

Centrální limitní věta a chyba metody Monte-Carlo

- metoda Monte-Carlo spočívá ve výpočtu určité veličiny (např. určitého integrálu) použí průměrování náhodných veličin X_i s určitou pravděpodobností i rozdělení. Příjemně lze očekávat, že chyba bude s rostoucím počtem n náhodných veličin X_i klesat jako $\frac{1}{\sqrt{n}}$. To souvisí s centrální limitní větou z teorie pravděpodobnosti.
- základní pojmy
 - pravděpodobností funkce náhodné diskrétní veličiny X
 $P[X=x] = P(x)$ udává pravděpodobnost, že náhodná veličina X bude mit hodnotu x
musí tedy být $P(x) \geq 0$ pro $\forall x$ a $\sum_x P(x) = 1$
 - distribuční funkce náhodné diskrétní veličiny X
 $F(x) = P[X \leq x] = \sum_{t \leq x} P(t)$, $P[x_1 \leq X \leq x_2] = F(x_2) - F(x_1)$
jde o neklesající funkci a $0 \leq F(x) \leq 1$
 - rozdělení pravděpodobnosti spojité náhodné veličiny X
se určuje použí hustoty pravděpodobnosti $g(x) \geq 0$
splňující $\int_R g(x) dx = 1$, kde R je definiční obor X
(obvykle určitý interval)
pravděpodobnost, že X bude mit hodnotu v $\langle x_1, x_2 \rangle$ je
 $P[x_1 \leq X \leq x_2] = \int_{x_1}^{x_2} g(x) dx = F(x_2) - F(x_1)$
kde $F(x)$ je distribuční funkce spojité náh. veličiny X
 $F(x) = P[X \leq x] = \int_{-\infty}^x g(t) dt$ a tedy $g(x) = \frac{dF(x)}{dx}$
pravděpodobnost nalezení přesné hodnoty x pro
spojitou náhodnou veličinu X je nula.

- střední hodnota (užívaný průměr daného rozdělení)

$$E(X) = \sum_{i=1}^n x_i p_i \quad \text{pro diskrétní náhodnou veličinu}$$

$$= \int x g(x) dx \quad \text{pro spojitou}$$

- rozptyl, neboli střední kvadratická odchylka (též variansa)

$$\begin{aligned}\sigma^2(X) &= D(X) = \text{var}(X) = \text{člen } -2x_i E(x) \text{ dle } -2E(x)^2 \\ &= \sum_{i=1}^n [x_i - E(X)]^2 p_i = \sum_i x_i^2 p_i - E(X)^2 \quad (\text{diskrétní}) \\ &= \int [x - E(X)]^2 g(x) dx = \int x^2 g(x) dx - E(X)^2 \\ &= E(X^2) - E(X)^2\end{aligned}$$

základní vlastnosti (pro dvě nezávislé veličiny X a Y)

$$E(X+c) = c + E(X)$$

$$D(X+c) = D(X)$$

$$E(cX) = c E(X)$$

$$D(cX) = c^2 D(X)$$

$$E(X+Y) = E(X) + E(Y)$$

$$D(X+Y) = D(X) + D(Y)$$

$$E(X \cdot Y) = E(X) \cdot E(Y)$$

napr. $D(X+Y) = \iint (x+y - E(X+Y))^2 g(x,y) dx dy =$ pro nezávislé X a Y

$$\begin{aligned}&= \iint (x-E(X) + y - E(Y))^2 g_x(x) g_y(y) dx dy = \\&= \iint [(x-E(X))^2 + 2\underbrace{(x-E(X))(y-E(Y))}_{0} + (y-E(Y))^2] g_x(x) g_y(y) dx dy \\&= D(X) + D(Y)\end{aligned}$$

- střední hodnota funkce náhodné veličiny

pokud $Y=f(X)$ je náhodná veličina, která je funkcí náhodné veličiny X , pak

$$E(Y) = E(f(X)) = \int f(x) g(x) dx \quad (\text{nebo } = \sum_i f(x_i) p_i)$$

a obecně $E(f(X)) \neq f(E(X))$

příklady rozdělení

- rovnoměrné rozdělení na intervalu $\langle a, b \rangle$

$$g(x) = \begin{cases} \frac{1}{b-a} & \text{na } \langle a, b \rangle \\ 0 & \text{jinde} \end{cases} \quad (\text{pro interval } \langle 0, 1 \rangle \text{ je } g(x)=1)$$

$$E(X) = \frac{a+b}{2}, \quad D(X) = \sigma^2(X) = \frac{(b-a)^2}{12}, \quad F(X) = \frac{x-a}{b-a} \text{ na } \langle a, b \rangle$$

- normální rozdělení $N(\mu, \sigma^2)$ na $\langle -\infty, \infty \rangle$

$$g(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad E(X) = \mu, \quad D(X) = \sigma^2$$

distribuční funkce

$$F(x) = \int_{-\infty}^x g(t) dt = \frac{1}{2} \left[1 + \operatorname{Erf} \left(\frac{x-\mu}{\sqrt{2\sigma^2}} \right) \right]$$

kde $\operatorname{Erf}(z) = \frac{2}{\sqrt{\pi}} \int_0^z e^{-t^2} dt$ (error function)

a tedy $P[x_1 \leq X \leq x_2] = \int_{x_1}^{x_2} g(t) dt = \frac{1}{2} \left[\operatorname{Erf} \left(\frac{x_2-\mu}{\sqrt{2\sigma^2}} \right) - \operatorname{Erf} \left(\frac{x_1-\mu}{\sqrt{2\sigma^2}} \right) \right]$

pravidlo 3σ

$$P[\mu - 3\sigma \leq X \leq \mu + 3\sigma] = \operatorname{Erf}\left(\frac{3}{\sqrt{2}}\right) = 0,9973$$

a tedy náhodné měření téměř vždy leží v intervalu $\langle \mu - 3\sigma, \mu + 3\sigma \rangle$

pravděpodobná chyba náhodné veličiny X

- máme-li tedy jedné měření, pak s velkou pravděpodobností bude v intervalu $\langle \mu - 3\sigma, \mu + 3\sigma \rangle$ a absolutní chyba (od střední hodnoty μ) bude rádu $r \approx 0,67\sigma$

kde r je dano $P[\mu - r \leq X \leq \mu + r] = 0,5$

r se nazývá pravděpodobná chyba měření X

• centrální limitní věta

- máme-li n nezávislých náhodných veličin X_1, \dots, X_n

které mají stejné rozdělení pravděpodobnosti

tj. jejich střední hodnoty $E(X_i) = m$ pro $i=1, \dots, n$

a rozptyly $D(X_i) = b^2$, $i=1, \dots, n$ jsou stejné

pak pro jejich součet $S_n = X_1 + \dots + X_n$ (opět náhodná veličina)

platí (viz $E(X+Y) = E(X) + E(Y)$ a $D(X+Y) = D(X) + D(Y)$)

$$E(S_n) = nm \quad \text{a} \quad D(S_n) = nb^2$$

- nějme dálé normální rozdělení $N(\mu, \sigma^2)$

kde $\mu = nm$ a $\sigma^2 = nb^2$ s hustotou pravděp. $g_n(x)$

- pak centrální limitní věta říká, že

$$P[S_1 \leq S_n \leq S_2] \approx \int_{S_1}^{S_2} g_n(x) dx$$

pro libovolné S_1 a S_2 a vžecna dostatečně velká n .

- neboli rozdělení součtu S_n velkého množství nezávislých náhodných veličin je přibližně normální

Pozn: existují zájemnější pro mnohem slabší podmínky

npr. veličiny X_1, \dots, X_n nemusí být nutně nezávislé

a stejného rozdělení

důležité je, aby jedna z veličin nehrála příliš důležitou roli

díky tomu se normální rozdělení často objevuje v přírodě

na výsledek má obvykle vliv řada náhodných faktorů

- to, že dostaneme normální rozdělení se $\sigma = \sqrt{n} b$

je důležité pro odhad chyby veličiny S_n v závislosti na n , pravděpodobná chyba S_n bude jiná než \sqrt{n} , a také pro odhad chyby MC

• odhad chyby metody Monte-Carlo

- počítáme-li určitou veličinu jako průměr náhodných veličin X_1, \dots, X_n , tj. $A_n = (X_1 + \dots + X_n)/n$,

se středními hodnotami $E(X_i) = m$ a rozptyly $D(X_i) = b^2$

pak pravděpodobností rozdělení A_n bude přibližně

$$\text{normalní s } E(A_n) = \sum_{i=1}^n \frac{E(X_i)}{n} = m$$

$$\text{a } D(A_n) = \sum_{i=1}^n \frac{D(X_i)}{n^2} = \frac{b^2}{n} \text{ neboli } \sigma = \frac{b}{\sqrt{n}}$$

a tedy $P[m - \frac{3b}{\sqrt{n}} \leq A_n \leq m + \frac{3b}{\sqrt{n}}] \approx 0,9973$

neboli $P\left[\left|\frac{1}{n} \sum_{i=1}^n X_i - m\right| < \frac{3b}{\sqrt{n}}\right] \approx 0,9973$

pravděpodobná chyba A_n tedy je přibližně $0,67 \sigma =$

$$= 0,67 \frac{b}{\sqrt{n}} \text{ a lze jí tedy očekávat rádu } \frac{1}{\sqrt{n}}$$

- pozor! nejde o omezení chyby shora, jde o pravděpodobnou chybu, chyba může být ve skutečnosti větší, ale pravděpodobnost toho je malá

Integrace pomocí metody Monte-Carlo

- jak si ukázalo má tato metoda předešlý velký význam

pro vícerozměrné integrály, protože její chyba nezávisí na dimenzi, ale pouze na počtu použitých bodů N ,

tj. chyba lze rády očekávat $\sim \frac{1}{\sqrt{N}}$ (zde budeme používat N pro celkový)

- pro $d \leq 4$ (dimenze integrálů) je obvykle mnohem přesnejší použít kvadraturní vzorce, např. Gaussovy kvadratury nebo Newton-Cotesovy vzorce

- pro jednoduchost si však ilustruje tuto metodu na jednorozněrném integrálu

$$I = \int_0^1 f(x) dx$$

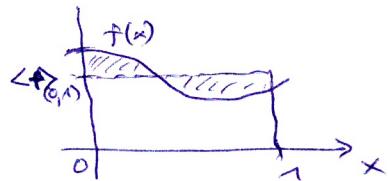
Pozn.: máme-li interval $\langle a, b \rangle$, lze v 1D použít lineární mapování $\int_a^b f(x) dx = (b-a) \int_0^1 \tilde{f}(y) dy$

kde $\tilde{f}(y) = f[(b-a)y - a]$

- tento integrál můžeme považovat za průměr funkce $f(x)$

na intervalu $\langle 0, 1 \rangle$ a approximovat

to pomocí $I \approx \frac{1}{N} \sum_{i=1}^N f(x_i)$



kde x_i jsou „rovnoměrné“ (zath.) rozložené v intervalu $\langle 0, 1 \rangle$, tj. jde o náhodné veličiny X_i s rovnoměrným rozdělením na intervalu $\langle 0, 1 \rangle$

- $f(X_i)$ je tedy náhodná veličina s určitým rozdělením, její střední hodnota je ovšem

$$E(f(X_i)) = \int_0^1 f(x) g(x) dx = \int_0^1 f(x) dx = I$$

a rozptyl σ_f je tím menší, čím je $f(x)$ „hladší“, nevisejí (pro konstantní $f(x)$ je $\sigma_f = 0$)

a díky centrální limithní větě dostaneme pro odhad

chyby $\sigma_I \approx \frac{\sigma_f}{\sqrt{N}}$

- v metódě MC se σ_f odhaduje pomocí

$$\sigma_f^2 \approx \langle f^2 \rangle - \langle f \rangle^2 = \frac{1}{N} \sum_{i=1}^N f(x_i)^2 - \left(\frac{1}{N} \sum_{i=1}^N f(x_i) \right)^2$$

- porovnání s khaběžníkovým pravidlem s chybou $O(b^3) \sim O(n^{-2})$ v 1D

dimenze	chyba	celkový počet bodů
1	$\sim n^{-2} \approx N^{-2}$	$N=n$
2	$\sim n^{-2} \approx N^{-1}$	$N=n^2$
3	$\sim n^{-2} \approx N^{-2/3}$	$N=n^3$
4	$\sim n^{-2} \approx N^{-1/2}$	$N=n^4$

srovnatelné s MC

n je počet bodů kladatury
v jedné dimenzi

N je celkový počet bodů,
kde je třeba vypočítat $f(x)$

- ve více rozměrném případu ovšem musíme obdržet složitější hranice oblasti, zde naša metoda Monte Carlo velkou výhodu, neboť stačí pouze určit, zda zvolený bod patří do oblasti přes kterou integrujeme
- protože chyba integrace pomocí metody Monte Carlo závisí na f_f , je výhodné, aby f_f byla co nejméně kladná a splňuje $\int_0^1 w(x) dx = 1$
 - výsledek lze zpravidla vypočítat vhodnou váhovou funkcí $w(x)$
 pak je kladná a splňuje $\int_0^1 w(x) dx = 1$

$$I = \int_0^1 dx w(x) \frac{f(x)}{w(x)} = \int_0^1 dy \frac{f(x(y))}{w(x(y))} = \int_0^1 dy g(y)$$

neboť $y(x) = \int_0^x dx' w(x')$, $\frac{dy}{dx} = w(x)$

$$y(0) = 0, y(1) = 1$$

problem je ovšem s inverzí funkce $y(x)$, i když najde se vhodné $w(x)$

Příklad: $f(x) = \frac{1}{1+x^2}$



Ize např. vztah $w(x) = \frac{4}{3} - \frac{2}{3}x$ (koeficienty zvoleny tak, aby $\int_0^1 w(x) dx = 1$)

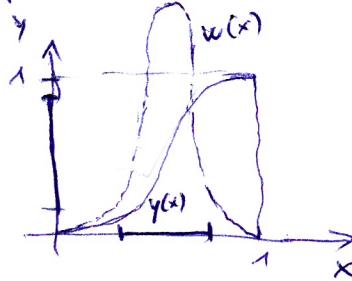
máme $y(x) = \frac{1}{3}x(4-x)$

neboli $x(y) = 2 - \sqrt{4-3y}$

tato jednoduchá transformace dá o následkem lepší výsledek!

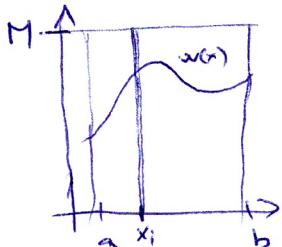
- protože najít $x(y)$ je u vícerozměrných integrálů prakticky nemožné ($w(x)$ je Jacobian $\left| \frac{\partial \vec{x}}{\partial \vec{y}} \right|$)
 je výhodnejší se na vztah $dy = w(x) dx$ dívat jako na základu rovnoučinného rozdělení $g(y)=1$ na rozdělení $w(x)$

- je-li y rovnouèerná na $\langle 0,1 \rangle$, pak x bude rozdelené podle $w(x)$
a tedy preferuje se x tam, kde má $w(x)$ velké hodnoty neboť tam $y(x)$ rychle roste
- místo změny proměnných tedy budeme stíle používat x , ale místo rovnouèerného rozdelení pro x potřebujeme umět generovat náhodnou veličinu s daným rozdelením pravděpodobnosti $w(x)$



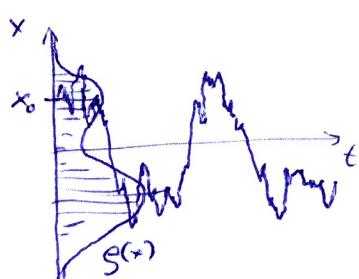
Generování náhodných veličin

- v počítaci vše většinou k dispozici pseudonáhodný generátor čísel na intervalu $\langle 0,1 \rangle$ s rovnouèerným rozdelením
- existují speciální postupy pro určitá rozdelení např. pro normální rozdelení ve 2D přechodem k polárním souřadnicím (viz notebook v Mathematica)
- obecně lze použít von Neumannova metodu či Metropolisov-Hastingsov algoritmus
- von Neumannova metoda
 - generuje dvojici náhodných čísel s rovnouèerným a podle $w(x)$ rozdelením, zda vygenerované x_i použije - nejde $\xi_i, M_i \in \langle 0,1 \rangle$ náhodná ξ_i namapuje na $\langle a,b \rangle$, $x_i = a + \xi_i(b-a)$ M_i namapuje na $\langle 0,M \rangle$, $y_i = M_i y_i$ (M zvoleno tak, aby $w(x) < M$ pro $x \in \langle a,b \rangle$)
 - pokud $y_i \leq w(x_i)$, pak x_i použije
 - pokud $y_i > w(x_i)$, pak x_i zahodi
 - pravděpodobnost přijetí hodnoty v $\langle x, x+dx \rangle$ je j-érna $w(x)dx$ neboť ξ_i je rovnouèerně rozložené na $\langle a,b \rangle$



- Metropolisův - Hastingsův algoritmus

- základní myšlenka: jít „náhodnou“ procházkou v prostoru



vzávislosti na hustotě pravděpodobnosti, tj. jsme-li v bodě x_k , učiníme náhodný krok o délce $\delta x_k = \xi_k \delta$, kde $\xi_k \in \{-1, 1\}$ a δ je fixní parametr procházkы (musí být vhodně nastaven, aby bylo případné možno překročit oblasti s malou hustotou),

a tento krok učiníme, pokud $g(x_k + \delta x_k) \geq g(x_k)$, nebo je-li $g(x_k + \delta x_k) < g(x_k)$, pak vygenerujeme náhodné $y \in (0, 1)$ a krok učiníme, pokud $g(x_k + \delta x_k) / g(x_k) > y$.

- pokud krok neděláme, pak zůstaneme a $x_{k+1} = x_k$.
- v d -rozměrném prostoru bude algoritmus vypadat takto:

procedura
vracející'
další'
bod, jsme-li
v $x(:)$

```
r(:) = random(:)    ← generátor na intervalu  $(0, 1)$ 
xtrial(:) = x(:) + δ(2r(:) - 1)    ← posun na  $(-1, 1)$ 
ratio = g(xtrial) / g(x)
if (ratio ≥ 1 or ratio > random) then
    x = xtrial
return x(:)
```

- v limitě dlouhé náhodné procházky tento algoritmus opravdu generuje body rozložené podle $g(x)$.

- označme-li $N_n(x)$ hustotu náhodných nezávislých chodců startujících v různych bodech po n krocích, pak počet chodců jdoucích z x do y v následujícím kroku je

$$\Delta N(x) = N_n(x) P(x \rightarrow y) - N_n(y) P(y \rightarrow x) =$$

$$= N_n(y) P(x \rightarrow y) \left[\frac{N_n(x)}{N_n(y)} - \frac{P(y \rightarrow x)}{P(x \rightarrow y)} \right]$$

kde $P(x \rightarrow y)$ je pravděpodobnost, že chodec přejde z x do y .

- rovnováha nastává pro

$$\frac{N_n(x)}{N_n(y)} = \frac{P(y \rightarrow x)}{P(x \rightarrow y)}$$

a lze ukázat, že pro velká n jde $N_n(x) \rightarrow N_e(x)$,

což je rovnovážné rozdělení chodců v x

- pro M-H algoritmus je $N_e(x) \sim g(x)$, neboť
 $P(x \rightarrow y)$ lze vyjádřit jako

$$P(x \rightarrow y) = T(x \rightarrow y) A(x \rightarrow y)$$

kde $T(x \rightarrow y)$ je pravděpodobnost náhodného výkročení
z x do y a $A(x \rightarrow y)$ je pravděpodobnost přijetí
tohoto kroku

- pokud tedy bude

$$\frac{g(x)}{g(y)} = \frac{P(y \rightarrow x)}{P(x \rightarrow y)} = \frac{T(y \rightarrow x) A(y \rightarrow x)}{T(x \rightarrow y) A(x \rightarrow y)}$$

pak $N_e(x)$ bude $\sim g(x)$

- pro vhodně zvolené δ , kdy lze y dosíhnout v jediném
kroku, tj. $y \in (x-\delta, x+\delta)$, pak $T(x \rightarrow y) = T(y \rightarrow x)$

a dále pro $g(x) \geq g(y)$ bude

$$A(y \rightarrow x) = 1 \quad \text{a} \quad A(x \rightarrow y) = \frac{g(y)}{g(x)}$$

a pro $g(x) \leq g(y)$ bude

$$A(x \rightarrow y) = 1 \quad \text{a} \quad A(y \rightarrow x) = \frac{g(x)}{g(y)}$$

tedy celkově v každém případě

$$\frac{A(y \rightarrow x)}{A(x \rightarrow y)} = \frac{g(x)}{g(y)}$$

a tedy vzhledem $\frac{N_e(x)}{N_e(y)} = \frac{g(x)}{g(y)}$

Generitory (pseudo) náhodých čísel (random numbers)

- deterministické programy generující posloupnost čísel, které splňují jisté kritéria (testy) náhodnosti
- při stejných počítačích generují tedy náhodná čísla
 - poč. pod. = random seed (náhodné "seinko")
- vzhledem ke konečné aritmetice počítaců ~~je~~ a deterministickou charakterem nejde o náhodná čísla, ale o pseudonáhodná
- většinou generována jako posloupnost čísel $\{x_0, x_1, x_2, \dots\}$ → vlna periodicitu, kvůli konstrukci aritmetice \Rightarrow perioda generatoru = nejmenší k takové, že $x_i = x_{i+k}$ pro každou

Základní generitory

- lineární kongruentní generátor (LCG) - velmi rychlý algoritmus
 - maximální (optimální) perioda = větší, než menší

speciální případ - multiplicativní LCG, když $c=0$ (navrhl D.H. Lehmer)

byl velmi často používán v historii, ale má často zásadní nedostatky pro aplikace ve více dimenzích, neboť matici $(x_i, x_{i+1}, \dots, x_{i+n-1})$ se nachází na ~~je~~ relativně malém počtu $(n-1)$ -rozm. nadplochách

Př. Nejznámější generátor RANDU

$$a = 65539, c = 0, m = 2^{31} \text{ splňuje rekur. vztah } x_{i+2} = (2^{16} + 3)x_{i+1} \bmod 2^{31}$$

$$x_{i+2} = (2^{32} + 6 \cdot 2^{16} + 9)x_i \bmod 2^{31} = (6x_{i+1} - 9x_i) \bmod 2^{31}$$

a lze ukázat, že k trojice bodů se nachází na 15 paralelních rovinách

Marsaglia: Random Numbers Fall Mainly in the Planes, PNAS 61 (1968) 25

Teoretyčtí pro MLGG: Pokud c_1, \dots, c_n jsou přirozená čísla taková, že

$$c_1 + c_2 a + c_3 a^2 + \dots + c_n a^{n-1} \equiv 0 \pmod{m}$$

pak k třem body $P_1 = (x_1, x_2, \dots, x_n), P_2 = (x_2, \dots, x_{n+1}), \dots$ leží v jedné z paralelních nadploch

$$c_1 z_1 + c_2 z_2 + \dots + c_n z_n = 0, \pm 1, \pm 2, \dots$$

Navíc je nejvýše $|c_1| + |c_2| + \dots + |c_n|$ těchto rovin protká krychli

$0 < z_1 < 1, \dots, 0 < z_n < 1$ a vždy lze vybrat taková c_1, \dots, c_n , že

tři body budou vzdálené než $(n!m)^{1/n}$ nadplochách.