

Markov Inequality

$$\Pr(X \geq \alpha) \leq \frac{E(X)}{\alpha}$$

↓
nonnegative

Chebyshev Inequality

$$\Pr(|Y - \mu| \geq \alpha) \leq \frac{\text{Var}(Y)}{\alpha^2}$$

Example: $X \sim \text{Bin}(n, \frac{1}{2})$

$$\Pr(X \geq \frac{3}{4}n) \leq \frac{\frac{1}{2}n}{\frac{3}{4}n} = \frac{2}{3}$$

Markov

$$\Pr(X \geq \frac{3}{4}n) \leq \Pr(|X - \frac{1}{2}n| \geq \frac{1}{4}n) \leq \frac{\text{Var}(X)}{(\frac{1}{4}n)^2} = \frac{\frac{1}{4}n}{(\frac{1}{4}n)^2} = \frac{4}{n}$$

↓
0

Chebyshev

Sampling & Polling

What fraction of people approve of president?

Poll: call up n random people

$$X = X_1 + X_2 + \dots + X_n$$

Define $\bar{X} = \frac{X}{n}$ as our estimate

Questions:

What should n be?

how confident are we?

How good an estimate?

can we say my polling estimate 100% guaranteed to be within $\pm 2\%$ of truth.

Question: Given Θ , $1 - \epsilon$: how large does n need to be

so

$$\Pr(|\bar{X} - p| \leq \Theta) \geq 1 - \epsilon$$

0.04 0.95

↑
margin
of error

↑
confidence

Ex: how big does n need to be so 95% confident that \bar{X} is within

4% of truth (p)?

$$\text{Want } \Pr(|\bar{X} - p| < 0.04) \geq 0.95$$

$$\equiv \Pr(|\bar{X} - p| \geq 0.04) \leq 0.05$$

$$\leq \frac{\text{Var}(\bar{X})}{(0.04)^2} \quad \text{by Chebyshev}$$

$$\leq \frac{1}{0.04^2 4n}$$

$$\text{Suffices } \frac{1}{(0.04)^2 4n} \leq 0.05 \quad \Rightarrow \frac{1}{0.05 (0.04)^2 4} \leq n$$

$$n \geq 3125$$

$$\begin{aligned} \text{Var}(\bar{X}) &= \text{Var}\left(\frac{X}{n}\right) \\ &= \frac{1}{n^2} \text{Var}(X) \\ &= \frac{np(1-p)}{n^2} \\ &= \frac{p(1-p)}{n} \\ &\leq \frac{1}{4n} \end{aligned}$$

Notes:

- # of samples n doesn't depend on size of total population
- $[p - 0.04, p + 0.04]$ sometimes called confidence interval

100,000 computers

each indep sends packet w/ prob $q=0.01$ each sec

Router processes its buffer each sec

How many packet buffers so it drops a packet:

never: 100,000

With prob $\leq 10^{-6}$ each hour?

$$X_{it} = \begin{cases} 1 & \text{if computer } i \text{ sends a packet in } t^{\text{th}} \text{ second} \\ 0 & \text{otherwise} \end{cases}$$

$$X_t = \sum_{1 \leq i \leq n} X_{it} \quad \# \text{ of packets sent in a second}$$

B : size of buffer needed so that in T secs prob of overflow $\leq 10^{-6}$

for what B is

$$\Pr(\exists t \ 1 \leq t \leq T \text{ s.t. } X_t \geq B) \leq \epsilon$$

First, $\Pr(X_t \geq B)$

By Chernoff bound:

$$\Pr(X > \underbrace{(1+\delta)\mu}_{\text{\# of buffers } B}) \leq e^{-\frac{\delta^2 \mu}{2}}$$

$$B = (1+\delta)\mu = (1+\delta)(1000)$$

$$\Pr(\text{overflow in } T \text{ secs}) = \Pr(\exists t, 1 \leq t \leq T \text{ s.t. } X_t > (1+\delta)\mu)$$

$$\leq T e^{-\frac{\delta^2 \mu}{2}}$$

$$\leq T e^{-\delta^2 \mu / 2}$$

$$\text{We want } T e^{-\delta^2 \mu / 2} \leq \epsilon$$

$$\text{where } \mu = 1000$$

$$\epsilon = 10^{-6}$$

T : # secs in year

Recipe: solve for δ

use that to determine $B = (1+\delta)\mu$

$$T e^{-\delta^2 \mu / 2} \leq \epsilon$$

$$e^{\frac{\delta^2 \mu}{2}} \geq \frac{T}{\epsilon}$$

$$\frac{\delta^2 \mu}{2} \geq \ln\left(\frac{T}{\epsilon}\right)$$

$$\delta^2 \geq \frac{2}{\mu} \ln\left(\frac{T}{\epsilon}\right)$$

$$\delta \geq \sqrt{\frac{2}{\mu} \ln\left(\frac{T}{\epsilon}\right)}$$

$$B = (1+\delta)\mu$$

Example: $\mu = 1000$

$$T = 60 \cdot 60$$

min sec

$$\epsilon = 10^{-6}$$

$$\delta = \sqrt{\frac{2}{1000} \ln\left(\frac{3600}{10^{-6}}\right)} = 0.2097$$

$$\text{Buffer size} = 1.2097 \times 1000 \approx 1210$$

Twitter processes $\approx 600,000,000$ tweets per day 6×10^8

Assume tweets independent

$$\# \text{ chars/tweet} \sim U[10, 140]$$

What is \approx prob that twitter processes between 3.899×10^{10}
characters? Use CLT & 3.90001×10^{10}

$$X = X_1 + X_2 + \dots + X_n \quad n = 6 \times 10^8 \quad X_i \sim U[10, 140]$$

$$E(X) = n \cdot 80 = 6 \times 10^8 \cdot 80 = 3.9 \times 10^{10}$$

$$\text{Var}(X) = \text{Var}(X_1 + X_2 + \dots + X_n) = n \text{Var}(X_1) = 6 \times 10^8 \cdot 1430$$

$$\sigma = 926283$$

$$\frac{X - 3.9 \times 10^{10}}{926283} \sim N(0, 1)$$