1}{2} \|Y - X\beta\|^2 + \lambda \|z\|_1$$

oraz

$$\mathcal{L}_{\rho}(\beta, z, \mu) = \frac{1}{2} \|Y - X\beta\|^2 + \lambda \|z\|_1 + \mu^T(\beta - z) + \frac{\rho}{2} \|\beta - z\|^2$$

ADMM

- 1 $\beta_{k+1} = (X^T X + \rho I)^{-1} (X^T Y - \mu_k + \rho z_k)$
- 2 $z_{k+1} = \text{prox}_{\frac{\lambda}{\rho} \|\cdot\|_1}(\beta_{k+1} + \frac{1}{\rho} \mu_k) =$
 $\text{sign}(\beta_{k+1} + \frac{1}{\rho} \mu_k) (|\beta_{k+1} + \frac{1}{\rho} \mu_k| - \frac{\lambda}{\rho})_+$
- 3 $\mu_{k+1} = \mu_k + \rho(\beta_{k+1} - z_{k+1})$

ADMM dla optymalizacji z ograniczeniami wypukłymi I

Problem:

$$\min_{x \in \mathcal{K}} f(x)$$

f - wypukła, \mathcal{K} domknięty wypukły podzbiór \mathbb{R}^d .

$$\min_{x=z} f(x) + 1_{\mathcal{K}}(z)$$

gdzie

$$1_{\mathcal{K}}(z) = \begin{cases} 1 & z \in \mathcal{K} \\ +\infty & z \notin \mathcal{K} \end{cases}$$

$$\mathcal{L}_{\rho}(x, z, \mu) = f(x) + 1_{\mathcal{K}}(z) + \mu^T(x - z) + \frac{\rho}{2} \|x - z\|^2$$

ADMM

- 1 $x_{k+1} = \arg \min_x \mathcal{L}_{\rho}(x, z_k, \mu_k)$
- 2 $z_{k+1} = \Pi_{\mathcal{K}}(x_{k+1} + \frac{1}{\rho} \mu_k)$
- 3 $\mu_{k+1} = \mu_k + \rho(x_{k+1} - z_{k+1})$

ADMM dla optymalizacji z ograniczeniami wypukłymi II

Jeżeli f funkcja kwadratowa czyli

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T\mathbf{P}\mathbf{x} + \mathbf{q}^T\mathbf{x} + r$$

to

$$\arg \min_{\mathbf{x}} \mathcal{L}_{\rho}(\mathbf{x}, \mathbf{z}_k, \mu_k) = (\mathbf{P} + \rho\mathbf{I})^{-1}(\rho\mathbf{z}_k - \mathbf{q} - \mu_k)$$

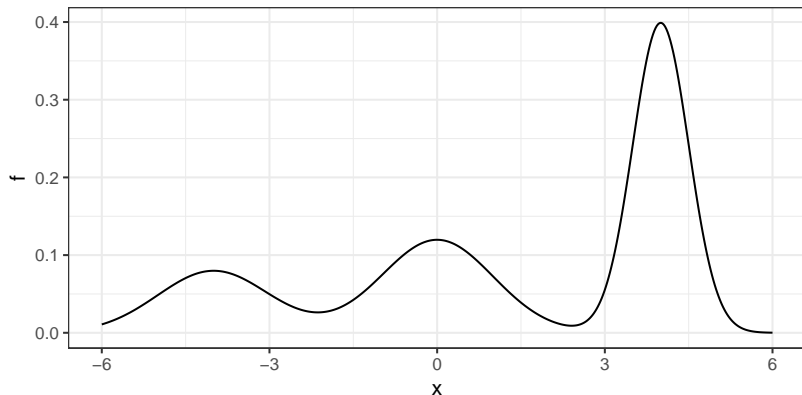
1 Optymalizacja

- Wstęp
- Metody gradientowe
- Algorytm Newtona-Raphsona
- Iterowane ważone najmniejsze kwadraty
- Proximal gradient
- ADMM
- Expectation Maximization

Mieszanka rozkładów gaussowskich

Niech Y_1, \dots, Y_n będą zmiennymi losowymi o rozkładzie

$$f_{\theta}(y) = \sum_{i=1}^k \alpha_i \phi_{\mu_i, \sigma_i}(y), \quad \theta = (\alpha_i, \mu_i, \sigma_i)_{i=1}^k$$



Jak obliczyć estymator θ ?

Mieszanka jako model z brakującymi danymi

Rozważmy parę (X, Y) o rozkładzie

$$P(X = i) = \alpha_i, \quad Y|X = i \sim \mathcal{N}(\mu_i, \sigma_i^2)$$

Wówczas rozkład brzegowy Y ma gęstość

$$f_{\theta}(y) = \sum_{i=1}^k \alpha_i \phi_{\mu_i, \sigma_i}(y)$$

Expectation Maximization I

Rozkład łączny

$$f_{\theta}(x, y)$$

Obserwujemy Y z rozkładu brzegowego

$$f_{\theta}(y) = \int f_{\theta}(x, y) dx$$

Chcemy policzyć MLE

$$\arg \max_{\theta} \ell(\theta), \quad \ell(\theta) = \log(f_{\theta}(y))$$

Wprowadźmy funkcję pomocniczą

$$Q(\theta, \vartheta) = \int \log(f_{\theta}(x, y)) f_{\vartheta}(x|y) dx = \mathbb{E}_{\vartheta}(\log(f_{\theta}(X, Y)) | Y)$$

Expectation Maximization II

Krok Expectation:

$$Q(\theta, \theta_k) = \mathbb{E}_{\theta_k}(\log(f_\theta(X, Y))|Y)$$

Krok Maximization:

$$\theta_{k+1} = \arg \max_{\theta} Q(\theta, \theta_k)$$

Własności funkcji Q

Oznaczmy przez

$$H(\theta, \vartheta) = - \int \log(f_{\theta}(x|y)) f_{\vartheta}(x|y) dx$$

- 1 $Q(\theta, \vartheta) = \ell(\theta) - H(\theta, \vartheta)$
- 2 $\ell(\theta) - \ell(\vartheta) \geq Q(\theta, \vartheta) - Q(\vartheta, \vartheta)$
- 3 $\nabla_{\theta} Q(\theta, \vartheta)|_{\theta=\vartheta} = \nabla \ell(\vartheta)$

Powyższe własności wynikają z faktu, że

$$\begin{aligned} H(\theta, \vartheta) - H(\vartheta, \vartheta) &= - \int \log(f_{\theta}(x|y)) f_{\vartheta}(x|y) + - \int \log(f_{\vartheta}(x|y)) f_{\vartheta}(x|y) \\ &= - \int \log \left(\frac{f_{\theta}(x|y)}{f_{\vartheta}(x|y)} \right) f_{\vartheta}(x|y) dx \geq 0 \end{aligned}$$

EM dla rodziny wykładniczej

Niech

$$f_{\theta}(x, y) = \exp\{\theta^T S(x, y) - b(\theta)\}$$

Zatem

$$Q(\theta, \vartheta) = \mathbb{E}_{\vartheta}(\log(f_{\theta}(X, Y))|Y) = \theta^T \mathbb{E}_{\vartheta}(S(X, Y)|Y) - b(\theta)$$

EM dla mieszanek I

Y_1, \dots, Y_n z rozkładu

$$f_{\theta}(y) = \sum_{i=1}^k \alpha_i \phi_{\mu_i, \sigma_i}(y)$$

Łączny rozkład $P(X_j = i) = \alpha_i$, $f_{\theta}(y_j | x_j = i) = \phi_{\mu_i, \sigma_i}(y_j)$

$$f_{\theta}(y, x) = \prod_{j=1}^n \prod_{i=1}^k [\alpha_i \phi_{\mu_i, \sigma_i}(y_j)]^{1(x_j=i)}$$

Zatem

$$Q(\theta, \theta') = \sum_{j=1}^n \sum_{i=1}^k \mathbb{E}_{\theta'}(1(X_j = i) | Y_j = y_j) [\log(\alpha_i) + \log(\phi_{\mu_i, \sigma_i}(y_j))]$$

EM dla mieszanek II

Oznaczmy przez

$w_{ij} = \mathbb{E}_{\theta'}(1(X_j = i) | Y_j = y_j) = \mathbb{P}_{\theta'}(X_j = i | Y_j = y_j)$ Ze wzoru Bayesa mamy

$$w_{ij} = \frac{\alpha'_i \phi_{\mu'_i, \sigma'_i}(y_j)}{\sum_{i=1}^k \alpha'_i \phi_{\mu'_i, \sigma'_i}(y_j)}$$

Stąd

$$Q(\theta, \theta') = \sum_{j=1}^n \sum_{i=1}^k w_{ij} [\log(\alpha_i) + \log(\phi_{\mu_i, \sigma_i}(y_j))]$$

Maksymalizacja po α :

Chcemy zmaksymalizować $Q(\theta, \theta')$ względem α przy ograniczeniu $\sum \alpha_i = 1$. Lagrangian dla tego problemu to

$$\mathcal{L}(\alpha, \lambda) = \sum_{j=1}^n \sum_{i=1}^k w_{ij} \log(\alpha_j) - \lambda \left(\sum_{i=1}^k \alpha_j - 1 \right)$$

EM dla mieszanek III

Stąd

$$\frac{\partial \mathcal{L}(\alpha, \lambda)}{\partial \alpha_j} = \frac{\sum_{j=1}^n w_{ij}}{\alpha_j} - \lambda = 0$$

Czyli

$$\hat{\alpha}_i = \frac{\sum_{j=1}^n w_{ij}}{n}$$

Maksymalizacja po μ i σ^2 :

$$g(\mu, \sigma^2) = \sum_{j=1}^n \sum_{i=1}^k w_{ij} \left[\frac{-(y_j - \mu_i)^2}{2\sigma_i^2} - \frac{1}{2} \log(\sigma_i^2) \right]$$

$$\frac{\partial g}{\partial \mu_i} = \sum_{j=1}^n w_{ij} \frac{y_j - \mu_i}{\sigma_i^2} = 0$$

To

$$\hat{\mu}_i = \frac{\sum_{j=1}^n w_{ij} y_j}{\sum_{j=1}^n w_{ij}}$$

EM dla mieszanek IV

oraz ($\tau_i = \sigma_i^2$)

$$\frac{\partial g}{\partial \tau_i} = \sum_{j=1}^n w_{ij} \left[\frac{(y_j - \mu_i)^2}{\tau_i^2} - \frac{1}{2\tau_i} \right] = 0$$

Czyli

$$\hat{\sigma}_i^2 = \frac{\sum_{j=1}^n w_{ij} (y_j - \hat{\mu}_i)^2}{\sum_{j=1}^n w_{ij}}$$

Podsumowanie EM dla mieszanek gaussowskich

Krok E:

$$w_{ij}^{[t+1]} = \frac{\alpha_i^{[t]} \phi_{\mu_i^{[t]}, \sigma_i^{[t]}}(y_j)}{\sum_{i=1}^k \alpha_i^{[t]} \phi_{\mu_i^{[t]}, \sigma_i^{[t]}}(y_j)}$$

Krok M:

$$\alpha_i^{[t+1]} = \frac{\sum_{j=1}^n w_{ij}^{[t+1]}}{n}$$
$$\mu_i^{[t+1]} = \frac{\sum_{j=1}^n w_{ij}^{[t+1]} y_j}{\sum_{j=1}^n w_{ij}^{[t+1]}}$$
$$(\sigma_i^2)^{[t+1]} = \frac{\sum_{j=1}^n w_{ij}^{[t+1]} (y_j - \mu_i^{[t+1]})^2}{\sum_{j=1}^n w_{ij}^{[t+1]}}$$

EM dla wielowymiarowych gaussowskich

Krok E:

$$w_{ij}^{[t+1]} = \frac{\alpha_i^{[t]} \phi_{\mu_i^{[t]}, \Sigma_i^{[t]}}(y_j)}{\sum_{i=1}^k \alpha_i^{[t]} \phi_{\mu_i^{[t]}, \Sigma_i^{[t]}}(y_j)}$$

Krok M:

$$\alpha_i^{[t+1]} = \frac{\sum_{j=1}^n w_{ij}^{[t+1]}}{n}$$

$$\mu_i^{[t+1]} = \frac{\sum_{j=1}^n w_{ij}^{[t+1]} y_j}{\sum_{j=1}^n w_{ij}^{[t+1]}}$$

$$\Sigma_i^{[t+1]} = \frac{\sum_{j=1}^n w_{ij}^{[t+1]} (y_j - \mu_i^{[t+1]})(y_j - \mu_i^{[t+1]})^T}{\sum_{j=1}^n w_{ij}^{[t+1]}}$$

Uwagi: Estymowanie k macierzy kowariancji jest trudne (szczególnie dla większego wymiaru). W takiej sytuacji najczęściej zakłada się, że macierz kowariancji jest wspólna dla wszystkich rozkładów lub zakłada się szczególną postać macierzy kowariancji: $\Sigma_i = \sigma_i^2 I$, Σ_i diagonalna, itp.

EM jako proximal algorytm

Algorytm EM można zapisać w postaci

$$\theta_{k+1} = \arg \max_{\theta} Q(\theta, \theta_k)$$

Ponieważ

$$Q(\theta, \theta_k) + H(\theta_k, \theta_k) = \ell(\theta) - \text{KL}(f_{\theta_k}(x|y) || f_{\theta}(x|y))$$

Czyli EM równoważnie można zapisać jako

$$\theta_{k+1} = \arg \min_{\theta} \{-\ell(\theta) + \text{KL}(f_{\theta_k}(x|y) || f_{\theta}(x|y))\}$$

Podsumowanie I

- 1 Spadek po gradientie
 - 1 Wymaga obliczania tylko gradientu, koszt kroku $\mathcal{O}(dn)$
 - 2 Dla funkcji wypukłych i L gładkich tempo zbieżności $\mathcal{O}(k^{-1})$
 - 3 Przyspieszenie Nesterova poprawia tempo zbieżności do $\mathcal{O}(k^{-2})$
- 2 Algorytm Newtona-Raphsona
 - 1 Wymaga obliczania odwrotności hessianu, koszt kroku $\mathcal{O}(d^2n)$.
 - 2 Lokalna kwadratowa zbieżność, wrażliwy na punkt startowy.
 - 3 W przypadku GLM, można zapisać jako iterowane ważone najmniejsze kwadraty.
- 3 Proximal gradient
 - 1 Pozwala obliczać funkcję niegładkie o ile obliczenie kroku jest możliwe.
 - 2 Takie same tempo zbieżności co spadek po gradientie. Istnieje przyspieszenie Nesterova.

Podsumowanie II

④ ADMM

- ① Również pozwala na niegładkie funkcje.
- ② Często szybszy od proximal gradient w szczególności gdy poszczególne kroki można wyznaczyć analitycznie.

⑤ EM

- ① Pozwala efektywnie obliczać estymatory dla modeli z ukrytymi zmiennymi (brakującymi danymi).
- ② Krok E polega na imputacji ukrytych zmiennych.
- ③ EM jest monotoniczny. Dla funkcji wypukłych zbiega nie wolniej niż spadek po gradiencie.

2 Ukryte modele Markowa

- Dyskretne modele Markowa
- Obliczanie prawdopodobieństwa obserwacji
- Forward Backward
- Algorytm Viterbiego
- Algorytm FFBS
- Algorytm Bauma-Welcha
- Modele gaussowskie i filtr Kalmana

3 Algorytmy MCMC

- Podstawowe pojęcia statystyki bayesowskiej
- Metody Monte Carlo
- Algorytm Metropolisa-Hastingsa
- Próbnik Gibbsa
- Własności teoretyczne MCMC
- Adaptacyjne algorytmy MCMC
- Hamiltonian Monte Carlo

Łańcuchy Markowa

\mathcal{S} - dyskretna przestrzeń stanów

Definicja

Ciąg zmiennych losowych X_1, \dots, X_n, \dots jest łańcuchem Markowa wtedy i tylko wtedy gdy

$$\mathbb{P}(X_k = s_k | X_{k-1} = s_{k-1}, X_{k-2}, \dots, X_1) = \mathbb{P}(X_k = s_k | X_{k-1} = s_{k-1})$$

łańcuch Markowa jest jednorodny jeżeli

$$\mathbb{P}(X_k = s_k | X_{k-1} = s_{k-1}) = \mathbb{P}(X_2 = s_k | X_1 = s_{k-1}).$$

Macierz przejścia

$$\mathbb{P}(X_{k+1} = i | X_k = j) = P_k(j, i)$$

- Jeśli łańcuch jest jednorodny to dla każdego k $P_k = P$
- $\sum_{i \in \mathcal{S}} P_k(j, i) = 1$
- Jeżeli $X_k \sim v$ to

$$X_{k+1} \sim vP_k$$

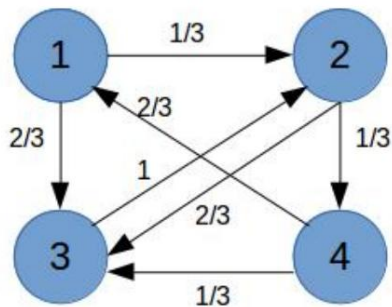
-

$$\mathbb{E}(f(X_{k+1}) | X_k = s) = (P_k f)(s)$$

- Równanie Chapmana - Kołmogorowa

$$\mathbb{P}(X_{k+2} | X_k = s) = \sum_{s'} \mathbb{P}(X_{k+2} | X_{k+1} = s') \mathbb{P}(X_{k+1} = s' | X_k = s)$$

Przykład



$$P = \begin{bmatrix} 0 & 1/3 & 2/3 & 0 \\ 0 & 0 & 2/3 & 1/3 \\ 0 & 1 & 0 & 0 \\ 2/3 & 0 & 1/3 & 0 \end{bmatrix}$$

Łańcuchy Markowa na ciągłej przestrzeni stanów

Jądro przejścia

$$\mathbb{P}(X_{k+1} \in dy | X_k = x) = P_k(x, dy)$$

- $P(x, \cdot)$ - miara probabilistyczna
- $P(\cdot, A)$ - funkcja mierzalna

-

$$\mu P(dy) = \int P(x, dy) \mu(dx)$$

-

$$Pf(x) = \int f(y) P(x, dy) = \mathbb{E}(f(X_2) | X_1 = x)$$

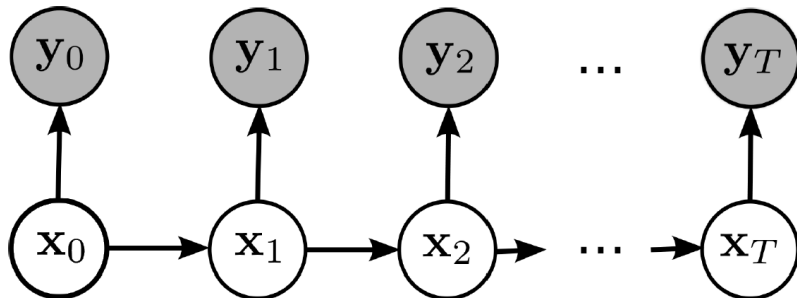
Przykłady

- Proces autoregresji

$$X_{k+1} = \rho X_k + Z_k$$

- Niech $Y = \rho X_k + Z_k$ i z prawdopodobieństwem $\alpha(X_k, Y)$
 $X_{k+1} = Y$ z pozostałym prawdopodobieństwem $X_{k+1} = X_k$

Ukryte modele Markowa



- 1 Policzyć prawdopodobieństwo obserwacji
- 2 Jak znaleźć najbardziej prawdopodobną ukrytą trajektorię?
- 3 Jak generować ukrytą trajektorię?
- 4 Jak wyestymować parametry modelu?

Ukryte modele Markowa

$$\mathbb{P}(X_{k+1} = i | X_k = j) = p(j, i)$$

$$\mathbb{P}(Y_k = y | X_k = j) = g(j, y)$$

$$\mathbb{P}(X_1 = x) = q(x)$$

- Rozkład łączny

$$\mathbb{P}(X_{1:T} = x_{1:T}, Y_{1:T} = y_{1:T}) = q(x_1)g(x_1, y_1)$$

$$\prod_{k=2}^T [g(x_k, y_k)p(x_{k-1}, x_k)]$$

- Rozkład wygładzania (smoothing)

$$\mathbb{P}(X_{1:T} = x_{1:T} | Y_{1:T} = y_{1:T}) = \frac{\mathbb{P}(X_{1:T} = x_{1:T}, Y_{1:T} = y_{1:T})}{\mathbb{P}Y_{1:T} = y_{1:T}}$$

- Rozkład filtru (filtering)

$$\mathbb{P}(X_{1:k} = x_{1:k} | Y_{1:k} = y_{1:k}) = \frac{\mathbb{P}(X_{1:k} = x_{1:k}, Y_{1:k} = y_{1:k})}{\mathbb{P}Y_{1:k} = y_{1:k}}$$

2 Ukryte modele Markowa

- Dyskretne modele Markowa
- Obliczanie prawdopodobieństwa obserwacji
- Forward Backward
- Algorytm Viterbiego
- Algorytm FFBS
- Algorytm Bauma-Welcha
- Modele gaussowskie i filtr Kalmana

3 Algorytmy MCMC

- Podstawowe pojęcia statystyki bayesowskiej
- Metody Monte Carlo
- Algorytm Metropolisa-Hastingsa
- Próbnik Gibbsa
- Własności teoretyczne MCMC
- Adaptacyjne algorytmy MCMC
- Hamiltonian Monte Carlo

Obliczanie prawdopodobieństwa obserwacji

$$\mathbb{P}(Y_{1:T} = y_{1:T}) = \sum_{\mathbf{x}_{1:T}} \mathbb{P}(Y_{1:T} = y_{1:T}, X_{1:T} = \mathbf{x}_{1:T})$$

Oznaczmy przez

$$\alpha_t(\mathbf{k}) = \mathbb{P}(X_t = \mathbf{k}, Y_{1:t} = y_{1:t})$$

Wówczas

$$\mathbb{P}(Y_{1:T} = y_{1:T}) = \sum_{\mathbf{k}} \alpha_T(\mathbf{k})$$

oraz

$$\alpha_1(\mathbf{k}) = q(\mathbf{k})g(\mathbf{k}, y_1)$$

$$\alpha_t(\mathbf{k}) = \mathbb{P}(X_t = \mathbf{k}, Y_{1:T} = y_{1:T})$$

$$= \sum_l \mathbb{P}(X_{t-1} = l, X_t = \mathbf{k}, Y_{1:T} = y_{1:T})$$

$$= \sum_l \alpha_{t-1}(l)p(l, \mathbf{k})g(\mathbf{k}, y_t)$$

Przykład

Mamy dwie monety S -symetryczną i N - niesymetryczną z prawdopodobieństwem $O = 0.8$. Rzucający może w sposób niezuważalny dla nas podmienić monetę, robi to z prawdopodobieństwem 0.3. Obserwujemy wyniki kolejnych rzutów.

O, O, O, R

Przykład

Mamy dwie monety S -symetryczną i N - niesymetryczną z prawdopodobieństwem $O = 0.8$. Rzucający może w sposób niezuważalny dla nas podmienić monetę, robi to z prawdopodobieństwem 0.3. Obserwujemy wyniki kolejnych rzutów.

O, O, O, R

	α_1	α_2	α_3	α_4
S	0,25	0,1475	0,094225	0,06214475
N	0,4	0,284	0,19444	0,0328751

$$\mathbb{P}(O, O, O, R) = 0,095$$

2 Ukryte modele Markowa

- Dyskretne modele Markowa
- Obliczanie prawdopodobieństwa obserwacji
- Forward Backward
- Algorytm Viterbiego
- Algorytm FFBS
- Algorytm Bauma-Welcha
- Modele gaussowskie i filtr Kalmana

3 Algorytmy MCMC

- Podstawowe pojęcia statystyki bayesowskiej
- Metody Monte Carlo
- Algorytm Metropolisa-Hastingsa
- Próbnik Gibbsa
- Własności teoretyczne MCMC
- Adaptacyjne algorytmy MCMC
- Hamiltonian Monte Carlo

Forward Backward I

$$\mathbf{x}_t^* = \arg \max_{\mathbf{x}_t} \mathbb{P}(\mathbf{X}_t = \mathbf{x}_t | \mathbf{Y}_{1:T} = \mathbf{y}_{1:T})$$

Równoważnie

$$\mathbf{x}_t^* = \arg \max_{\mathbf{x}_t} \mathbb{P}(\mathbf{X}_t = \mathbf{x}_t, \mathbf{Y}_{1:T} = \mathbf{y}_{1:T})$$

Oznaczmy przez

$$\beta_t(\mathbf{k}) = \mathbb{P}(\mathbf{Y}_{t+1:T} | \mathbf{X}_t = \mathbf{k}), \quad \beta_T(\mathbf{k}) = 1$$

Wówczas

$$\begin{aligned} \mathbb{P}(\mathbf{X}_t = \mathbf{x}_t, \mathbf{Y}_{1:T} = \mathbf{y}_{1:T}) &= \mathbb{P}(\mathbf{X}_t = \mathbf{x}_t, \mathbf{Y}_{1:t} = \mathbf{y}_{1:t}) \mathbb{P}(\mathbf{Y}_{t+1:T} | \mathbf{X}_t = \mathbf{x}) \\ &= \alpha_t(\mathbf{x}) \beta_t(\mathbf{x}) \end{aligned}$$

Forward Backward II

Dodatkowo

$$\begin{aligned}\beta_t(\mathbf{k}) &= \mathbb{P}(Y_{t+1:T} | X_t = \mathbf{k}) \\ &= \sum_l \mathbb{P}(Y_{t+1:T}, X_{t+1} = l | X_t = \mathbf{k}) \\ &= \sum_l \mathbb{P}(X_{t+1} = l | X_t = \mathbf{k}) \mathbb{P}(Y_{t+2:T}, | X_{t+1} = l) \mathbb{P}(Y_{t+1} | X_{t+1} = l) \\ &= \sum_l p(\mathbf{k}, l) g(l, y_{t+1}) \beta_{t+1}(l)\end{aligned}$$

Przykład

Mamy dwie monety S -symetryczną i N - niesymetryczną z prawdopodobieństwem $O = 0.8$. Rzucający może w sposób niezuważalny dla nas podmienić monetę, robi to z prawdopodobieństwem 0.3. Obserwujemy wyniki kolejnych rzutów.

O, O, O, R

Przykład

Mamy dwie monety S -symetryczną i N - niesymetryczną z prawdopodobieństwem $O = 0.8$. Rzucający może w sposób niezuważalny dla nas podmienić monetę, robi to z prawdopodobieństwem 0.3. Obserwujemy wyniki kolejnych rzutów.

O, O, O, R

	α_1	α_2	α_3	α_4
S	0,25	0,1475	0,094225	0,06214475
N	0,4	0,284	0,19444	0,0328751

	β_1	β_2	β_3	β_4
S	0,128321	0,21315	0,41	1
N	0.157349	0,2239	0,29	1

2 Ukryte modele Markowa

- Dyskretne modele Markowa
- Obliczanie prawdopodobieństwa obserwacji
- Forward Backward
- Algorytm Viterbiego
- Algorytm FFBS
- Algorytm Bauma-Welcha
- Modele gaussowskie i filtr Kalmana

3 Algorytmy MCMC

- Podstawowe pojęcia statystyki bayesowskiej
- Metody Monte Carlo
- Algorytm Metropolisa-Hastingsa
- Próbnik Gibbsa
- Własności teoretyczne MCMC
- Adaptacyjne algorytmy MCMC
- Hamiltonian Monte Carlo

Algorytm Viterbiego

$$\mathbf{x}_{1:T}^* = \arg \max_{\mathbf{x}_{1:T}} \mathbb{P}(X_{1:T} = \mathbf{x}_{1:T} | Y_{1:T} = y_{1:T})$$

Wystarczy

$$\mathbf{x}_{1:T}^* = \arg \max_{\mathbf{x}_{1:T}} \mathbb{P}(X_{1:T} = \mathbf{x}_{1:T}, Y_{1:T} = y_{1:T})$$

Zdefiniujmy:

$$v_1(j) = \mathbb{P}(X_1 = j, Y_1 = y_1) = q(j)g(j, y_1)$$

$$v_t(j) = \max_{\mathbf{x}_{1:(t-1)}} \mathbb{P}(X_{1:(t-1)} = \mathbf{x}_{1:(t-1)}, x_t = j, Y_{1:t} = y_{1:t})$$

$$\Psi_t(j) = \arg \max_{\mathbf{x}_{1:(t-1)}} \mathbb{P}(X_{1:(t-1)} = \mathbf{x}_{1:(t-1)}, x_t = j, Y_{1:t} = y_{1:t})$$

Przy tych oznaczeniach jeśli $j^* = \arg \max_j v_T(j)$ to

$$\mathbf{x}_{1:T}^* = [\Psi_T(j^*), j^*]$$

Obliczanie v_t

$$v_t(j) = g(j, y_t) \max_k [v_{t-1}(k)p(k, j)]$$

Przykład

Mamy dwie monety S -symetryczną i N - niesymetryczną z prawdopodobieństwem $O = 0.8$. Rzucający może w sposób niezuważalny dla nas podmienić monetę, robi to z prawdopodobieństwem 0.3. Obserwujemy wyniki kolejnych rzutów.

O, O, O, R

Przykład

Mamy dwie monety S -symetryczną i N - niesymetryczną z prawdopodobieństwem $O = 0.8$. Rzucający może w sposób niezuważalny dla nas podmienić monetę, robi to z prawdopodobieństwem 0.3. Obserwujemy wyniki kolejnych rzutów.

O, O, O, R

N, N, N, S

2 Ukryte modele Markowa

- Dyskretne modele Markowa
- Obliczanie prawdopodobieństwa obserwacji
- Forward Backward
- Algorytm Viterbiego
- Algorytm FFBS
- Algorytm Bauma-Welcha
- Modele gaussowskie i filtr Kalmana

3 Algorytmy MCMC

- Podstawowe pojęcia statystyki bayesowskiej
- Metody Monte Carlo
- Algorytm Metropolisa-Hastingsa
- Próbnik Gibbsa
- Własności teoretyczne MCMC
- Adaptacyjne algorytmy MCMC
- Hamiltonian Monte Carlo

Forward Filtering Backward Sampling

Cel: Generowanie $X_{1:T}|Y_{1:T}$

Przypomnijmy, że

$$\alpha_t(j) = \mathbb{P}(X_t = j, Y_{1:t} = y_{1:t})$$

Stąd

$$\mathbb{P}(X_T = x|Y_{1:T}) = \frac{\alpha_T(x)}{\sum_j \alpha_T(j)}$$

oraz

$$\mathbb{P}(X_{t-1}|X_t, Y_{1:t}) \propto \mathbb{P}(X_{t-1}, X_t, Y_{1:t})$$

gdzie

$$\begin{aligned} & \mathbb{P}(X_{t-1}, X_t, Y_{1:t}) \\ &= \alpha_{t-1}(X_{t-1})\mathbb{P}(X_t|X_{t-1}, Y_{1:t-1})\mathbb{P}(Y_{t:T}|X_t, X_{t-1}, Y_{1:t-1}) \\ &\propto \alpha_{t-1}p(X_{t-1}, X_t) \end{aligned}$$

2 Ukryte modele Markowa

- Dyskretne modele Markowa
- Obliczanie prawdopodobieństwa obserwacji
- Forward Backward
- Algorytm Viterbiego
- Algorytm FFBS
- Algorytm Bauma-Welcha
- Modele gaussowskie i filtr Kalmana

3 Algorytmy MCMC

- Podstawowe pojęcia statystyki bayesowskiej
- Metody Monte Carlo
- Algorytm Metropolisa-Hastingsa
- Próbnik Gibbsa
- Własności teoretyczne MCMC
- Adaptacyjne algorytmy MCMC
- Hamiltonian Monte Carlo

Algorytm Bauma-Welcha I

$$\theta = [q(j), p(i, j), g(j, y)]$$

I chcemy znaleźć

$$\arg \max_{\theta} \mathbb{P}(Y_{1:T})$$

Algorytm EM

Expectation:

$$\begin{aligned} \mathbb{P}(X_{1:T}, Y_{1:T}) &= q(x_1)g(x_1, y_1) \prod_{t=2}^T p(x_{t-1}, x_t)g(x_t, y_t) \\ &= \prod_j q(j)^{1(x_1=j)} \prod_i \prod_j p(i, j)^{\sum_{t=2}^T 1(x_{t-1}=i, x_t=j)} \\ &\quad \times \prod_j \prod_y g(j, y)^{\sum_{t=1}^T 1(x_t=j)1(y_t=y)} \end{aligned}$$

Algorytm Bauma-Welcha II

Stąd

$$\begin{aligned} Q(\theta, \theta') &= \sum_j \mathbb{P}_{\theta'}(X_1 = j | Y_{1:T}) \log(q(j)) \\ &+ \sum_i \sum_j \left[\sum_{t=2}^T \mathbb{P}_{\theta'}(X_{t-1} = i, X_t = j | Y_{1:T}) \right] \log(p(i, j)) \\ &+ \sum_y \sum_j \left[\sum_{t=1}^T 1(y_t = y) \mathbb{P}_{\theta'}(X_t = j | Y_{1:T}) \right] \log(g(j, y)) \end{aligned}$$

Ponadto

$$\gamma_t(j) = \mathbb{P}(X_t = j | Y_{1:T}) = \frac{\alpha_t(j)\beta_t(j)}{\sum_1 \alpha_t(1)\beta_t(1)}$$

Algorytm Bauma-Welcha III

oraz

$$\begin{aligned}\xi_t(i, j) &= \mathbb{P}(X_{t-1} = i, X_t = j | Y_{1:T}) \\ &= \mathbb{P}(X_{t-1} = i | Y_{1:T}) \mathbb{P}(X_t = j | X_{t-1} = i, Y_{t:T}) \\ &= \frac{\gamma_{t-1}(i) \mathbb{P}(X_t = j, Y_{t:T} | X_{t-1} = i)}{\mathbb{P}(Y_{t:T} | X_{t-1} = i)} \\ &= \frac{\gamma_{t-1}(i) \mathbb{P}(X_t = j | X_{t-1} = i) \mathbb{P}(Y_t | X_t = j) \mathbb{P}(Y_{t+1:T} | X_t = j)}{\beta_{t-1}(i)} \\ &= \frac{\gamma_{t-1}(i) p(i, j) g(j, y_t) \beta_t(j)}{\beta_{t-1}(i)}\end{aligned}$$

Algoritm Bauma-Welcha IV

Maximization:

$$\max_q \sum_j \gamma_1(j) \log(q(j)), \quad \sum_j q(j) = 1$$

$$\max_p \sum_i \sum_j \sum_{t=2}^T \xi_t(i, j) \log(p(i, j)), \quad \sum_j p(i, j) = 1$$

$$\max_g \sum_j \sum_y \sum_{t=1}^T 1(y_t = y) \gamma_t(j) \log(g(j, y)), \quad \sum_y g(j, y) = 1$$

Algorytm Bauma-Welcha V

Czyli

$$\hat{q}(j) = \frac{\gamma_1(j)}{\sum_j \gamma_1(j)}$$

$$\hat{p}(i, j) = \frac{\sum_{t=2}^T \xi_t(i, j)}{\sum_j \sum_{t=2}^T \xi_t(i, j)}$$

$$\hat{g}(j, y) = \frac{\sum_{t=1}^T \mathbf{1}(y_t = y) \gamma_t(j)}{\sum_{t=1}^T \gamma_j(t)}$$

2 Ukryte modele Markowa

- Dyskretne modele Markowa
- Obliczanie prawdopodobieństwa obserwacji
- Forward Backward
- Algorytm Viterbiego
- Algorytm FFBS
- Algorytm Bauma-Welcha
- Modele gaussowskie i filtr Kalmana

3 Algorytmy MCMC

- Podstawowe pojęcia statystyki bayesowskiej
- Metody Monte Carlo
- Algorytm Metropolisa-Hastingsa
- Próbnik Gibbsa
- Własności teoretyczne MCMC
- Adaptacyjne algorytmy MCMC
- Hamiltonian Monte Carlo

Linear state space model

$$\mathbf{X}_k = \mathbf{F}_k \mathbf{X}_{k-1} + \mathbf{B}_k \mathbf{u}_k + \mathbf{w}_k$$

$$\mathbf{Y}_k = \mathbf{H}_k \mathbf{X}_k + \mathbf{v}_k$$

$$\mathbf{w}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_k), \quad \mathbf{v}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_k), \quad \mathbf{X}_1 \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$

Rozkład filtracji i predykcji

$$p(\mathbf{x}_{k+1} | \mathbf{Y}_{1:k}) \sim \mathcal{N}(\hat{\mathbf{X}}_{k+1|k}, \boldsymbol{\Sigma}_{k+1|k}), \quad p(\mathbf{x}_k | \mathbf{Y}_{1:k}) \sim \mathcal{N}(\hat{\mathbf{X}}_{k|k}, \boldsymbol{\Sigma}_{k|k})$$

Stwierdzenie

Niech $X \sim N(\mu_x, \Sigma_x)$, $V \sim N(0, V)$ niezależne oraz

$$Y = HX + V$$

Wówczas

$$\begin{aligned}\mathbb{E}[X|Y] &= \mathbb{E}[X] + \text{Cov}(X, Y)\text{Cov}(Y)^{-1}(Y - \mathbb{E}[Y]) \\ &= \mu_x + \Sigma_X H^T [H\Sigma_X H^T + \Sigma_V]^{-1} (Y - H\mu_x)\end{aligned}$$

oraz

$$\begin{aligned}\text{Cov}(X|Y) &= \text{Cov}(X - \mathbb{E}[X|Y]) = \mathbb{E}[(X - \mathbb{E}[X|Y])X^T] \\ &= \Sigma_X - \Sigma_X H^T [H\Sigma_X H^T + \Sigma_V]^{-1} H\Sigma_X\end{aligned}$$

Dowód:

Oznaczmy prawą stronę przez \hat{X} i mamy

$$X - \hat{X} = X - \mathbb{E}[X] + \text{Cov}(X, Y)\text{Cov}(Y)^{-1}(Y - \mathbb{E}[Y])$$

Stąd

$$\text{Cov}(X - \hat{X}, Y) = \text{Cov}(X, Y) - \text{Cov}(X, Y)\text{Cov}(Y)^{-1}\text{Cov}(Y) = 0$$

Czyli $X - \hat{X}$ są niezależne oraz \hat{X} jest mierzalny względem $\sigma(Y)$ czyli

$$\mathbb{E}(X|Y) = \mathbb{E}(X - \hat{X} + \hat{X}|Y) = \mathbb{E}(X - \hat{X}) + \hat{X} = \hat{X}$$

Z definicji warunkowej wariancji

$$\text{Cov}(X|Y) = \mathbb{E}((X - \hat{X})(X - \hat{X})^T | Y) = \mathbb{E}((X - \hat{X})(X - \hat{X})^T) = \text{Cov}(X - \hat{X})$$

oraz

$$\mathbb{E}((X - \hat{X})(X - \hat{X})^T) = \mathbb{E}((X - \hat{X})X^T)$$

Filtr Kalmana

Inicjalizacja:

$$\hat{X}_{1|0} = \mu_0$$

$$\Sigma_{1|0} = \Sigma_0$$

Predykcja:

$$\hat{X}_{k|k-1} = F_k \hat{X}_{k-1|k-1} + B_k u_k$$

$$\Sigma_{k|k-1} = F_k \Sigma_{k-1|k-1} F_k^T + Q_k$$

Filtracja:

Inowacja

$$\hat{Y}_k = Y_k - H_k \hat{X}_{k|k-1}$$

Wariancja inowacji

$$S_k = R_k + H_k \Sigma_{k|k-1} H_k^T$$

Przyrost Kalmana

$$K_k = \Sigma_{k|k-1} H_k^T S_k^{-1}$$

Filtr

$$\hat{X}_{k|k} = \hat{X}_{k|k-1} + K_k \hat{Y}_k$$

Wariancja filtru

$$\Sigma_{k|k} = \Sigma_{k|k-1} - K_k H_k \Sigma_{k|k-1}$$

Wykładanie I

Chcemy policzyć rozkład $X_{k|T}$ wiedząc, że $X_{k+1|T} \sim N(\hat{X}_{k+1|T}, \Sigma_{k+1|T})$ Zauważmy, że

$$\begin{aligned}\mathbb{E}(X_k|Y_{1:T}) &= \mathbb{E}[\mathbb{E}(X_k|X_{k+1}, Y_{1:k})|Y_{1:T}] \\ \text{Cov}(X_k|Y_{1:T}) &= \mathbb{E}[\text{Cov}((X_k|X_{k+1}, Y_{1:k}))|Y_{1:T}] \\ &\quad + \text{Cov}[\mathbb{E}((X_k|X_{k+1}, Y_{1:k}))|Y_{1:T}]\end{aligned}$$

Czyli trzeba policzyć rozkład $X_k|X_{k+1}, Y_{1:k}$

$$\begin{aligned}\mathbb{E}[X_k, X_{k+1}|Y_{1:k}]^T &= [\hat{X}_{k|k}, \hat{X}_{k+1|k}]^T \\ \text{Cov}([X_k, X_{k+1}|Y_{1:k}]^T) &= \begin{bmatrix} \Sigma_{k|k} & \Sigma_{k|k}F_{k+1}^T \\ \Sigma_{k|k}F_{k+1}^T & \Sigma_{k+1|k} \end{bmatrix}\end{aligned}$$

Wykładanie II

Lemat

Jeśli

$$\begin{pmatrix} Z_1 \\ Z_2 \end{pmatrix} \sim N \left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{21} \\ \Sigma_{21}^T & \Sigma_{22} \end{pmatrix} \right)$$

to

$$Z_1 | Z_2 \sim \mathcal{N}(\mu_1 + \Sigma_{21} \Sigma_{22}^{-1} (Z_2 - \mu_2), \Sigma_{11} - \Sigma_{21} \Sigma_{22}^{-1} \Sigma_{21}^T)$$

Stąd

$$X_k | X_{k+1}, Y_{1:k} \sim N(\hat{X}_{k|k} + L_k (X_{k+1} - \hat{X}_{k+1|k}), \Sigma_{k|k} - L_k \Sigma_{k+1|k}^{-1} L_k^T)$$

$$L_k = \Sigma_{k|k} F_{k+1}^T \Sigma_{k+1|k}^{-1}$$

Wykładanie III

Czyli

$$\begin{aligned}\hat{X}_{k|T} &= \hat{X}_{k|k} + L_k(\hat{X}_{k+1|T} - \hat{X}_{k+1|k}) \\ \Sigma_{k|T} &= \Sigma_{k|k} - L_k \Sigma_{k+1|k}^{-1} L_k^T + L_k \Sigma_{k+1|T} L_k^T \\ &= \Sigma_{k|k} + L_k(\Sigma_{k+1|T} - \Sigma_{k+1|k}) L_k^T\end{aligned}$$

Rauch–Tung–Striebel

$$\begin{aligned}L_k &= \Sigma_{k|k} F_{k+1}^T \Sigma_{k+1|k}^{-1} \\ \hat{X}_{k|T} &= \hat{X}_{k|k} + L_k (\hat{X}_{k+1|T} - \hat{X}_{k+1|k}) \\ \Sigma_{k|T} &= \Sigma_{k|k} + L_k (\Sigma_{k+1|T} - \Sigma_{k+1|k}) L_k^T\end{aligned}$$

EM dla liniowych modeli

$$X_k = FX_{k-1} + w_k$$

$$Y_k = HX_k + v_k$$

$$w_k \sim N(0, Q), \quad v_k \sim N(0, R), \quad X_1 \sim N(\mu_0, \Sigma_0)$$

Chcemy estymować F, H, Q, R

Log wiarygodność

$$\begin{aligned} & \ell(\mathbf{F}, \mathbf{H}, \mathbf{Q}, \mathbf{R}) \\ &= \sum_{k=2}^T \left[-\frac{1}{2} \log(|\mathbf{Q}|) - \frac{1}{2} (\mathbf{x}_k - \mathbf{F}\mathbf{x}_{k-1})^T \mathbf{Q}^{-1} (\mathbf{x}_k - \mathbf{F}\mathbf{x}_{k-1}) \right] \\ &+ \sum_{k=1}^T \left[-\frac{1}{2} \log(|\mathbf{R}|) - \frac{1}{2} (\mathbf{y}_k - \mathbf{H}\mathbf{x}_k)^T \mathbf{R}^{-1} (\mathbf{y}_k - \mathbf{H}\mathbf{x}_k) \right] \\ &= \frac{T-1}{2} \log(|\mathbf{Q}|^{-1}) - \frac{1}{2} \text{tr} \left(\mathbf{Q}^{-1} \sum_{k=2}^T (\mathbf{x}_k - \mathbf{F}\mathbf{x}_{k-1})(\mathbf{x}_k - \mathbf{F}\mathbf{x}_{k-1})^T \right) \\ &+ \frac{T}{2} \log(|\mathbf{R}|^{-1}) - \frac{1}{2} \text{tr} \left(\mathbf{R}^{-1} \sum_{k=1}^T (\mathbf{y}_k - \mathbf{H}\mathbf{x}_k)(\mathbf{y}_k - \mathbf{H}\mathbf{x}_k)^T \right) \\ &= \frac{T-1}{2} \log(|\mathbf{Q}|^{-1}) - \frac{1}{2} \text{tr} \left(\mathbf{Q}^{-1} \sum_{k=2}^T \left[\mathbf{x}_k \mathbf{x}_k^T + \mathbf{F}\mathbf{x}_{k-1} \mathbf{x}_{k-1}^T \mathbf{F}^T - \mathbf{x}_k \mathbf{x}_{k-1}^T \mathbf{F}^T - \mathbf{F}\mathbf{x}_{k-1} \mathbf{x}_k^T \right] \right) \\ &+ \frac{T}{2} \log(|\mathbf{R}|^{-1}) - \frac{1}{2} \text{tr} \left(\mathbf{R}^{-1} \sum_{k=1}^T \left[\mathbf{y}_k \mathbf{y}_k^T + \mathbf{H}\mathbf{x}_k \mathbf{x}_k^T \mathbf{H}^T - \mathbf{y}_k \mathbf{x}_k^T \mathbf{H}^T - \mathbf{H}\mathbf{x}_k \mathbf{y}_k^T \right] \right) \end{aligned}$$

Expectation

$$\mathbb{E}(x_k | Y_{1:T}) = \hat{X}_{k|T}$$

$$\mathbb{E}(x_k x_k^T | Y_{1:T}) = \hat{X}_{k|T} \hat{X}_{k|T}^T + \Sigma_{k|T}$$

$$\begin{aligned}\mathbb{E}(x_k x_{k+1}^T | Y_{1:T}) &= \mathbb{E}(x_k x_{k+1}^T | x_{k+1}, Y_{1:k}) | Y_{1:T} \\ &= \mathbb{E}[(\hat{X}_{k|k} + L_k(x_{k+1} - \hat{X}_{k+1|k})) x_{k+1}^T | Y_{1:T}] \\ &= \hat{X}_{k|k} \hat{X}_{k+1|T}^T + L_k \left(\Sigma_{k+1|T} + (\hat{X}_{k+1|T} - \hat{X}_{k+1|K}) \hat{X}_{k+1|T}^T \right)\end{aligned}$$

Maximization I

Potrzebne fakty:



$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = \mathbf{x}^T (\mathbf{A}^T + \mathbf{A})$$



$$\frac{\partial \mathbf{B}^T \mathbf{A} \mathbf{B}}{\partial \mathbf{B}} = \mathbf{B}^T (\mathbf{A}^T + \mathbf{A})$$



$$\frac{\partial \text{tr}(\mathbf{A} \mathbf{B})}{\partial \mathbf{A}} = \frac{\partial \text{tr}(\mathbf{B} \mathbf{A})}{\partial \mathbf{A}} = \frac{\partial \text{tr}(\mathbf{B}^T \mathbf{A}^T)}{\partial \mathbf{A}} = \mathbf{B}^T$$



$$\frac{\partial \text{tr}(\mathbf{A} \mathbf{B} \mathbf{A}^T)}{\partial \mathbf{A}} = 2 \mathbf{A} \mathbf{B}$$



$$\frac{\partial \log(|\mathbf{A}|)}{\partial \mathbf{A}} = \mathbf{A}^{-T}$$

Maximization II

Maksymalizacja po F

$$\frac{\partial \ell}{\partial F} = \sum_{k=2}^T x_k x_{k-1}^T Q^{-1} - F \sum_{k=2}^T x_{k-1} x_{k-1}^T Q^{-1} = 0$$

Czyli

$$F = \sum_{k=2}^T x_k x_{k-1}^T \left(\sum_{k=2}^T x_{k-1} x_{k-1}^T \right)^{-1}$$

Maximization III

Maksymalizacja po H

$$\frac{\partial \ell}{\partial H} = \sum_{k=1}^T x_k y_k^T R^{-1} - H \sum_{k=1}^T x_k x_k^T R^{-1} = 0$$

Czyli

$$H = \sum_{k=1}^T x_k y_k^t \left(\sum_{k=1}^T x_k x_k^T \right)^{-1}$$

Maximization IV

Maksymalizacja po Q

$$\frac{\partial \ell}{\partial Q^{-1}} = \frac{T-1}{2} Q - \frac{1}{2} \left\{ \sum_{k=2}^T \left[\mathbf{x}_k \mathbf{x}_k^T + \mathbf{F} \mathbf{x}_{k-1} \mathbf{x}_{k-1}^T \mathbf{F}^T - \mathbf{x}_k \mathbf{x}_{k-1}^T \mathbf{F}^T - \mathbf{F} \mathbf{x}_{k-1} \mathbf{x}_k^T \right] \right\}^T$$

Czyli

$$Q = \frac{1}{T-1} \left\{ \sum_{k=2}^T \left[\mathbf{x}_k \mathbf{x}_k^T + \mathbf{F} \mathbf{x}_{k-1} \mathbf{x}_{k-1}^T \mathbf{F}^T - \mathbf{x}_k \mathbf{x}_{k-1}^T \mathbf{F}^T - \mathbf{F} \mathbf{x}_{k-1} \mathbf{x}_k^T \right] \right\}$$

Maximization V

Maksymalizacja po R

$$\frac{\partial \ell}{\partial \mathbf{R}^{-1}} = \frac{\mathbf{T}}{2} \mathbf{R} - \frac{1}{2} \left\{ \sum_{k=1}^{\mathbf{T}} \left[\mathbf{y}_k \mathbf{y}_k^{\mathbf{T}} + \mathbf{H} \mathbf{x}_k \mathbf{x}_k^{\mathbf{T}} \mathbf{H}^{\mathbf{T}} - \mathbf{y}_k \mathbf{x}_k^{\mathbf{T}} \mathbf{H}^{\mathbf{T}} - \mathbf{H} \mathbf{x}_k \mathbf{y}_k^{\mathbf{T}} \right] \right\}^{\mathbf{T}}$$

Czyli

$$\mathbf{R} = \frac{1}{\mathbf{T}} \left\{ \sum_{k=1}^{\mathbf{T}} \left[\mathbf{y}_k \mathbf{y}_k^{\mathbf{T}} + \mathbf{H} \mathbf{x}_k \mathbf{x}_k^{\mathbf{T}} \mathbf{H}^{\mathbf{T}} - \mathbf{y}_k \mathbf{x}_k^{\mathbf{T}} \mathbf{H}^{\mathbf{T}} - \mathbf{H} \mathbf{x}_k \mathbf{y}_k^{\mathbf{T}} \right] \right\}$$

2 Ukryte modele Markowa

- Dyskretne modele Markowa
- Obliczanie prawdopodobieństwa obserwacji
- Forward Backward
- Algorytm Viterbiego
- Algorytm FFBS
- Algorytm Bauma-Welcha
- Modele gaussowskie i filtr Kalmana

3 Algorytmy MCMC

- Podstawowe pojęcia statystyki bayesowskiej
- Metody Monte Carlo
- Algorytm Metropolisa-Hastingsa
- Próbnik Gibbsa
- Własności teoretyczne MCMC
- Adaptacyjne algorytmy MCMC
- Hamiltonian Monte Carlo

Model bayesowski

W statystyce bayesowskiej zakłada się, że parametry rozkładu są zmiennymi losowymi

$$\theta \sim \pi(\theta) \quad \text{rozkład a'priori.}$$

Obserwacje przy znanym θ mają rozkład

$$X \sim p(X|\theta) \quad \text{wiarygodność.}$$

Wnioskujemy o parametrach na podstawie rozkładu θ pod warunkiem zaobserwowanego X

$$\pi(\theta|X) = \frac{\pi(\theta)p(X|\theta)}{\int_X \pi(\theta)p(x|\theta)dx} \propto \pi(\theta)p(X|\theta) \quad \text{rozkład a'posteriori}$$

Przykład 1

S liczba wyrzuconych orłów w N rzutach monetą o prawdopodobieństwie wyrzucenia orła θ .

$$p(S|\theta) \propto \theta^S(1 - \theta)^{N-S}$$

Zakładamy, że rozkład a’piori jest rozkładem Beta(α, β) czyli

$$\pi(\theta) \propto \theta^{\alpha-1}(1 - \theta)^{\beta-1}.$$

Stąd

$$\pi(\theta|S) \propto \theta^{S+\alpha-1}(1 - \theta)^{N-S+\beta-1}$$

czyli $\pi(\theta|S) \sim \text{Beta}(S + \alpha, N - S + \beta)$.

Przykład 2

$X_1, \dots, X_n \sim N(\mu, \sigma^2)$, μ - nieznane, σ^2 znane

$$\pi(\mu) \sim N(m, v^2)$$

Rozkład a'posteriori

$$\begin{aligned}\pi(\mu|X_1, \dots, X_n) &\propto \exp \left\{ -\frac{\sum_{i=1}^n (X_i - \mu)^2}{2\sigma^2} - \frac{(\mu - m)^2}{2v^2} \right\} \\ &\propto \exp \left\{ -\frac{nv^2 + \sigma^2}{2v^2\sigma^2} \mu^2 + 2\frac{n\bar{X}v^2 + m\sigma^2}{2v^2\sigma^2} \mu \right\} \\ &\propto \exp \left\{ -\frac{1}{2} \frac{nv^2 + \sigma^2}{v^2\sigma^2} \left(\mu - \frac{n\bar{X}v^2 + m\sigma^2}{nv^2 + \sigma^2} \right)^2 \right\}\end{aligned}$$

Czyli

$$\pi(\mu|X_1, \dots, X_n) \sim N \left(\frac{n\bar{X}v^2 + m\sigma^2}{nv^2 + \sigma^2}, \frac{v^2\sigma^2}{nv^2 + \sigma^2} \right)$$

Przykład 3

$X_1, \dots, X_n \sim N(\mu, 1/\tau)$, μ - znane, τ - nieznanne

$$\pi(\tau) \sim \text{Gamma}(a, b)$$

Rozkład a'posteriori

$$\begin{aligned}\pi(\tau|X_1, \dots, X_n) &\propto \tau^{a-1} \tau^{n/2} \exp \left\{ -\tau \frac{\sum_{i=1}^n (X_i - \mu)^2}{2} - \tau b \right\} \\ &\propto \tau^{a+n/2-1} \exp \left\{ -\tau \left(\frac{\sum_{i=1}^n (X_i - \mu)^2}{2} + b \right) \right\}\end{aligned}$$

Czyli

$$\pi(\tau|X_1, \dots, X_n) \sim \text{Gamma} \left(a + n/2, b + \frac{\sum_{i=1}^n (X_i - \mu)^2}{2} \right)$$

Estymatory bayesowskie I

Funkcja straty:

$$L(\theta, a)$$

Strata gdy prawdziwa wartość jest θ a my przewidujemy a . L wypukła, równa 0 wtedy i tylko wtedy gdy $\theta = a$. Standardowo używa się kwadratowej funkcji straty $(\theta - a)^2$.

Ryzyko estymatora:

$\delta(X)$ -estymator

$$R(\delta) = \int_{\Theta} \int_{\mathcal{X}} L(\theta, \delta(X)) p(X|\theta) \pi(\theta) dX d\theta.$$

Ryzyko bayesowskie:

$$R_B(\delta(X)) = \int_{\Theta} L(\theta, \delta(X)) \pi(\theta|X) d\theta.$$

Estymatory bayesowskie II

Jeżeli zminimalizujemy ryzyko bayesowskie to minimalizujemy również zwykłe ryzyko

$$\begin{aligned}R(\delta) &= \int_{\Theta} \int_{\mathcal{X}} L(\theta, \delta(X)) p(X|\theta) \pi(\theta) dX d\theta \\ &= \int_{\mathcal{X}} \int_{\Theta} L(\theta, \delta(X)) \pi(\theta|X) d\theta p(X) dX \\ &= \int_{\mathcal{X}} R_B(\delta(X)) dX\end{aligned}$$

Estymator bayesowski:

$$\hat{\theta}_B = \arg \min_{\delta} R_B(\delta(X))$$

W przypadku kwadratowej funkcji straty

$$\hat{\theta}_B = E(\theta|X)$$

Przykłady estymatorów bayesowskich

Model Dwumianowy-Beta:

$$\pi(\theta|S) \sim \text{Beta}(S + \alpha, N - S + \beta)$$

$$\hat{\theta}_B = \frac{S + \alpha}{N + \alpha + \beta} = \frac{S}{N} \frac{N}{N + \alpha + \beta} + \frac{\alpha}{\alpha + \beta} \frac{\alpha + \beta}{N + \alpha + \beta}$$

Model Normalny-Normalny:

$$\pi(\mu|X_1, \dots, X_n) \sim N\left(\frac{n\bar{X}v^2 + m\sigma^2}{nv^2 + \sigma^2}, \frac{v^2\sigma^2}{nv^2 + \sigma^2}\right)$$

$$\hat{\mu}_B = \frac{n\bar{X}v^2 + m\sigma^2}{nv^2 + \sigma^2} = \bar{X} \frac{nv^2}{nv^2 + \sigma^2} + m \frac{\sigma^2}{nv^2 + \sigma^2}$$

Model Normalny-Gammy:

$$\pi(\tau|X_1, \dots, X_n) \sim \text{Gamma}\left(a + n/2, b + \frac{\sum_{i=1}^n (X_i - \mu)^2}{2}\right)$$

$$\hat{\tau}_B = \frac{a + n/2}{b + \sum_{i=1}^n (X_i - \mu)^2/2}$$

Obrobine bardziej złożony model

$X_1, \dots, X_n \sim N(\mu, 1/\tau)$, μ - nieznane, τ -nieznane

$$\pi(\mu) \sim N(m, v^2), \quad \pi(\tau) \sim \text{Gamma}(a, b)$$

μ i τ niezależne Rozkład a'posteriori

$$\pi(\mu, \tau | X_1, \dots, X_n) \propto \tau^{a+n/2-1} \exp \left\{ -\tau \frac{\sum_{i=1}^n (X_i - \mu)^2}{2} - \frac{(\mu - m)^2}{2v^2} - \tau b \right\}$$

Czyli

$$\pi(\mu, \tau | X_1, \dots, X_n) \sim \text{???}, \quad \hat{\tau}_B = \text{???} \quad \hat{\mu}_B = \text{???}$$

Znamy $\pi(\mu, \tau | X_1, \dots, X_n)$ z dokładnością do stałej normującej

Hierarchiczny model bayesowski

Mamy k grup (np powiatów), w j -tej grupie mamy N_j obserwacji

$$X_{1j}, \dots, X_{N_j, j} \sim N(\mu_j, 1/\tau_j)$$

Następnie

$$\mu_1, \dots, \mu_k \sim N(\mu, 1/\tau)$$

$$\tau_1, \dots, \tau_k \sim \text{Gamma}(a, b)$$

μ i τ odpowiadają wartości oczekiwanej i wariancji badanej cechy w całym kraju. Są to z reguły wartości nieznane i na nie również nakładamy rozkłady a'priori

$$\mu \sim N(m, v^2)$$

$$\tau \sim \text{Gamma}(\alpha, \lambda)$$

Rozkład a'posteriori

$$\pi(\mu_1, \dots, \mu_k, \tau_1, \dots, \tau_k, \mu, \tau | X, a, b, \alpha, \lambda, m, v^2) \propto \prod_{j=1}^k \left\{ \prod_{i=1}^{N_j} [p(X_{ij} | \mu_j, \tau_j)] \pi(\mu_j | \mu, \tau) \pi(\tau_j | a, b) \right\} \pi(\mu | m, v^2) \pi(\tau | \alpha, \lambda)$$

Inne elementy wnioskowania bayesowskiego

Przedziały ufności:

Szukamy najkrótszego przedziału $[\underline{\theta}, \bar{\theta}]$ takiego, że

$$P(\theta \in [\underline{\theta}, \bar{\theta}] | X) = 1 - \alpha$$

Prognozowanie:

$$\pi(X^* | X) = \int_{\Theta} p(X^* | \theta) \pi(\theta | X) d\theta$$

Dla kwadratowej funkcji straty predykcja bayesowska

$$E(X^* | X)$$

Problemy obliczeniowe

Potrzebujemy metod pozwalających:

- Aproksymować/ generować z rozkładu $\pi(\theta|X)$ znanego z dokładnością do stałej.
- Obliczać całki względem rozkładu jak powyżej.

2 Ukryte modele Markowa

- Dyskretne modele Markowa
- Obliczanie prawdopodobieństwa obserwacji
- Forward Backward
- Algorytm Viterbiego
- Algorytm FFBS
- Algorytm Bauma-Welcha
- Modele gaussowskie i filtr Kalmana

3 Algorytmy MCMC

- Podstawowe pojęcia statystyki bayesowskiej
- Metody Monte Carlo
- Algorytm Metropolisa-Hastingsa
- Próbnik Gibbsa
- Własności teoretyczne MCMC
- Adaptacyjne algorytmy MCMC
- Hamiltonian Monte Carlo

Proste Monte Carlo

Chcemy obliczyć

$$I = \int f(x)\pi(x)dx,$$

gdzie π gęstość prawdopodobieństwa. Czyli z MPWL

$$I = \text{Ef}(X) \approx \frac{1}{n} \sum_{i=1}^n f(X_i) =: \hat{I}_n,$$

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

Własności estymatora MC:

- $E\hat{I}_n = I$
- Jeśli $\text{Var}(f(X_i)) = \sigma^2$ to

$$E(I - \hat{I}_n)^2 = \frac{\sigma^2}{n}$$

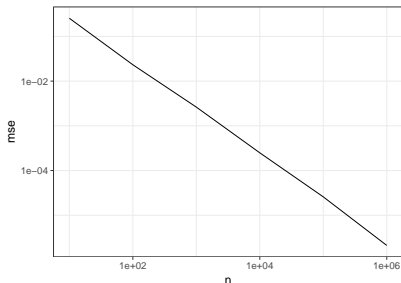
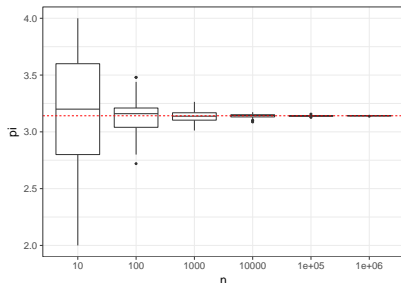
Przykład oblicznie liczby π

Jeśli $X \sim \text{Unif}([0, 1] \times [0, 1])$ to

$$\pi = 4E1(X_1^2 + X_2^2 \leq 1)$$

Jeśli $X^1, \dots, X^n \sim \text{Unif}([0, 1] \times [0, 1])$ to

$$\pi \approx \frac{1}{n} \sum_{i=1}^n 1((X_1^i)^2 + (X_2^i)^2 \leq 1)$$



Rejection Sampling I

Chcemy wylosować zmienną losową o gęstości f , ale potrafimy generować z g oraz wiemy, że

$$f(x) \leq Mg(x)$$

Algorytm:

- 1 $Y \sim g$
- 2 $U \sim \text{Unif}([0, 1])$
- 3 Jeśli $U \leq \frac{f(Y)}{Mg(Y)}$ to zwracamy Y w.p.p wracamy do kroku 1

Poprawność: Niech $Z_i = (U_i, Y_i)$, $A = \{(y, u) : u \leq \frac{f(y)}{Mg(y)}\}$,

$$N = \min\{n : Z_i \in A\}$$

Wówczas N i Z_N są niezależne oraz

$$P(Z_N \in B) = P(Z \in B | Z \in A)$$

Rejection Sampling II

Obliczmy najpierw

$$P(N = n) = P(Z_1 \notin A, \dots, Z_{n-1} \notin A, Z_n \in A) = q^{n-1}p,$$

gdzie $p = P(Z_1 \in A)$ oraz $q = 1 - p$ Teraz

$$\begin{aligned} P(N = n, Z_N \in B) &= P(Z_1 \notin A, \dots, Z_{n-1} \notin A, Z_n \in A \cap B) \\ &= P(Z_1 \notin A)^{n-1} P(Z_n \in B) \frac{P(Z_n \in A \cap B)}{P(Z_n \in B)} \end{aligned}$$

Teraz wystarczy obliczyć

$$P(Y \in B | Z \in A) = \frac{\int_B \int_0^{f(y)/Mg(y)} g(y) du dy}{\int \int_0^{f(y)/Mg(y)} g(y) du dy} = \frac{\int_B f(y) dy}{\int f(y) dy}$$

Oczekiwana liczba kroków algorytmu:

$$p = P(Z_1 \in A) = \int \frac{f(y)}{Mg(y)} g(y) dy = \frac{1}{M}$$

Stąd

$$EN = M$$

Przykład

Chcemy generować z X rozkładu gaussowskiego $N(0, 1)$ uciętego do odcinka $[3, 4]$ czyli

$$f(x) = \frac{\exp(-x^2/2)}{\int_3^4 \exp(-x^2/2)} 1(x \in [3, 4])$$

I korzystamy z dwóch gęstości pomocniczych g_1 gęstość $N(0, 1)$ i $g_2 \propto x \exp(-x^2/2) 1(x \in [3, 4])$ W pierwszym przypadku do wygenerowania próbki długości 100 potrzeba 73936 zmiennych z rozkładu g_1 natomiast dla g_2 wystarczyło 303 zmiennych.

Własności rejection sampling

- Jeśli potrafimy dobrze, przybliżyć **globalnie** gęstość to algorytm jest efektywny.
- Generuje próbkę i.i.d
- W przypadku użycia do obliczania całek “marnujemy” znaczną część wygenerowanych zmiennych losowych.

Importance sampling

Chcemy obliczyć

$$I = \int f(x)p(x)dx$$

p gęstość rozkładu. Dysponujemy próbką X_1, \dots, X_n i.i.d z rozkładu q . Zauważmy, że

$$I = \int f(x)p(x)dx = \int f(x)\frac{p(x)}{q(x)}q(x) = E_q f(X)w(X),$$

gdzie $w = \frac{p}{q}$.

$$I_{IS} = \frac{1}{n} \sum_{i=1}^n f(X_i)w(X_i)$$

Własności IS

- $E(I_{IS}) = I$
- $\text{Var}(I_{IS}) = \frac{\text{Var}_q(f(X)w(X))}{n}$
- Wymaga znajomości stałej normującej f

Samonormujący IS

Zauważmy, że

$$E_q(w(X)) = 1$$

Zastąpmy n przez estymator $\hat{n} = \sum_{i=1}^n w(X_i)$. Czyli

$$\bar{I}_{IS} = \frac{\sum_i f(X_i)w(X_i)}{\sum_{i=1}^n w(X_i)}$$

Teraz w występuje w liczniku i mianowniku.

- $E(\bar{I}_{IS}) \neq I$ ale $\lim_{n \rightarrow \infty} E\bar{I}_{IS} = I$
- $\text{Var}(\bar{I}_{IS}) = \mathcal{O}(1/n)$
- Można stosować do gęstości znanej z dokładnością do stałej normującej.

Efektywność IF

Proste MC

$$\text{MSE}(\hat{I}) = \frac{\sigma^2}{n}$$

Dla IS

$$\text{MSE}(\bar{I}_{\text{IS}}) = \frac{c}{n}$$

Effective Sample Size

Dla jakiej długości próbki i.i.d. dostaniemy ten sam błąd co nasz estymator. Dla IS ESS przybliża się następująco:

Traktujemy wagi $w_i = w(X_i)$ jako deterministyczne wówczas wariancja estymatora IS jest postaci

$$\frac{\sum w_i^2}{(\sum w_i)^2} \text{Var}_q(f(X_i))$$

I ESS definiujemy jako

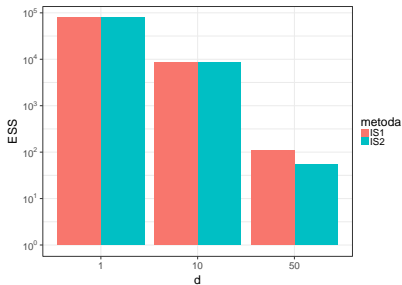
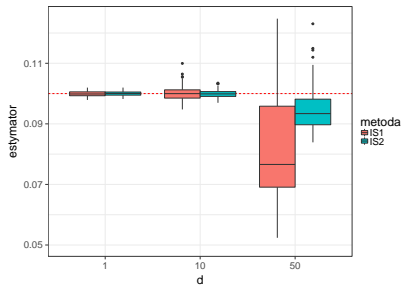
$$\text{ESS} = \frac{(\sum w_i)^2}{\sum w_i^2}$$

Przykład

$Y \sim N(0.1, 0.2^2 I_d)$ i chcemy obliczyć

$$I = \frac{1}{d} \sum_{i=1}^d EY_i = 0.1$$

Bedziemy korzystać z (samonormującej) IS z rozkładem propozycji $N(0, 0.2^2 I_d)$.



Własności IS

- Wykorzystuje całą próbkę do estymacji.
- Wrażliwa na dobór rozkładu pomocniczego. Trzeba dobrze przybliżyć rozkład globalnie.
- Nie wymaga znajomości stałej normującej

MCMC

Wadą standardowych MC jest to, że kolejne iteracje niezależą od poprzednich. Co powoduje, że każdy krok musi być dopasowany globalnie do rozkładu. Chcielibyśmy eksplorować rozkład lokalnie, np. w okolicy wyniku z poprzedniego kroku.

MCMC:

Chcemy wygenerować ergodyczny łańcuch Markowa X_1, \dots, X_n o rozkładzie stacjonarnym π . Wówczas na podstawie MPWL dla łańcuchów Markowa otrzymujemy

$$I = \int f(x)\pi(x)dx \approx \frac{1}{n} \sum_{i=1}^n f(X_i)$$

Z uwagi na to, że asymptotycznie mamy właściwy rozkład to stosuje się “burn-in” time

$$\hat{I}_{n_0, n} = \frac{1}{n} \sum_{i=n_0+1}^{n+n_0} f(X_i)$$

Rozkład stacjonarny i odwracalność I

Łańcuch Markowa z jądrem przejścia $P(x, dy)$ ma rozkład stacjonarny π wtedy i tylko wtedy gdy

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

Łańcuch Markowa z jądrem przejścia $P(x, dy)$ jest odwracalny względem π wtedy i tylko wtedy gdy

$$\int_{A \times B} P(x, dy)\pi(x)dx = \int_{A \times B} P(y, dx)\pi(y)dy$$

lub w skrócie

$$P(x, dy)\pi(dx) = P(y, dx)\pi(dy)$$

Rozkład stacjonarny i odwracalność II

Łańcuch Markowa jest odwracalny względem $\pi \implies \pi$ jest rozkładem stacjonarnym

$$\begin{aligned}\pi P(A) &= \int P(x, A) \pi(x) dx = \int \int_A \pi(x) P(x, dy) dx \\ &= \int_A \int P(y, dx) \pi(y) dy = \int_A \pi(y)\end{aligned}$$

2 Ukryte modele Markowa

- Dyskretne modele Markowa
- Obliczanie prawdopodobieństwa obserwacji
- Forward Backward
- Algorytm Viterbiego
- Algorytm FFBS
- Algorytm Bauma-Welcha
- Modele gaussowskie i filtr Kalmana

3 Algorytmy MCMC

- Podstawowe pojęcia statystyki bayesowskiej
- Metody Monte Carlo
- Algorytm Metropolisa-Hastingsa
- Próbnik Gibbsa
- Własności teoretyczne MCMC
- Adaptacyjne algorytmy MCMC
- Hamiltonian Monte Carlo

Algorytm Metropolisa-Hastingsa

- $Q(x, dy)$ jądro przejścia o gęstości $q(x, y)$
- π rozkład docelowy

Algorytm MH

- 1 Generujemy $Y \sim Q(X_n, \cdot)$
- 2 Obliczamy prawdopodobieństwo akceptacji $\alpha(X_n, Y)$ gdzie

$$\alpha(x, y) = \frac{\pi(y)q(y, x)}{\pi(x)q(x, y)} \wedge 1$$

3

$$X_{n+1} = \begin{cases} Y & \text{z prawdopodobieństwem } \alpha(X_n, Y) \\ X_n & \text{z prawdopodobieństwem } 1 - \alpha(X_n, Y) \end{cases}$$

Algorytm Metropolisa-Hastingsa

Algorytm MH generuje łańcuch Markowa odwracalny względem π Czyli

$$P_{MH}(x, y)\pi(x) = P_{MH}(y, x)\pi(y)$$

Jeśli $x \neq y$ to

$$P_{MH}(x, y) = \alpha(x, y)q(x, y)$$

Założmy, że $\alpha(x, y) < 1$ czyli $\alpha(x, y) = \frac{\pi(y)q(y, x)}{\pi(x)q(x, y)}$ to $\alpha(y, x) = 1$
(drugi przypadek symetryczny)

$$\begin{aligned} P_{MH}(x, y)\pi(x) &= \alpha(x, y)q(x, y)\pi(x) = \pi(y)q(y, x) \\ &= \pi(y)q(y, x)\alpha(y, x) = P_{MH}(y, x)\pi(y) \end{aligned}$$

Przykłady MH I

- Random Walk Metropolis $q(x, y) = q(|x - y|)$ np.
 $Q(x, dy) \sim N(x, \Sigma)$ wówczas

$$\alpha(x, y) = \frac{\pi(y)}{\pi(x)} \wedge 1$$

- Independent MH $q(x, y) = q(y)$

$$\alpha(x, y) = \frac{\pi(y)q(x)}{\pi(x)q(y)} \wedge 1$$

Przykłady MH II

- Metropolis Adjusted Langevin Algorithm

$$dX_t = \nabla \log(\pi(X_t))dt + \sqrt{2}dW_t$$

Dyskretyzacja Eulera- Maruyamy

$$X_{k+1} = X_k + \tau \nabla \log(\pi(X_k)) + \sqrt{2\tau}Z_k,$$

gdzie $Z_k \sim N(0, I)$.

MALA

$$Y = X_k + \tau \nabla \log(\pi(X_k)) + \sqrt{2\tau}Z_k$$

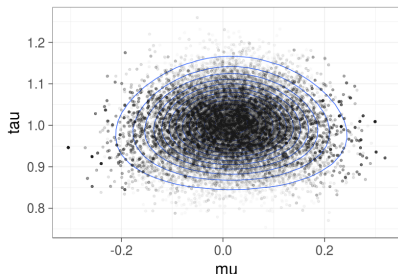
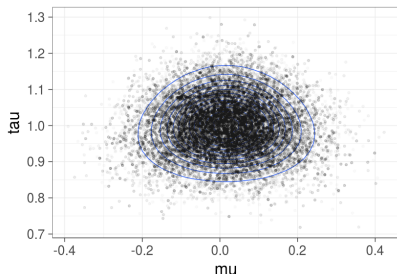
oraz

$$\alpha(x, y) = \frac{\pi(y)\phi(x; y + \tau \nabla \log(\pi(y)); 2\tau I)}{\pi(x)\phi(y; x + \tau \nabla \log(\pi(x)); 2\tau I)} \wedge 1$$

Przykład

X_1, \dots, X_n i.i.d $N(\mu, 1/\tau)$,

$$\pi(\mu) \sim N(0, 10), \quad \pi(\tau) \sim \text{Gamma}(1, 1)$$



	RWM	MALA
μ	-0.0013 (\pm 0.094)	0.0630 (\pm 0.150)
τ	1.0000 (\pm 0.070)	0.9700 (\pm 0.060)

2 Ukryte modele Markowa

- Dyskretne modele Markowa
- Obliczanie prawdopodobieństwa obserwacji
- Forward Backward
- Algorytm Viterbiego
- Algorytm FFBS
- Algorytm Bauma-Welcha
- Modele gaussowskie i filtr Kalmana

3 Algorytmy MCMC

- Podstawowe pojęcia statystyki bayesowskiej
- Metody Monte Carlo
- Algorytm Metropolisa-Hastingsa
- Próbnik Gibbsa
- Własności teoretyczne MCMC
- Adaptacyjne algorytmy MCMC
- Hamiltonian Monte Carlo

Próbnik Gibbsa I

Mały krok próbnika Gibbsa:

$$X_i \sim \pi(X_i | X_{-i})$$

Oznaczmy przez $P_i(x, dy)$ jego jądro przejścia, jego gęstość

$$p_i(x, y) = \begin{cases} \pi(y_i | x_{-i}) & \text{jeśli } x_{-i} = y_{-i} \\ 0 & \text{w p.p.} \end{cases}$$

Systematic Scan

Wykonujemy kroki P_i w ustalonej kolejności czyli np.

$$X_1^n \sim \pi(\cdot | X_{-1}^{n-1})$$

$$X_2^n \sim \pi(\cdot | X_1^n, X_2^{n-1}, \dots, X_d^{n-1})$$

\vdots

$$X_k^n \sim \pi(\cdot | X_1^n, \dots, X_{k-1}^n, X_{k+1}^{n-1}, \dots, X_d^{n-1})$$

\vdots

Próbnik Gibbsa II

Random Scan

W każdym kroku losujemy współrzędna która ma być zmieniana

- 1 Losujemy $j \sim \text{Unif}\{1, \dots, d\}$
- 2

$$X_j^n \sim \pi(\cdot | X_{-j}^{n-1})$$

$$X_k^n = X_k^{n-1}, \quad \text{dla } k \neq j$$

Poprawność próbnika Gibbsa I

Mały krok próbnika Gibbsa jest odwracalny. Dla $\mathbf{x}_{-i} = \mathbf{y}_{-i}$

$$\begin{aligned} P_i(\mathbf{x}, \mathbf{y})\pi(\mathbf{x}) &= \pi(y_i | \mathbf{x}_{-i})\pi(\mathbf{x}_i, \mathbf{x}_{-i}) \\ &= \frac{\pi(y_i, \mathbf{x}_{-i})\pi(\mathbf{x}_i | \mathbf{x}_{-i})}{\pi(\mathbf{x}_{-i})} \\ &= \frac{\pi(y_i, \mathbf{y}_{-i})\pi(\mathbf{x}_i | \mathbf{y}_{-i})}{\pi(\mathbf{y}_{-i})} \\ &= \pi(\mathbf{x}_i | \mathbf{y}_{-i})\pi(y_i, \mathbf{y}_{-i}) = P_i(\mathbf{y}, \mathbf{x})\pi(\mathbf{y}) \end{aligned}$$

Poprawność próbnika Gibbsa II

Próbnik Gibbsa z losowym wyborem współrzędnych jest odwracalny.

$$P_{RS}(x, y) = \sum_{i=1}^d \frac{1}{n} P_i(x, y) \pi(x)$$

Każdy z P_i odwracalny to kombinacja wypukła jest odwracalna. Próbnik Gibbsa z deterministycznym wyborem współrzędnych jest nieodwracalny a zachowuje rozkład stacjonarny.

$$P_{SS}(x, y) = P_1((x_1, \dots, x_d), (y_1, x_2, \dots, x_d)) \cdots P_i((y_1, \dots, y_{i-1}, x_i, \dots, x_d), (y_1, \dots, y_i, x_{i+1}, \dots, x_d)) \cdots P_d((y_d, \dots, y_{d-1}, x_d), (y_1, \dots, y_d))$$

Poprawność próbnika Gibbsa III

$$\begin{aligned}\pi P_{SS}(y) &= \int P_{SS}(x, y) \pi(x) dx \\ &= \int P_2 \cdots P_d \int P_1((x_1, \dots, x_d), (y_1, x_2, \dots, x_d)) \pi(x) dx_1 dx_{-1} \\ &= \int P_2((y_1, x_2, \dots, x_d), (y_1, y_2, \dots, x_d)) \cdots P_d \pi((y_1, x_2, \dots, x_d)) dx_{-1} \\ &= \cdots = \pi(y)\end{aligned}$$

Przykład 1

$$\pi(\mu, \tau | X_1, \dots, X_n) \propto \tau^{a+n/2-1} \exp \left\{ -\tau \frac{\sum_{i=1}^n (X_i - \mu)^2}{2} - \frac{(\mu - m)^2}{2v^2} - \tau b \right\}$$

Pełne rozkłady warunkowe:

$$\pi(\mu | \tau, X_1, \dots, X_n) \sim N \left(\frac{n\bar{X}v^2 + m\tau^{-1}}{nv^2 + \tau^{-1}}, \frac{v^2\tau^{-1}}{nv^2 + \tau^{-1}} \right)$$

$$\pi(\tau | \mu, X_1, \dots, X_n) \sim \text{Gamma} \left(a + n/2, \frac{\sum_{i=1}^n (X_i - \mu)^2}{2} + b \right)$$

Próbnik Gibbsa dla modelu hierarchicznego

$$\pi(\mu_1, \dots, \mu_k, \tau_1, \dots, \tau_k, \mu, \tau | \mathbf{X}, \mathbf{a}, \mathbf{b}, \alpha, \lambda, \mathbf{m}, \mathbf{v}^2) \propto \prod_{j=1}^k \left\{ \prod_{i=1}^{N_j} [p(X_{ij} | \mu_j, \tau_j)] \pi(\mu_j | \mu, \tau) \pi(\tau_j | \mathbf{a}, \mathbf{b}) \right\} \pi(\mu | \mathbf{m}, \mathbf{v}^2) \pi(\tau | \alpha, \lambda)$$

Pełne rozkłady warunkowe:

$$\pi(\mu_j | \tau_1, \dots, \tau_k, \tau, \mu, \mathbf{X}_1, \dots, \mathbf{X}_n) \sim N \left(\frac{N_j \bar{X}_j \tau^{-1} + \mu \tau_j^{-1}}{N_j \tau^{-1} + \tau_j^{-1}}, \frac{\tau^{-1} \tau_j^{-1}}{N_j \tau^{-1} + \tau_j^{-1}} \right)$$

$$\pi(\tau_j | \mu_1, \dots, \mu_k, \tau, \mu, \mathbf{X}_1, \dots, \mathbf{X}_n) \sim \text{Gamma} \left(\mathbf{a} + N_j/2, \frac{\sum_{i=1}^{N_j} (X_{ij} - \mu_j)^2}{2} + \mathbf{b} \right)$$

$$\pi(\mu | \tau_1, \dots, \tau_k, \mu_1, \dots, \mu_k, \tau, \mathbf{X}_1, \dots, \mathbf{X}_n) \sim N \left(\frac{k \bar{\mu} \mathbf{v}^2 + \mathbf{m} \tau^{-1}}{k \mathbf{v}^2 + \tau^{-1}}, \frac{\mathbf{v}^2 \tau^{-1}}{k \mathbf{v}^2 + \tau^{-1}} \right)$$

$$\pi(\tau | \tau_1, \dots, \tau_k, \mu_1, \dots, \mu_k, \mu, \mathbf{X}_1, \dots, \mathbf{X}_n) \sim \text{Gamma} \left(\alpha + k/2, \frac{\sum_{i=1}^k (\mu_i - \mu)^2}{2} + \lambda \right)$$

Metropolis wewnątrz Gibbsa

Próbnik Gibbsa składa się z małych kroków postaci

$$P_i(x, dy) = \pi(dy_i | x_{-i}).$$

Może się zdarzyć, że generowanie z jednego (wielu) rozkładów warunkowych $\pi(\cdot | x_{-i})$ jest trudne. Wówczas krok P_i można zastąpić krokiem Metropolisa Hastingsa o rozkładzie docelowym $\pi(\cdot | x_{-i})$.

2 Ukryte modele Markowa

- Dyskretne modele Markowa
- Obliczanie prawdopodobieństwa obserwacji
- Forward Backward
- Algorytm Viterbiego
- Algorytm FFBS
- Algorytm Bauma-Welcha
- Modele gaussowskie i filtr Kalmana

3 Algorytmy MCMC

- Podstawowe pojęcia statystyki bayesowskiej
- Metody Monte Carlo
- Algorytm Metropolisa-Hastingsa
- Próbnik Gibbsa
- Własności teoretyczne MCMC
- Adaptacyjne algorytmy MCMC
- Hamiltonian Monte Carlo

Potrzebujemy następujących własności łańcuchów Markowa

- Ergodyczność $P(X_n \in A | X_1 = x) \rightarrow \pi(A)$
- Prawo wielkich liczb $\hat{I}_n = \frac{1}{n} \sum_{i=1}^n f(X_i) \rightarrow \pi f$
- Centralne twierdzenie graniczne $\sqrt{n}(\hat{I}_n - \pi f) \rightarrow N(0, \sigma^2)$

Ergodyczność dla skończonej przestrzeni stanów

TWIERDZENIE

Jeśli łańcuch Markowa o rozkładzie stacjonarnym π jest nieprzywiedlny i nieokresowy to jest ergodyczny.

π -nieprzywiedlność

Definicja

Łańcuch Markowa nazywamy π -nieprzywiedlnym, jeżeli dla każdego zbioru A , takim że $\pi(A) > 0$, oraz dla każdego x istnieje n

$$P^n(x, A) > 0$$

Jeśli łańcuch Markowa jest π nieprzywiedlny to rozkład stacjonarny jest wyznaczony jednoznacznie.

Ergodyczność I

Definicja

Zbiór α nazywamy atomem jeżeli $\pi(\alpha) > 0$ oraz

$$\forall x \in \alpha, \quad P(x, A) = \nu(A).$$

Definicja

Dla miar probabilistycznych μ, ν odległością pełnego wahania nazywamy

$$\|\mu - \nu\|_{\text{tv}} = \sup_A |\mu(A) - \nu(A)|$$

Jeśli $X \sim \mu$, $Y \sim \nu$ to

$$\|\mu - \nu\|_{\text{tv}} \leq P(X \neq Y)$$

Ergodyczność II

TWIERDZENIE

Jeśli łańcuch Markowa o rozkładzie stacjonarnym π jest π nieprzywydlny i nieokresowy to

$$\|P^n(x, \cdot) - \pi(\cdot)\|_{tv} \rightarrow 0$$

Dowód:

Najpierw założymy, że łańcuch posiada atom α oraz zdefiniujemy parę łańcuchów Markowa X_n, \tilde{X}_n : $X_1 = x, \tilde{X}_1 \sim \pi$. Jeżeli $X_n \neq \tilde{X}_n$ to łańcuchy generowane są niezależnie zgodnie z $P(x, dy)$, natomiast jeśli $X_n = \tilde{X}_n$ to dalej trajektorie pozostają sklezione. Wówczas

$$\|P^n(x, \cdot) - \pi(\cdot)\|_{tv} \leq P(X_n \neq \tilde{X}_n) \leq 1 - P(\tau_{\alpha \times \alpha} \leq n),$$

gdzie

$$\tau_{\alpha \times \alpha} = \min\{k \geq 1 : X_k = \tilde{X}_k = \alpha\}.$$

Ergodyczność III

Z nieprzywiedlności i nieokresowości $P(\tau_{\alpha \times \alpha} \leq n) \rightarrow 1$

Jeżeli łańcuch nie posiada atomu to konstruujemy sztuczny atom:

Definicja

Zbiór C nazywamy małym zbiorem jeżeli istnieje n takie, że

$$P^n(x, A) \geq \beta \nu(A), \quad \forall x \in C$$

Konstrukcja rozszczepienia łańcucha dla $n = 1$. Definiujemy nowy łańcuch Markowa X_n, Γ_n na przestrzeni $X \times \{0, 1\}$

$$P(X_{n+1} \in A | X_n = x, \Gamma_n = 0) = \frac{P(x, A) - \beta \nu(A) 1(x \in C)}{1 - \beta 1(x \in C)}$$

$$P(X_{n+1} \in A | X_n = x, \Gamma_n = 1) = \nu(A)$$

$$P(\Gamma_{n+1} = 1 | X_{n+1} = x) = \beta 1(x \in C)$$

Tak skonstruowany łańcuch ma atom $C \times \{1\}$ oraz rozkład brzegowy X_n jest taki sam jak pierwotnego łańcucha.

Prawo wielkich liczb I

TWIERDZENIE

Jeśli łańcuch Markowa z rozkładem stacjonarnym π jest π -nieprzywiedlny oraz posiada atom α to

$$\pi(A) \propto \mathbb{E}_\alpha \sum_{k=1}^{\tau_\alpha-1} 1(X_k \in A),$$

gdzie

$$\tau_\alpha = \min\{k > 1 : X_k \in \alpha\}$$

Dowód dla α jednoelementowego

Prawo wielkich liczb II

Oznaczmy przez $\nu(A)$ prawą stronę. Najpierw założmy, że $A \cap \alpha = \emptyset$

$$\begin{aligned}\nu P(A) &= \int_X \mathbb{E}_\alpha \sum_{k=1}^{\tau_\alpha - 1} 1(X_k \in dx) P(x, A) \\ &= \int_X \sum_{k=1}^{\infty} P_\alpha(X_k \in dx, \tau_\alpha > k) P(x, A) \\ &= \sum_{k=1}^{\infty} \int_X P_\alpha(X_k \in dx, \tau_\alpha > k) P(x, A) \\ &= \sum_{k=1}^{\infty} P_\alpha(X_{k+1} \in A, \tau_\alpha > k + 1) \\ &= \nu(A)\end{aligned}$$

Prawo wielkich liczb III

$$\begin{aligned}\nu P(\alpha) &= \int_{\mathbf{X}} \mathbb{E}_{\alpha} \sum_{k=1}^{\tau_{\alpha}-1} 1(X_i \in dx) P(x, \alpha) \\ &= \int_{\mathbf{X}} \sum_{k=1}^{\infty} P_{\alpha}(X_k \in dx, \tau_{\alpha} > k) P(x, \alpha) \\ &= \sum_{k=1}^{\infty} \int_{\mathbf{X}} P_{\alpha}(X_k \in dx, \tau_{\alpha} > k) P(x, \alpha) \\ &= \sum_{k=1}^{\infty} P_{\alpha}(X_{k+1} = \alpha, \tau_{\alpha} = k + 1) \\ &= 1 = \nu(\alpha)\end{aligned}$$

Prawo wielkich liczb IV

TWIERDZENIE

Jeśli πf istnieje oraz π oraz łańcuch jest π -nieprzywiedlny to

$$\frac{\sum_{k=1}^n f(X_k)}{n} \rightarrow \pi f, \quad \text{p.n.}$$

Niech T_i oznaczają kolejne czasy powrotu do α . Wówczas kawałki trajektorii pomiędzy kolejnymi czasami T_i

$$\Xi_k(f) = \sum_{i=T_{k-1}}^{T_k-1} f(X_i)$$

są niezależne i z wyjątkiem pierwszej mają ten sam rozkład.

Niech

$$T_{N(n)} \leq n \leq T_{N(n)+1}$$

Z PWL

$$T_k/k \rightarrow \nu(X)$$

Prawo wielkich liczb V

Z twierdzenia o trzech ciągach

$$\frac{n}{N(n)} \rightarrow \nu(X)$$

Dodatkowo jeśli $f \geq 0$ to

$$\sum_{k=1}^{N(n)} \Xi_k(f) \leq \sum_{i=1}^n f(X_i) \leq \sum_{k=1}^{N(n)+1} \Xi_k(f)$$

Ponownie z twierdzenia o trzech ciągach i PWL dla i.i.d

$$\frac{\sum_{k=1}^n f(X_i)}{n} \rightarrow \frac{\nu f}{\nu(X)} = \pi f$$

Centralne twierdzenie graniczne I

Centralne twierdzenie graniczne również można dowodzić techniką regeneracyjną. Jednak w tym przypadku, założenia oraz postać asymptotycznej wariancji łatwiej wyjaśnić korzystając z techniki martyngałowej.

Równanie Poissona Funkcja \hat{f} spełnia równanie Poissona jeśli

$$\hat{f}(x) - P\hat{f}(x) = f(x) - \pi f$$

Jeśli \hat{f} rozwiązaniem równania Poissona to

$$\begin{aligned} \sum_{k=1}^n f(X_k) - \pi f &= \sum_{k=1}^n \hat{f}(X_k) - P\hat{f}(X_k) \\ &= \sum_{k=2}^n \hat{f}(X_k) - P\hat{f}(X_{k-1}) + \hat{f}(X_1) + P\hat{f}(X_n) \end{aligned}$$

Ostatnie kawałki po podzieleniu przez \sqrt{n} są pomijalne a pierwszy tworzy martyngał. I chcemy skorzystać z CTG dla

Centralne twierdzenie graniczne II

martyngałów. Dodatkowo można pokazać, że CTG nie zależy od rozkładu początkowego, można założyć $X_1 \sim \pi$

$$\mathbb{E}((\hat{f}(X_i) - P\hat{f}(X_{i-1}))^2 | X_{i-1}) = P\hat{f}^2(X_{i-1}) - (P\hat{f}(X_{i-1}))^2$$

Stąd

$$\sigma^2 = \pi(\hat{f}^2 - (P\hat{f})^2)$$

Oraz

$$\frac{1}{\sqrt{n}} \sum_{k=1}^n f(X_k) - \pi f \rightarrow N(0, \sigma^2)$$

Do precyzyjnego sformułowania będziemy potrzebowali geometrycznej ergodyczności. Najpierw jednak przyjmiemy się asymptotycznej wariancji.

Asymptotyczna wariancja I

Jedynym kandydatem na rozwiązanie równania Poissona jest

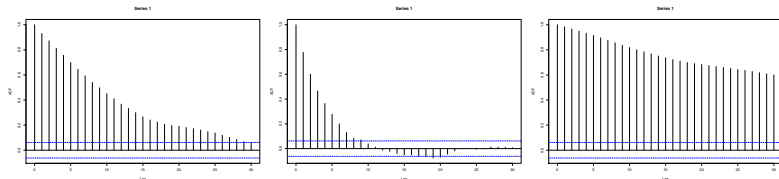
$$\hat{f} = \sum_{k=0}^{\infty} P^k f(x) - \pi f$$

Oznaczmy $\bar{f} = f - \pi f$ wówczas

$$\begin{aligned}\sigma^2 &= \pi(\hat{f}^2 - (P\hat{f})^2) \\ &= \pi(\hat{f}^2 - (\hat{f} - \bar{f})^2) \\ &= \pi(2\bar{f}\hat{f} - \bar{f}^2) \\ &= 2\pi\left(\sum_{k=0}^{\infty} \bar{f}P^k\bar{f} - \bar{f}^2\right) \\ &= \pi(\bar{f}^2) + \sum_{k=1}^{\infty} \pi(\bar{f}P^k\bar{f}) \\ &= \text{Var}_{\pi}(f(X)) + \sum_{k=2}^{\infty} \text{cov}_{\pi}(X_1, X_k)\end{aligned}$$

Asymptotyczna wariancja II

Stąd ważna w diagnostyce MCMC jest funkcja autokorelacji.



Tempo zbieżności

Jednostajna ergodyczność:

$$\|P^n(x, \cdot) - \pi(\cdot)\|_{\text{tv}} \leq M\rho^n, \quad \rho < 1$$

Geometryczna ergodyczność:

$$\|P^n(x, \cdot) - \pi(\cdot)\|_{\text{tv}} \leq M(x)\rho^n, \quad \rho < 1$$

Jednostajna ergodyczność

TWIERDZENIE

Jednostajana ergodyczność jest równoważna warunkowi, że cała przestrzeń jest małym zbiorem czyli

$$P^n(x, A) \geq \beta \mu(A)$$

$$P(X_{kn} \neq Y_{kn}) \leq P(\tau > kn) = (1 - \beta^2)^k = \rho^{kn},$$

gdzie $\rho = (1 - \beta^2)^{1/n}$

Geometryczna ergodyczność

TWIERDZENIE

Łańcuch Markowa jest geometrycznie ergodyczny wtedy i tylko wtedy gdy istnieje mały zbiór C

$$P^n(x, A) \geq \beta \mu(A) 1(x \in C)$$

oraz istnieje funkcja $V \geq 1$ taka, że

$$PV(x) \leq \lambda V(x) + K 1(x \in C)$$

dla $\lambda < 1$

Powrót do CTG

CTG jest spełnione jeżeli łańcuch jest geometrycznie ergodyczny oraz jeden z warunków jest spełniony:

- łańcuch jest odwracalny i $\pi(f^2) < \infty$
- $\pi(|f|^{2+\delta}) < \infty$
- $\sup_x \frac{f^2(x)}{V(x)} < \infty$

2 Ukryte modele Markowa

- Dyskretne modele Markowa
- Obliczanie prawdopodobieństwa obserwacji
- Forward Backward
- Algorytm Viterbiego
- Algorytm FFBS
- Algorytm Bauma-Welcha
- Modele gaussowskie i filtr Kalmana

3 Algorytmy MCMC

- Podstawowe pojęcia statystyki bayesowskiej
- Metody Monte Carlo
- Algorytm Metropolisa-Hastingsa
- Próbnik Gibbsa
- Własności teoretyczne MCMC
- Adaptacyjne algorytmy MCMC
- Hamiltonian Monte Carlo

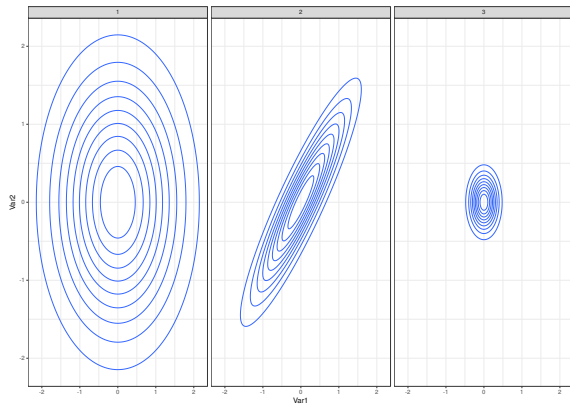
Problemy z algorytmem Metropolisa

RWM:

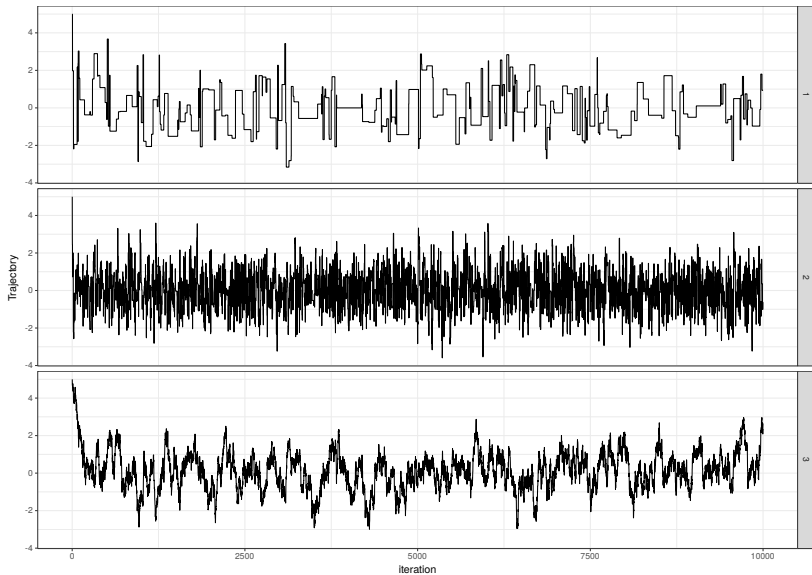
Propozycja $Y \sim N(X_n, \Sigma)$

Jak wybrać Σ ?

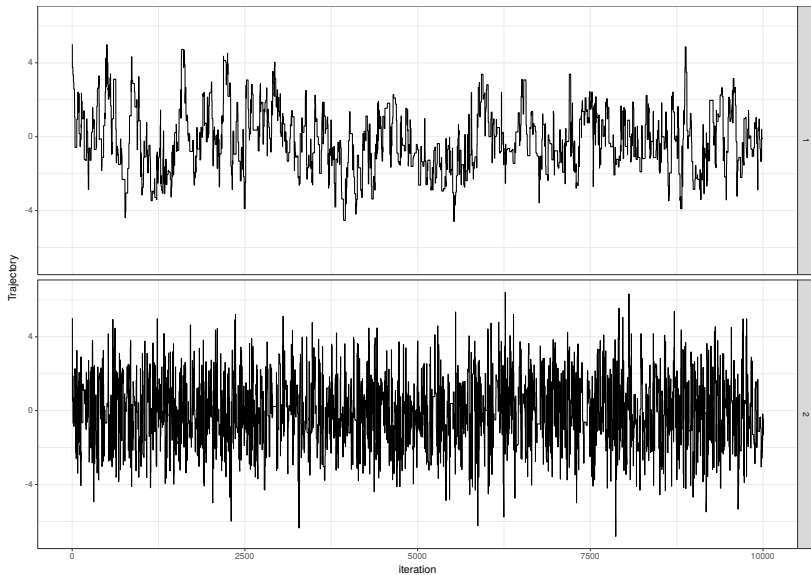
- Wielkość.
- Kształt.



Wielkość



Kształt



“Optymalny” wybór Σ

- Optymalana $\Sigma = \ell^2 \text{Cov}_\pi(\mathbf{X})$, gdzie ℓ^2 można zdefiniować na dwa sposoby.
- $\ell^2 \approx \frac{2.38^2}{d}$, gdzie d wymiar.
- ℓ^2 takie, że średnie prawdopodobieństwo akceptacji wynosi 0.234

Uzasadnienie teoretyczne I

Jeśli rozpatrzy się

$$\pi(\mathbf{x}) = \prod_{i=1}^d f(x_i)$$

To można pokazać że dla $\Sigma = \frac{\ell^2}{d}I$ łańcuch Markowa generowany przez algorytm Metropolis'a w poprzekalowaniu zbiega wraz z wymiarem do dyfuzji. Dokładniej

$$X_{\lfloor dt \rfloor} \rightarrow Y_t$$

gdzie

$$dY_t = \sqrt{h(l)}dW_t - h(l)\frac{1}{2}V'(Y_t)dt,$$

$$V(x) = -\log(f(x)), \quad h(l) = 2l^2\Phi\left(-\frac{l\sqrt{I}}{2}\right), \quad I = \int (V')^2 \exp(-V).$$

Uzasadnienie teoretyczne II

Dodatkowo prawdopodobieństwo akceptacji zbiega do

$$a(l) = 2\Phi\left(-\frac{\ell\sqrt{I}}{2}\right)$$

Funkcja $h(l)$ odpowiada za szybkość dyfuzji i chcemy ją zmaksymalizować.

$$\ell^* = \frac{2.38}{\sqrt{I}} \text{ oraz } a(\ell^*) = 0.234$$

Adaptacyjny RWM

Chcemy skonstruować algorytm, który automatycznie znajduje optymalną macierz Σ . W tym celu w każdym kroku algorytmu będziemy estymować macierz kowariancji oraz odpowiednio ją skalować. niech

$$\mu_n = \frac{n-1}{n}\mu_{n-1} + \frac{1}{n}X_n$$

$$\Gamma_n = \frac{n-1}{n}\Gamma_{n-1} + \frac{1}{n}(X_n - \mu_n)(X_n - \mu_n)^T$$

Wersja I:

$$\Sigma_n = \frac{2.38^2}{d}\Gamma_n + \epsilon I$$

Wersja II:

$$T_n = T_{n-1} + \frac{1}{n}(\alpha(X_n, Y_n) - 0.234)$$

$$\Sigma_n = \exp(T_n)\Gamma_n + \epsilon I$$

Uwagi o implementacji adaptacyjnego RWM I

Generowanie wielowymiarowego rozkładu gaussowskiego

Jeśli $X = (X_1, \dots, X_d)$ i.i.d $N(0,1)$ oraz $\Sigma = LL^T$ to

$$LX \sim N(0, \Sigma)$$

Typowo używa się rozkładu Choleskiego. Dla którego istnieje one-rank-update Jeśli

$$M_{n+1} = M_n + xx^T$$

To odpowiadające im rozkłady Choleskiego

$$L_{n+1} = \text{ORU}(L_n, x).$$

W adaptacyjnym RWM wystarczy pamiętać L_n rozkład Choleskiego Γ_n i

$$L_{n+1} = \text{ORU}\left(\frac{\sqrt{n}}{\sqrt{n+1}}L_n, \frac{1}{\sqrt{n+1}}(X_{n+1} - \mu_{n+1})\right)$$

Uwagi o implementacji adaptacyjnego RWM II

Generowanie propozycji:

$$Y_n = X_n + \exp(T_n/2)L_n Z_n + \sqrt{\epsilon}Z'_n$$

gdzie Z_n i Z'_n niezależne $N(0, I)$.

2 Ukryte modele Markowa

- Dyskretne modele Markowa
- Obliczanie prawdopodobieństwa obserwacji
- Forward Backward
- Algorytm Viterbiego
- Algorytm FFBS
- Algorytm Bauma-Welcha
- Modele gaussowskie i filtr Kalmana

3 Algorytmy MCMC

- Podstawowe pojęcia statystyki bayesowskiej
- Metody Monte Carlo
- Algorytm Metropolisa-Hastingsa
- Próbnik Gibbsa
- Własności teoretyczne MCMC
- Adaptacyjne algorytmy MCMC
- Hamiltonian Monte Carlo

Układ Hamiltonowski

$$H(p, x) = K(p) + U(x)$$

K - energia kinetyczna, U- energia potencjalna Układ Hamiltonowski

$$\frac{dx_i}{dt} = \frac{\partial H}{\partial p_i}$$
$$\frac{dp_i}{dt} = -\frac{\partial H}{\partial x_i}$$

Rozwiązanie zachowuje energię.

Związek z rozkładami prawdopodobieństwa

Rozkład Gibbsa

$$\pi(p, x) \propto \exp(-H(p, x))$$

Rozkład brzegowy dla dowolnego K

$$\pi(x) \propto \exp(-U(x))$$

Rozkład warunkowych

$$\pi(p|x) \propto \exp(-K(p)) = \pi(p)$$

Idealne Hamiltonowskie Monte Carlo

$$X_n \rightarrow X_{n+1}$$

- Wylosuj $P \sim \pi(p)$.
- Zgodnie z Hamiltonowską dynamiką

$$(P, X_n) = (x_0, p_0) \rightarrow (x_T, p_T) = (X_{n+1}, P')$$

Potrzeba K zadające rozkład łatwy do generowania, oraz
potrzeba umieć obliczyć rozwiązanie układu hamiltonowskiego
w czasie T .

Hamiltonian MCMC I

$$K(p) = \frac{1}{2}p^T p$$

Czyli $\pi(p) = N(0, I)$. Dodatkowo potrzeba numerycznie rozwiązywać układ hamiltonowski.

Leapfrog integrator:

Dla $1 \leq n \leq T/\epsilon$

$$p_{n+1/2} = p_n - \frac{\epsilon}{2} \frac{\partial U}{\partial x}(x_n)$$

$$x_{n+1} = x_n + \epsilon p_{n+1/2}$$

$$p_{n+1} = p_{n+1/2} - \frac{\epsilon}{2} \frac{\partial U}{\partial x}(x_{n+1})$$

Hamiltonian MCMC II

Jeżeli zamienimy znak $p_T \rightarrow -p_T$ to leapfrog integrator jest odwracalny. I żeby usunąć błąd numerycznego rozwiązania wystarczy dodać akceptację Metropolis

$$\alpha((x_0, p_0), (x_T, -p_T)) = \exp(H(x_T, p_T) - H(x_0, p_0)) \wedge 1$$

Hamiltonian MCMC III

Algorytm

- 1 Generujemy $p \sim N(0, I)$
- 2 Przybliżamy rozwiązanie układu hamiltonowskiego (leapfrog integrator) $X_n = x_0, p = p_0 \rightarrow x_T, p_T$
- 3

$$X_{n+1} = \begin{cases} x_T & \text{z prawdopodobieństwem } \alpha((x_0, p_0), (x_T, -p_T)) \\ X_n & \text{z prawdopodobieństwem } 1 - \alpha((x_0, p_0), (x_T, -p_T)) \end{cases}$$