

Recall: CLT:

X_1, \dots, X_n, \dots iid (independent + identically distrib)

R.V. and $S_n = \sum_{i=1}^n X_i$. Assume $EX_i = \mu$ and

$\text{Var}(X_i) = \sigma^2$. Then $\frac{S_n - n\mu}{\sigma\sqrt{n}} \xrightarrow{\text{dist}} N(0,1)$

idea: $S_n \approx N(n\mu, n\sigma^2)$

⊗ Binomial is sum of Bernoulli R.V.s so CLT applies.

Q: Half the people support a claim.

A poll asks 1000 people their opinion.

what is $P(\text{at least } 55\% \text{ Yes})$?

$X \sim \text{Bin}(1000, \frac{1}{2})$ want $P(X \geq 550)$

$$Z = \frac{X - EX}{\sqrt{\text{Var } X}} = \frac{X - 500}{\sqrt{1000 \cdot \frac{1}{2} \cdot \frac{1}{2}}} \text{ is roughly } N(0,1)$$

$$X \geq 550 \Leftrightarrow \frac{X - 500}{\sqrt{250}} \geq \frac{50}{\sqrt{250}} = 3.18\dots$$

$$P(X \geq 550) \approx P(N(0,1) \geq 3.18) = 1 - \Phi(3.18) < 0.001$$

e.g. Each job takes a server $\text{Exp}(1)$ to finish.

what is $P(50 \text{ jobs take } \geq 60 \text{ time})$?

Sol: $X_i = \text{time for job } i$. ~~Each~~ Each is $\text{Exp}(1)$

need $P(S_{50} \geq 60)$

note: S_n is $\text{Gamma}(n, 1)$.

Here $\mu=1$ and $\sigma^2=1$

$$\text{Exact: } P(S_{50} \geq 60) = \int_{60}^{\infty} \frac{x^{49}}{49!} e^{-x} dx = 0.084$$

$$\text{CLT: } P(S_{50} \geq 60) = P\left(\frac{S_{50} - \mu n}{\sqrt{n} \sigma} \geq \frac{60 - n\mu}{\sqrt{n} \sigma}\right)$$

$$\approx P(N(0,1) \geq \underbrace{\frac{60 - n\mu}{\sqrt{n} \sigma}}_{1.4}) = 1 - \Phi(1.4) = 0.079$$

$$\frac{60 - 50 \cdot 1}{\sqrt{50} \cdot 1} \approx 1.4$$

Statistics

def: Sample mean of samples X_1, \dots, X_n is

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

If (X_i) are iid samples from some distrib with expectation μ , then \bar{X} is an estimator for μ .

Claim: \bar{X} is an unbiased estimator: $E\bar{X} = \mu$

eg $(X_i) = \{3, 1, 4, 1, 5, 9\}$ $\bar{X} = \frac{23}{6}$

def: Sample variance is $s^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2$

Claim: s^2 is unbiased estimator for $\text{Var}(X_i)$

$$\text{i.e. } E s^2 = \text{Var}(X_i)$$

Proof: Let $S = \sum_{i=1}^n X_i$ so $\bar{X} = \frac{S}{n}$.

$$E S = n\mu \quad \text{where } \mu = E X_i$$

$$E s^2 = E\left(\sum_{i=1}^n X_i^2 + \sum_{i \neq j} X_i X_j\right) = n E(X_i^2) + 2 \binom{n}{2} E(X_i)(E X_j) \underbrace{=}_{\mu^2}$$

$$\begin{aligned} E\left(\sum_{i=1}^n (X_i - \bar{X})^2\right) &= E\left(\sum_{i=1}^n X_i^2 - 2 \sum_{i=1}^n X_i \bar{X} + n \bar{X}^2\right) \\ &= n E(X_i^2) - E \sum_{i=1}^n X_i \bar{X} + n \bar{X}^2 \\ &= n E(X_i^2) - E \sum_{i=1}^n X_i \frac{S}{n} + n \left(\frac{S}{n}\right)^2 \end{aligned}$$

$$= n E(X_i^2) - E(X_i) E(S) + n \frac{E(S^2)}{n} = (n-1) E(X_i^2) - (n-1) \mu^2 = (n-1) (E(X_i^2) - \mu^2)$$

$$= (n-1) \text{Var}(X) \quad \square$$

Hypothesis testing:

Get some data. Hypothesis: This comes from a given distrib. Determine the likelihood of data (or better, given hypothesis. (This is p-value).
If p-value small, reject hypothesis

$$\text{Likelihood} = P(\text{Data} | \text{Hyp})$$

e.g. ~~For~~ Hyp: Coin is fair.

Data: 550 Heads, 450 Tails in 1000 tosses

$$P(\text{Data} | \text{Hyp}) < 10^{-3}$$

Note: we want p-value = $P(\geq 550 \text{ Heads} \mid H_{\text{hyp}})$

Rule out the H_{hyp} . if $X = \# \text{heads}$ satisfies

$|X - 500| \geq l$ where l is such that

$$P(|X - 500| \geq l \mid \text{fair coin}) = 5\% \quad (\text{level of certainty})$$

For 1000 coins, this l is 1.96 std dev.

$$\text{So } l = 31$$