

Tvorování dat pro experimenty s algoritmy

- a) náhodující generování (např. číslo průsopnosti, graf)
- b) použitím existujících dátových souborů (jmenem osoby, telefon, adresy)

Náhodné generování

Pro algoritmus přísluší jednotlivě k binárním číslům jiné početnosti
 pak obecný postup ještě může ještě vypadat, že $P(a_i = 1) = \frac{1}{2}$,
 $P(a_i = 0) = \frac{1}{2}$. Jak je možné vytvořit?

1) budeme hledat menší - když v i-ém bode padne číslo,

bude $a_i = 1$, když padne nečíslo, bude $a_i = 0$

2) použijeme počítac, na kterém máme generátor náhod, jde
 o řešení komplementního rozdělení $\mathcal{U}(0, 1)$

když bude i-é generované číslo menší než 1/2, bude $a_i = 1$,
 jinak $a_i = 0$

$$\text{dále } P(X_i < \frac{1}{2}) = F\left(\frac{1}{2}\right) = \frac{1}{2}$$

$$P(X_i \geq \frac{1}{2}) = 1 - P(X_i < \frac{1}{2}) = \frac{1}{2}$$

jak ale dostanu náhodné číslo z $\mathcal{U}(0, 1)$?

2.1) Kvalitní se generuje číslo 11 hodin některé
 na listci majíš postupně několik čísel od 0 do 14-1,
 složíš do bloků, namíchat a se rovinou
 očima jeden listec vyhlášne - nech je na něm číslo X

$$P(X=k) = \frac{1}{14} \text{ pro m. } k=0, 1, \dots, 14-1$$

položíš Y = $\frac{X}{14-1}$ a to je moje náhodné číslo z $\mathcal{U}(0, 1)$

2.2) použije po výpočtu x následující řešení:
 řešenek by měl mít své vlastnosti, jehož hodnota je čísla taká
 k blokovku

čísla posízena límkem výpočtem se nazývají pseudorandom

jich výhody: dají se snadno a rychle generovat
 dají se s pravděpodobností vygenerovat několik

výhody: nejsou výpočtem národního počítače
 nemusí být periodická

mohou se periodicky opakovat

jich pravděpodobnosti nemusí písničky soubasit s pořadovým
 nepravidelností f - liniární homogenní generátory jsou
 rozdílením

$$X_m = (a_0 + a_1 X_{m-1} + \dots + a_k X_{m-k}) \bmod M$$

hde konstanty M, a_0, \dots, a_k je třeba vhodně volit -
 má nich některé kvalita generátoru

pořadatel má generátor: maximální možná perioda M

minimální serialní korelace

shoda s pořadovým rozdílením

kvalita se dá obecně dokazovat teoreticky

ovíruje se statistickými testy

našel jsem některé konkrétní používání generátoru.

(přehled je z roku 1950)

Tabulka 1. Přehled vybraných typů generátorů z překladačů a statistických programových systémů.

Číslo generátoru

Zdroj

Tvar

Poznámka

1 Překladač jazyka APL.
 $V_{i+1} = (16807V_i) \bmod (2^{31}-1)$

Tento generátor implicitně užívá i většina programových systémů napsaných v APL, např. Statgraphics.

2 Počítacího typu ATARI ST
 (v OS ROM).

$V_{i+1} = (3141592621V_i + 1) \bmod (2^{32})$

3 Počítacího typu IBM
 360/370, resp. operační
 systém VMS pro DEC VAX.

$V_{i+1} = ((2^{16} + 3)V_i) \bmod (2^{31})$

Tzv. generátor RANDU. Jedná se o jeden z nejméně kvalitních generátorů jež byl, a leckde bohužel doposud je, široce používán

4

a) $V_{i+1} = (950706376V_i) \bmod (2^{31}-1)$

Nejlepší generátor typu
 $V_{i+1} = (aV_i) \bmod (2^{31}-1)$ (dle hodnocení
 Fishmana a Moorea (1982)).

Knihovna podprogramů
 IMSL, část Stat/Lib,
 v.1.4.

b) $V_{i+1} = (397204094V_i) \bmod (2^{31}-1)$

Viz LLRANDOMII.

c) $V_{i+1} = (16807V_i) \bmod (2^{31}-1)$

Viz LLRANDOMI.

LLRANDOMI
 (viz Lewis a další
 (1988)).

$V_{i+1} = (16807V_i) \bmod (2^{31}-1)$

Jeden z prvních úspěšných generátorů typu $V_{i+1} = (aV_i) \bmod (2^{31}-1)$, navržený již v roce 1969 Lewisem a kol.. Postupně byl převzat řadou dalších firem. Po čase mu byly prokázány neoptimální vlastnosti ve vyšších dimenzích.

6 LLRANDOMII
 viz Lewis a další (1988).

$V_{i+1} = (397204094V_i) \bmod (2^{31}-1)$

Modifikace generátoru LLRANDOMI
 navržená v roce 1974 Lewisem a kol.. pro odstraňení "špatného" chování generátoru LLRANDOMI ve vyšších dimenzích.

(719)

Číslo generátoru	Zdroj	Tvar	Poznámka
------------------	-------	------	----------

- 32 -
- 7 Knihovna podprogramů využitých v JET Propulsion Laboratories v Los Alamos.
- $$V_{i+1} = (5^{19} V_i) \bmod (2^{48})$$
- Generátor navržený a široce používaný v laboratořích v Los Alamos na sálových počítacích.
-
- 8 Knihovna podprogramů NAG, v.11.
- $$V_{i+1} = (13^{13} V_i) \bmod (2^{59})$$
- Značně obškurní volba konstant, jež nebyla nikdy žádnej vysvětlena. NAG ani nikdy neposkytl podrobnejší informaci o vlastnostech svého generátoru.
-
- 9 SASPC v.6.03.
- a) $V_{i+1} = (16807 V_i) \bmod (2^{31}-1)$
Viz LLRANDOMI.
- b) $V_{i+1} = (3977204094 V_i) \bmod (2^{31}-1)$
Viz LLRANDOMII.
-
- 10 Prekladač jazyka SIMULA.
- $$V_{i+1} = (5^{2p+1} V_i) \bmod (2^n)$$
- Konstanty p a n jsou voleny dle typu počítace, na němž je jazyk implementován.
-
- 11 SPSS/PC v.2, resp. SPSSX v.8.
- $$V_{i+1} = (16807 V_i) \bmod (2^{31}-1)$$
- Viz LLRANDOMI.
-
- 12 Operační systém UNIX.
- a) $V_{i+1} = (1103515245 V_i + 12345) \bmod (2^{31})$
Tzv. generátor "RAND".
- b) $V_{i+1} = (25214903917 V_i + 11) \bmod (2^{48})$
Tzv. generátor "DRAND".

Metody pro generování jiných rozdělení

a) obecně platné'

b) speciální

ada) Metoda inverzní transformace

Veta: Nechť $F(x)$ je distribuční funkce, $F^{-1}(y) = \inf\{x : F(x) \geq y\}$
je odpovídající invazní funkce a nechť Y má rozdělení $R(0,1)$.

Potom náhodná veličina $X = F^{-1}(Y)$ má rozdělení s distribucí $F(x)$.

Důkaz po případu, když F je spojita a rostoucí:

F^{-1} je invazní funkce k F

$$P(X < x) = P(F^{-1}(Y) < x) = P(Y < F(x)) = F(x)$$

o ostatních případech se musí opakovat - zdejší se z toho,
že F je vždy spojita a měkkosrající

metoda je obecně platná'

po měkkosrajení se hodi více, po jiné méně
hodi se např. po exponenciální rozdělení:

$$F(x) = 1 - e^{-\lambda x} \text{ pro } x \geq 0$$

= 0 jinak

$$\text{invazní funkce: } F^{-1}(y) = -\frac{1}{\lambda} \ln(1-y), \quad 0 < y < 1$$

$$F^{-1}(F(x)) = -\frac{1}{\lambda} \ln(1 - (1 - e^{-\lambda x})) = -\frac{1}{\lambda} \ln e^{-\lambda x} = x$$

Tabulkou: vygenerujte $Y \sim R(0,1)$

$$\text{potom } X = -\frac{1}{\lambda} \ln(1 - Y)$$

sloučíme následkem $X = -\frac{1}{\lambda} \ln Y$, protože $1 - Y$ má rozdělení $R(0,1)$

$$P(1-Y < x) = P(Y > 1-x) = 1 - P(Y \leq 1-x) = 1 - (1-x) = x$$

dále je vhodná po rozdělení, když F je unikátní (one-to-one) a dá se snadno invertovat

není např. vhodná po normální rozdělení, když známá je hustota a distribuce se nám počítá jako $F(x) = \int_{-\infty}^x f(t)dt$

pokud bývá možnou být u diskrétních rozdělení

Metoda ramilaci (von Neumann) - vhodná po příjmu, když znám hustotu a ta je matic obaricena

$$\text{Nechť } c = \max \{f(x), a \leq x \leq b\}$$

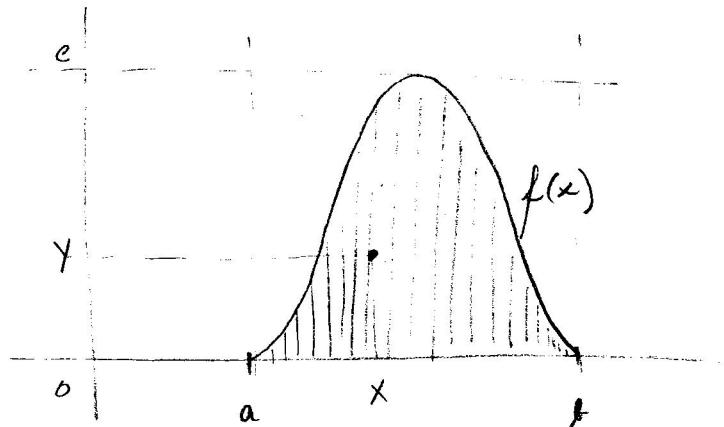
generujeme: X rovnoměrně na (a, b)

Y rovnoměrně na $(0, c)$

je-li $Y \leq f(x)$, uvedeme X jako náhodné číslo
s rozdílením $f(x)$

o opačním případě drojce X, Y ramilaceme a
generujeme novou

$$\text{ab: } P(Y \leq f(x)) = f(x)$$



generujeme rovnoměrně rozd.

na $(a, b) \times (0, c)$

vybíráme body, které jsou
pod $f(x)$

Pro diskritní máhodné relaci

Chci generovat mál. číslo x rozdělení p_0, p_1, \dots , kde $p_i = P(X=i)$

měloda reverzní transformace:

polozím $X=0$

$$S = p_0$$

vygeneruj $U \in R(0,1)$

while $U > S$ do $X = X+1$, $S = S + p_X$ od

X je mál. čísla k pořadování rozdělení

dle: $P(X=k) = P(\text{cyklus probíh k-kále}) =$

$$= P(U > p_0 + p_1 + \dots + p_{k-1} \wedge U < p_0 + p_1 + \dots + p_k) =$$

$$= \sum_{i=0}^k p_i - \sum_{i=0}^{k-1} p_i = p_k$$

nebo pomocí "projile" máhodné relaci s dist. fci H , pro

blízou $H(i+1) - H(i) = p_i$

vygenerujeme $U \in R(0,1)$

polozime $X = [H'(U)]$ (ala 'čáš')

např. po diskritní romomeréni rozdělení na $\{1, \dots, m\}$

$$H(x) = \frac{x-1}{m} \quad \text{pro } 1 \leq x \leq m+1$$

$$\text{bude } H^{-1}(y) = 1 + my$$

$$H^{-1}(H(x)) = 1 + m \cdot \frac{x-1}{m} = 1 + x - 1 = x$$

$$\text{doslaneme } X = [1 + mU]$$

ad b) Normalní rozdělení - maláda imenována transformace se rozděluje, protože její kvantile je pro jeho poměrně zdrobný rozdělení vlastně mít. velikou $\approx N(0,1)$.

Použije se centrální limitní věta, na součet nzávislých síticích rozdělených mít. velikou $\approx N(0,1)$.

X_1, \dots, X_n mítodná čísla $\approx N(0,1)$

$$E \sum_{i=1}^n X_i = \frac{n}{2}, \quad \text{var } \sum_{i=1}^n X_i = \frac{n}{12}$$

$Y_m = \sqrt{\frac{n}{m}} \left(\sum_{i=1}^m X_i - \frac{n}{2} \right) \sim N(0,1)$ asymptoticky pro $m \rightarrow \infty$
a pak se bere $m=12$, tedy $Y_{12} = \sum_{i=1}^{12} X_i - 6$.

Box-Müller: X_1, X_2 mítodn. mít. vel. $\approx N(0,1)$

$$Z_1 = \sqrt{-2 \ln X_1} \sin(2\pi X_2)$$

$$Z_2 = \sqrt{-2 \ln X_1} \cos(2\pi X_2)$$

Z_1, Z_2 jsou mítodn. mít. vel. $\approx N(0,1)$

dálej uvedená projekce rozdělení se daje generální rozložit. tabulkou.

χ^2_m -rozd: $X = \sum_{i=1}^m Z_i^2$, kde Z_i jsou nezávislé mít. $N(0,1)$

T_m -rozd: $X = \frac{Y_1}{\sqrt{\frac{Y_2}{m}}}$, kde Y_1, Y_2 jsou mítodn.,
 $Y_1 \sim N(0,1)$, $Y_2 \sim \chi^2_m$

$F_{m,m}$ -rozd: $X = \frac{\frac{Y_1}{m}}{\frac{Y_2}{m}}$, kde Y_1, Y_2 jsou mítodn.,
 $Y_1 \sim \chi^2_m$, $Y_2 \sim \chi^2_m$

Generování diskritních rozdělení

1) Binomické rozdělení $P(X=k) = \binom{m}{k} p^k (1-p)^{m-k}$

popisuje pravd. výskytu n pokusů s konstantními pravd. výskytu jednotlivých pokusů, a když je pak výskyt pr

postup: počítáme $X=0$

metoda generování: generujeme U náhodně mezi $(0,1)$
je-li $U < p$, počítáme $X \leftarrow X+1$

X je malé číslo a binom. rozd.

$$\text{dL: } P(X=k) = \binom{m}{k} p^k (1-p)^{m-k}$$

po náhodném rozd. $P(U < p) = p$

$$P(U > p) = 1-p$$

2) geometrické rozdělení $P(X=k) = (1-p)^k p$

o počítáme konstantní pravd. výskyt jednotlivých pokusů, až do když generujeme U náhodně mezi $(0,1)$

postup: počítáme $X=-1$

generujeme, dokud $U > p$: počítáme $X = X+1$

generujeme U náhodně mezi $(0,1)$

3) Poissonovo rozdělení - jeho vlastní k binomickému

Tyto metody reálně simulují náhodný počet

jou počítáme, ale často výčet

používají se mnohem jiné a lepši, a bývají se dát až někdy teoreticky těžko dokázat, že generují to, co mají.

Při generování čísel s pořadovým rozdělením se obvykle dletočky řeší.

χ^2 -test dobré šady (srovnávání)

obor možných hodnot málo nákladně řešícího rozdělení na k disjunktních intervalů A_1, A_2, \dots, A_k - např. $(-\infty, x_1), (x_1, x_2), \dots, (x_{k-1}, \infty)$

označím $p_j = P(X \in A_j)$ - byla pravděpodobností kromě mimo
je umím specifikat (např. pro rovnoměrné rozdělení) hodnoty

$$p_j = \frac{1}{n} \text{ pro } j=1, \dots, k, \text{ pro } N(0,1) \text{ je}$$

$$P(X \in (x_j, x_{j+1})) = \Phi(x_{j+1}) - \Phi(x_j)$$

označím m_j počet vygenerovaných hodnot, které padly do A_j
(empirické čísla)

$\frac{m_j}{n}$, kde $n = \sum_{j=1}^k m_j$, je empirická relativní čísla, jejím
teoretickým protějškem je p_j

Kol. je založen na odchylek $m_j - np_j$

konkrétně na statistice $\chi^2 = \sum_{j=1}^k \frac{(m_j - np_j)^2}{np_j}$

platí, že χ^2 má asymptoticky rozdělení χ^2_{k-1}

přednosti - liž χ^2 přesloužil kritickou hodnotu, hypotéza o
šadě rozdělení se vzdává

Poznámka: test se da' použít např. k ověření normality podle použitím t-testu, F-testu a podobných
 pokud je, že bývá předem známo některé parametry rozdělení,
 nejméně specifické parametry, které potřebují - meloda se pak
 komplikuje tím, že se při něj zároveň musí počítat nějaké odbady

další testy:

používání teoretičtějších momentů různých rádu příslušného rozdělení
 s jejich výslovnějšími požadavky specifickými a generovanými hodnotami

$$E X^k \text{ používam } \rightarrow \frac{1}{n} \sum_{i=1}^n X_i^k$$

Kolmogorov-Gomringerův test - je rozdíl mezi teoretické
 a empirické distribuční fce

$$F_m(x) = \text{empirická d.f.} = \frac{\text{počet hodnot menších než } x}{\text{počet všech hodnot}} = P(X \leq x) = F(x)$$

konkrétně: používám $D_n = \sup_x |F_m(x) - F(x)|$

hypoteticku ramilnu, bývá $D_n \geq D_n(\alpha)$, kde

$$D_n(\alpha) = \sqrt{\frac{1}{2n} \ln \frac{2}{\alpha}}$$

(kritické hodnoty vyzájí tabulkou)

Generacií množiných posloupností

grafovaným použitím generátoru - musíme najít možnosť grafickej možnosti poriadania na posloupnosť.

- a) čísla sú normované grafom
- b) posloupnosť, kde sú usporiadane;
- c) sú usporiadane permutace
- d) ?

ad a) generuje posloupnosť čísla a kontroly, počle je už veľkou mierou - poludanos, keď teoreticky: nahradím aly rýber a racinek pomocou policie: miesto čísla, akere sú mŕtve, vyzdvihujem

ad b) vyzdvihujem malodaný rýber, kdežto jah polohy
nebo použijem mellerou metodu počítačom usporiadanie rýber

ad c) jedna z možností metod:

vyhľadávanie usporiadania pole $\alpha(1)=1, \alpha(2)=2, \dots, \alpha(n)=n$

po $j = n, n-1, \dots, 2$

vyzdvihujem čísla i korektnosť na $\{1, \dots, j\}$

nahradím $\alpha(j) \leftarrow \alpha(i)$

Generování náhodných posloupností - jiná metoda

Předpokládajme, že $b = (b_0, \dots, b_{N-1})$ je nějaká posloupnost
čísel $0, \dots, N-1$.

Okamžitě a_k počít si čísla b_i , která v posloupnosti málojiž nejsou
a jsou menší než k .

(př.: $b = (0, 5, 7, 3, 1, 4, 2, 6)$, tj. $N=8$)

$$a_0 = 0, a_1 = 0, a_2 = 2, a_3 = 1, a_4 = 5, a_5 = 0, a_6 = 5,$$

Když $a_k \leq k$ po výčtu b .

† Kždej posloupnosti b přiřadíme čísla

$$\tau(b) = a_0 1! + a_1 2! + \dots + a_{N-1} (N-1)!$$

Tím je definována významně polynomálně rozhazující
množina všech posloupností na množinu čísel $0, 1, \dots, N^L-1$.

Postup:

- 1) Generujeme náhodné čísla z distribučního rozmezí
rozdělení $R(0, 1, \dots, N^L-1)$
- 2) učíme čísla a_0, \dots, a_{N-1}
- 3) a následujeme odpovídající posloupnost

Příklad: vygenerovaly jsou náhodná čísla 19885

a rozdělení $R(0, 1, \dots, 8^L-1)$

postupne doslavame:

$$19885 = 3 \cdot 7! + 4765 \quad a_7 = 3$$

$$4765 = 6 \cdot 6! + 495 \quad a_6 = 6$$

$$495 = 3 \cdot 5! + 85 \quad a_5 = 3$$

$$85 = 3 \cdot 4! + 13 \quad a_4 = 3$$

$$13 = 2 \cdot 3! + 1 \quad a_3 = 2$$

$$1 = 0 \cdot 2! + 1 \quad a_2 = 0$$

$$1 = 1 \cdot 1! + 0 \quad a_1 = 1$$

pri konstrukcií posuvnej doslavame od alla 0

doslavame 0

1 0

1 0 2

1 3 0 2

1 4 3 0 2

1 4 5 3 0 2

6 1 4 5 3 0 2

6 1 4 5 4 3 0 2

generování grafů

graf $G = (V, E)$ $V \dots$ množina vrcholů, $|V| = n$

$E \dots$ množina hrani, tj. dvojic (i, j) , $i \in V, j \in V$

graf je s počtem hrani:

vygeneruje pětset šestnáct dvojic rombového na čtverci $V \times V$

je (i, i) hrana? - pokud ne, tyto dvojice vyhodím a vygeneruju mišle nichžíne'

da' se to použít pro orientované i neorientované grafy -

pro neorientované grafy výhodné dvojici (i, j) , když
už mám ne výběr (j, i)

polohu naslouhnu po dalších poradancích - jak nazv. magický
asymetrický graf?

náhodný graf $G(n, p)$:

pro každou dvojici (i, j) generuje náhodné čísla s alternativním

dělením s poli $p = P(\text{node mana } \& \text{ i do } j)$

$1-p = P(\text{node } \dots \text{ --})$

tj. generuje $0 \& R(0, 1)$

když $U < p$, pak node mana

$\&$ opačně, případě node

pro $p = \frac{1}{2}$ je to stejné, jako když náhodné (rombové)

výběr a množství nich menších grafů s daným počtem vrcholů

Nabodné grafy generování limitním přirozenem ale často upřadají tak, že celý graf je jedna silně soudě komponenta, původní graf obsahuje 1 nebo několik silně soudě komponent a mnoho maliných (jednobodových) komponent.

Pokud potřebujeme graf, který má vše malináčky SSK, můžeme použít model $G(m, p, \ell)$. Parametr ℓ (který lokalita) určuje oblast vrcholu, do kterého může mít hrana. Konkrétně (i, j) může být hrana, když $(i - \ell + m) \bmod m \leq j \leq (i + \ell) \bmod m$.
 ℓ by mělo být malej (5, 10), pro nějaký ℓ je $G(m, p, \ell)$ podobný grafu $G(m, p)$.

Acyklické orientované grafy můžeme generovat pomocí jeho topologického uspořádání, tj. směr se na dvojici (i, j) , kde $i < j$.