

PSK - zajęcia zaplanowane na 14 & 21.05

Analiza statystyczna wyników symulacji

Nsbp. Badania symulacyjne przeprowadza się najczęściej po to, aby ustalić wartości  $\theta$  pewnej wielkości zmiennej z określonym modelem stochastycznym, który prowadzi do wykonania tej symulacji.

Wtedy wynikiem takiej symulacji jest zmienna losowa  $X$ , której wartość oczekiwaną  $EX = \theta$ .

Takie potrzeby symulacji powstają w sposób naturalny np. do pewnej ilości,  $k$ -probiegów.

Być może ciąg niezależnych zmiennych losowych

$X_1, X_2, \dots, X_k$  o tym samym rozkładzie co  $X$ .

Wtedy (na mocy moichy prawa wielkich liczb Bernoulliego) średnia

$$\bar{X}(k) = \frac{1}{k} (X_1(k) + \dots + X_k(k)) \approx \theta \quad \text{jestymatorem tej wartości } \theta.$$

Pojawia się problem: kiedy zatrzymać proces symulacji, czyli jak ustalić  $k$ , aby pomiar był precyzyjny

lepto adekwatne do badanej sygnali metode symulacji.

Idac dalej, czy dla danych licy  $\alpha \in (0, 1)$ ,

da mi znalesci temu predmiot losowy obrazy!

przez dane statystyki

$$Z_1 = f_1(X_{n-1}, X_n) < Z_2 = f_2(X_{n-1}, X_n),$$

$$\textcircled{*} \text{ie } P(\text{Kwadrat: } Z_1(\omega) < \theta < Z_2(\omega)) \geq 1 - \alpha$$

Srednia i wariancja z podly jako estymatory

Nech  $X$  oznacza zadeny rozklad. Dk  $n \geq 2$ ,

niech  $X_1, X_2, \dots, X_n$  - zmienna losowa, takie sa

(i) niezalezne

(ii) każda ma rozklad jak  $X$

Jeli  $EX = \theta$ ,  $\text{var}(X) = \sigma^2$ , b

$E(X_j) = \theta$ ,  $\text{var}(X_j) = \sigma^2$ ;

a wielken  $(X_1, X_2, \dots, X_n)$  nazwiemy PRÓBA,

POW SREDNIA, z PRÓBY nominy

$$\bar{X} = \frac{1}{n} \sum_{j=1}^n X_j$$

Pokazemy, że wyliczeni  $\bar{X}$  do oszacowania (nieznanej) wartości  $\theta$ .

Zauważmy, że

$$E\bar{X} = E\left(\frac{1}{n} \sum_{j=1}^n X_j\right) = \frac{1}{n} \sum_{j=1}^n EX_j = \frac{n\theta}{n} = \theta$$

Oblinaj ten średnie kwadratowe odchylenie pomijając wartości  $\bar{X}$  a wartości  $\theta$ , czyli

$$E[(\bar{X} - \theta)^2] = \text{Var}(\bar{X}), \text{ bo } E\bar{X} = \theta$$

Ali z niezależności  $X_1, \dots, X_n$ :

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum_{j=1}^n X_j\right) =$$

$$= \frac{1}{n^2} \sum_{j=1}^n \text{Var}(X_j) = \frac{1}{n^2} \cdot n\sigma^2 = \frac{\sigma^2}{n},$$

co oznacza, że jeśli  $n$  jest  $m$  p. dostatecznie duże,

to  $E[(\bar{X} - \theta)^2]$  p. jest małe (jak chcemy),

czyli  $\bar{X}$  w tym sensie dobrze estymuje

$\theta$ .

-h-

Co więcej, korzystając z nierówności Czebyszewa  
możemy oszacować miarę rozpręgnięcia wariancji  $\bar{X}$   
względem  $\theta$ . Dokładniej, dla  $c > 0$

$$P(\text{złoty}: |\bar{X}(n) - \theta| \geq c \sqrt{\text{var} \bar{X}}) \leq \frac{1}{c^2}$$

||

$$P(\text{złoty}: |\bar{X}(n) - \theta| \geq \frac{c\sigma}{\sqrt{n}}) \leq \frac{1}{c^2},$$

co równoważnie możemy zapisać

$$P(\text{złoty}: |\bar{X}(n) - \theta| < \frac{c\sigma}{\sqrt{n}}) > 1 - \frac{1}{c^2}$$

||

$$(*) P(\text{złoty}: \theta \in (\bar{X}(n) - \frac{c\sigma}{\sqrt{n}}, \bar{X}(n) + \frac{c\sigma}{\sqrt{n}})) > 1 - \frac{1}{c^2}$$

co pozwala określić nam przedział losowy,  
o którym wspominalyśmy w (\*) (str. 2).

Problem z wyliczaniem miary (xx) polega na tym, iż na ogół wartość  $\sigma^2$  nie jest znana.

Tymczasem jest zatem oszacować (estymować).

Bezpośrednim odpowiednikiem  $\sigma^2 = E(X - \theta)^2$  jest natomiast błąd estymacji oznaczony jako

$$S^2 = \frac{1}{n} \sum_{j=1}^n (X_j - \bar{X})^2.$$

Ma on jednak pewną wady, mianowicie

$$ES^2 \neq \sigma^2.$$

Jest dokonany jego korekty, c' wzmocny

$$\frac{n}{n-1} S^2 = \frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X})^2,$$

o otrzy estymator oznaczony  $\hat{S}^2$ , taki iż

$$\underline{E \hat{S}^2 = \sigma^2} \rightarrow \text{WARIANCJA} \\ \geq \text{PRÓBY}$$

$$\text{Wtedy } \hat{S} = \sqrt{\hat{S}^2} \rightarrow \text{DYSPERSIJA} \\ \geq \text{PRÓBY}$$

Mozemy domniemywać, że  $(x, x)$ . Pomyślmy, że  $N(x, x)$ :

(i) chcemy, aby

$$1 - \frac{1}{c^2} \geq 0,95 \equiv (c > 0)$$

$$\frac{5}{100} \geq \frac{1}{c^2} \equiv c^2 \geq 20 \equiv$$

czyli  $c \geq \sqrt{20}$

(ii) chcemy mieć przedział

$$\left( \bar{X} - \frac{c\sigma}{\sqrt{n}}, \bar{X} + \frac{c\sigma}{\sqrt{n}} \right), \text{ czyli}$$

$\frac{2c\sigma}{\sqrt{n}}$  była mniejsza od zadanej liczby

$$2d \quad (d > 0), \text{ czyli}$$

$$\frac{2c\sigma}{\sqrt{n}} \leq 2d \equiv \frac{c\sigma}{\sqrt{n}} \leq d$$

$$\equiv \sqrt{n} \geq \frac{c\sigma}{d} \equiv$$

$$n \geq \frac{c^2 \sigma^2}{d^2} \geq 20 \frac{\sigma^2}{d^2} \approx 20 \frac{5^2}{d^2}$$

Zatem dla  $n \gg 20 \frac{s^2}{d^2}$ , warunki (i) i (ii) spełniony, przy wymaganiach zadanych w (i'), (ii').

Podsumowanie

Uwaga. Przewidywalność wystąpienia w (i) i (ii)  $P(\theta \in (\bar{X}(n) - \frac{c\sigma}{\sqrt{n}}, \bar{X}(n) + \frac{c\sigma}{\sqrt{n}}))$

można oszacować lepiej, a jeżeli wyraża to z nierówności Czebyszewa.

W tym celu zauważamy, że

$$\frac{\bar{X} - E\bar{X}}{\sqrt{\text{var}(\bar{X})}} = \frac{\bar{X} - \theta}{\frac{\sigma}{\sqrt{n}}} \quad \text{oraz dla}$$

odpowiednio dostajemy  $n (n \gg 200) = \text{C.T.G.}$

$$\frac{\bar{X} - \theta}{\frac{\sigma}{\sqrt{n}}} \in \mathcal{N}(0, 1)$$

Pomocniczo

$$\theta \in \left( \bar{X}(n) - \frac{c\sigma}{\sqrt{n}}, \bar{X}(n) + \frac{c\sigma}{\sqrt{n}} \right) \equiv$$

8 -

$$|\bar{X}(n) - \theta| < \frac{c\sigma}{\sqrt{n}} \equiv$$

$$\left| \frac{\bar{X}(n) - \theta}{\frac{\sigma}{\sqrt{n}}} \right| < c, \text{ nje dla takich } n,$$

$$P(\text{zawor: } \theta \in (\bar{X}(n) - \frac{c\sigma}{\sqrt{n}}, \bar{X}(n) + \frac{c\sigma}{\sqrt{n}})) =$$

$$= P(\text{zawor: } \left| \frac{\bar{X}(n) - \theta}{\frac{\sigma}{\sqrt{n}}} \right| < c) \approx$$

$$\Phi(c) - \Phi(-c) = 2\Phi(c) - 1,$$

gdzie  $\Phi$  oznaczamy dystrybucję  $N(0,1)$

Jaki ten ma być szacownik zakresu (1)

$$2\Phi(c) - 1 = 0,95 \Leftrightarrow$$

$$\Phi(c) = 0,975 \equiv \underline{\underline{c = 1,96}}$$

(tablice rozkładu  $N(0,1)$ !)



-9-

Dlatego wtedy

$$|\bar{X}(n) - \theta| < \frac{c\sigma}{\sqrt{n}} = 1,96 \frac{\sigma}{\sqrt{n}},$$

gdzie  $\sigma$  jest estymowaną statystyką  $\hat{S} = \sqrt{S^2}$

Zatem warunkiem na margines rozproszenia ma postać:

$$1,96 \frac{\hat{S}}{\sqrt{n}} \leq d, \text{ gdzie } d -$$

akceptowalne rozproszenie.

Jest to n-linowa postać symulacji spektralnej  
ten pomiar warunków, do dla  
tego n konieczny proces symulacji  
i mamy

$$\bar{X} \approx \theta \approx \text{prawdop. } 0,95.$$

## Przykład 1

Zadany, w przedmiotem analizy p. firma  
usługowa F. O F wiadomo, w po godzinie 17.00  
nie powinny pojawiać się nowi klienci.

Chcąc oszacować oczekiwany czas, w którym klient  
F opóźni się do godziny 17.00.

Wymagamy, aby pewność tego oszacowania była  
na poziomie 95% oraz że szacunkowa wartość  
średnia różniła się od rzeczywistej mniej niż  
15 sekund.

Aby to zrobić, musimy generować (w symulacji  
symulacji) zdania odnośnie do oszacowania  
syberji. Ze względu na wymagania C.T.G.

liczba ta k powinna być  $\geq 100$  oraz  
powinno spełniać nierówność

$$1.96 \frac{\hat{S}}{\sqrt{k}} \leq 15, \text{ gdzie } \hat{S} \text{ miernik}$$

o w sekundach. Wtedy  $\theta \approx \bar{X}$ .

Podsumujmy:

Z teoretycznego punktu widzenia potrafią obsługiwać  
linijny cykl symulacji  $n$ , a więc zgodnie  
z oczekiwaniem esymacji wartości składowej  $\theta$   
zmiany losowej  $X$  będącej wynikiem <sup>tej</sup> symulacji.

Z praktycznego punktu widzenia (chociaż  
o przetwarzaniu na maszynie) byłoby  
lepiej, abyśmy do tej wartości  $n$  dochodzili  
metodą rekurencyjną, a nie tak jak

polecamy powyższą metodę — poprzez generowanie  
za każdym razem kilku wartości, a 7

wariantach  $1,96 \frac{s}{\sqrt{k}} < d$  były spełnione.

Kolejne  
próby

{  
• • • NIE  
• • • NIE

} NIE

• • • • • TAIL

Pokazujemy jak to można zrobić!

Niech:

(i)  $X$  - wyniki symulacji  
 $EX = \theta$ ,  $\text{var}(X) = \sigma^2$

(ii)  $X_1, X_2, \dots, X_j, \dots$

niezależny ciąg zm. losowych o rozkładzie  $X$  oraz niezależnych

Dla każdego  $j \geq 2$ , możemy

$$\bar{X}_j = \frac{1}{j} \sum_{i=1}^j X_i \quad \left( = \frac{X_1 + X_2 + \dots + X_j}{j} \right)$$

oraz  $\hat{S}_j^2 = \frac{1}{j-1} \sum_{i=1}^j (X_i - \bar{X}_j)^2$

Zauważmy, że

Zdefiniujmy  $(\hat{S}_1^2) = \hat{S}_1^2 = 0$ , oraz

-13-

$$\bar{X}_{j+1} = \frac{1}{j+1} \sum_{i=1}^{j+1} X_i = \frac{1}{j+1} \left( \sum_{i=1}^j X_i + X_{j+1} \right)$$

$$= \frac{1}{j+1} \sum_{i=1}^j X_i + \frac{X_{j+1}}{j+1} =$$

$$= \left\{ \frac{j}{j+1} \cdot \frac{1}{j} \sum_{i=1}^j X_i + \frac{X_{j+1}}{j+1} \right\} =$$

$$= \left( 1 - \frac{1}{j+1} \right) \bar{X}_j + \frac{X_{j+1}}{j+1}, \quad \text{czyli}$$

$$\boxed{\bar{X}_{j+1} = \bar{X}_j + \frac{X_{j+1} - \bar{X}_j}{j+1}} \quad (\#)$$

Podobnie

$$\boxed{\hat{S}_{j+1}^2 = \left( 1 - \frac{1}{j+1} \right) \hat{S}_j^2 + (j+1) (\bar{X}_{j+1} - \bar{X}_j)^2} \quad (\#\#)$$

Zależności (#) i (\#\#) daje nam de rekurencję.

Pomyłka 2

Pomyłka 1, 4  $X_1 = 5, X_2 = 14, X_3 = 9$ .

Wtedy (#) i (\#\#) mają postać:

-1h-

$$\bar{X}_1 = 5$$

$$\bar{X}_2 = \bar{X}_1 + \frac{X_2 - \bar{X}_1}{2} = 5 + \frac{14 - 5}{2} = 5 + \frac{9}{2} = \frac{19}{2}$$

$$\begin{aligned} \hat{S}_2^2 &= \left(1 - \frac{1}{2}\right) \hat{S}_1^2 + (1+1) (\bar{X}_2 - \bar{X}_1)^2 \\ &= 2 \left(\frac{19}{2} - 5\right)^2 = \frac{81}{2} \end{aligned}$$

ifd.

Sytema przedstawiona wyżej wygląda mielo  
prościej, kiedy wartości danych są binarne,  
czyli  $X$  ma rozkład: 

0	1
1-p	p

, wtedy

$$EX = p, \text{ Var } X = p(1-p),$$

gdzie  $p$   $p'$  przedmiotem estymacji ( $\theta = p$ )

Każdy Zadań, i potrzebny wygenerować ciąg  
merokomach zm. losowych  $X_i$  o rozkładzie  
kardas 

0	1
1-p	p

Mamy mamy  $X_1, X_2, \dots, X_n$  o własności j.v.,  $p$

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \quad p \text{ estymator } p = E(X)$$

W takim razie dla  $\sigma^2 = \text{var}(X) = p(1-p)$ ,

$$\bar{X}_n(1 - \bar{X}_n) \quad p \text{ estymator } \sigma^2.$$

Można, w takim (diskretnym) przypadku procedurę  $p$  następująco:

- 1) Ustalić akceptowalną miarę rozproszenia  $d$
- 2) Generujemy co najmniej 1000 losów wyroków symulacji (CTG)

3) Kontynuujemy aż do chwili  $k$ , kiedy

$$\frac{\bar{X}_k(1 - \bar{X}_k)}{k} \leq d^2$$

Punkt 3

Pomysł, w obserwacji proces obsługi klientów firmy  $F$  jak w przykładzie 1.

Interesuje nas prawdop. że o 17.30 była obsługiwany klient.

Aby to uzyskać możemy w kolejnych dniach obserwować ruchy firmy. Robimy to za pomocą symulacji.

$$N_{i\text{dn}} \quad X_i = \begin{cases} 1, & i\text{-tyo dnia o 17.30} \\ & \text{p. klient} \\ 0, & \text{nie ma żadnego} \end{cases}$$

Dla  $i \geq 100$  bierzemy takie  $K = \frac{i}{21}$ .

$$\frac{P_K(1-P_K)^K}{K} \leq d, \text{ gdzie:}$$

$d$  - akceptowalna wartość approximationa

$P_K = \bar{X}_K$  - frekwencja dni dni, że klient był w  $F$  o 17.30.

### ZADANIE 1

Zaprogramować rekurencyjnie (#/i) (##).

### ZADANIE 2

Aby estymować wartość  $\theta$ , wygenerujemy 20 wartości zmiennych losowych o średniej  $\theta$ .



- 17 -

102, 112, 131, 107, 114, 95, 133, 145, 139, 117  
93, 117, 124, 122, 136, 141, 119, 122, 117, 143.

Spravidil čaj ta ilost' je dostatenno, jsh  
z pervoju 99% onelkup, y' voprosi' eshmerena  
votni n' od necymnih o me k'cej omizeli' o, j.