

Stat 98 lec 33

Warmup 2:00 - 2:10

a) if $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$

What distribution is \bar{X} (include parameters)

$$\bar{X} \sim N\left(p, \frac{pq}{n}\right), q=1-p$$

b) If $X_1, \dots, X_{200} \sim \text{Bernoulli}(p_x)$ } independent
 $Y_1, \dots, Y_{300} \stackrel{iid}{\sim} \text{Bernoulli}(p_y)$ } samples

What distribution is $\bar{X} - \bar{Y}$ (include parameters)?

$$\bar{X} - \bar{Y} \sim N\left(p_x - p_y, \frac{p_x q_x}{200} + \frac{p_y q_y}{300}\right)$$

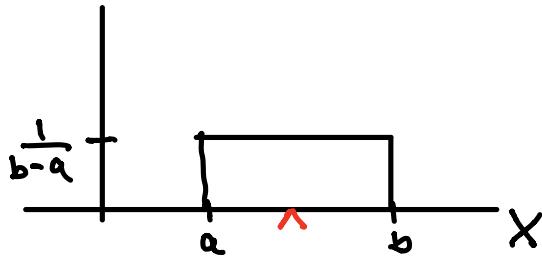
last time

Chap 10 continuous probability distributions

uniform

$$X \sim \text{Unif}(a, b)$$

$$f(x) = \begin{cases} \frac{1}{b-a} & a < x < b \\ 0 & \text{else} \end{cases}$$



$$E(X) = \frac{b+a}{2}$$

$$\text{SD}(X) = \frac{(b-a)}{\sqrt{12}}$$

* Gauss model for measurement error.

Let

$\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n \stackrel{iid}{\sim} \text{Unif}(-1, 1)$ are n independent measurement errors.

For example $-0.3, 0.7, -0.2$, etc

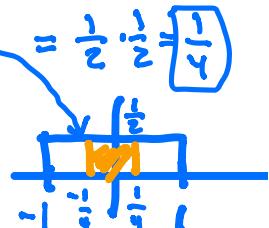
All numbers between -1 and 1 are "equally likely".

What is the chance the absolute value of the error is less than $\frac{1}{4}$? $\rightarrow P(|\varepsilon| < \frac{1}{4}) = \text{area} = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$

exponential

$$X \sim \text{Exp}(\lambda)$$

$$f(x) = \lambda e^{-\lambda x}, x > 0, \lambda > 0$$



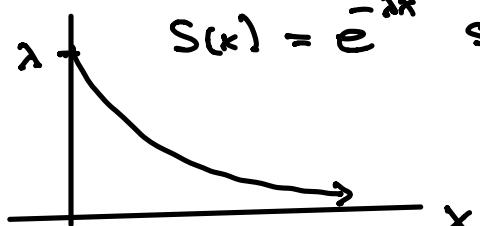
$$E(X) = \frac{1}{\lambda}$$

$$\text{SD}(X) = \frac{1}{\lambda}$$

$$h = \frac{\log(2)}{\lambda}$$

half life

$$S(x) = e^{-\lambda x} \quad \text{survival function}$$



X = time until failure of a mechanical device.

Normal

$$X \sim N(\mu, \sigma^2)$$

$$E(X) = \mu$$

$$SD(X) = \sigma$$

if $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu_x, \sigma_x^2)$ }
 $Y_1, \dots, Y_m \stackrel{iid}{\sim} N(\mu_y, \sigma_y^2)$ } indep

$$\bar{X} - \bar{Y} \sim N(\mu_x - \mu_y, \frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m})$$

or Test $H_0: \mu_x - \mu_y = 0$

$$H_1: \mu_x - \mu_y \neq 0$$

see whether $0 \in \bar{X} - \bar{Y} \pm 2\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}}$

or P-val $|\bar{X} - \bar{Y}|$ is $< .05$

Today

- ① Sec 10.4 CI for the difference between proportions and Test for the equality of two proportions
- ② extra practice problems
- ③ Sec 11.1 Bias and Variance,

① Sec 10.4 CI for the difference between proportions

This is a special case of before where now populations are 0 and 1s.

$$\begin{array}{l} X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p_x) \\ Y_1, \dots, Y_m \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p_y) \end{array} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{indep samples}$$

$\hat{X} - \hat{Y}$ is an unbiased estimator for $p_x - p_y$

$$\hat{X} - \hat{Y} \sim N(p_x - p_y, \frac{p_x q_x}{n} + \frac{p_y q_y}{m})$$

Now suppose we have independent samples from two cities, size $n=400$, and $m=600$ for city X and city Y .

- 37% of the City X sample are undecided about who they want as President
- 28% of the City Y sample are undecided about who they want as President

$$\text{we approximate } p_x \approx .37$$

$$p_y \approx .28$$

then a 95% CI for $p_x - p_y$ is

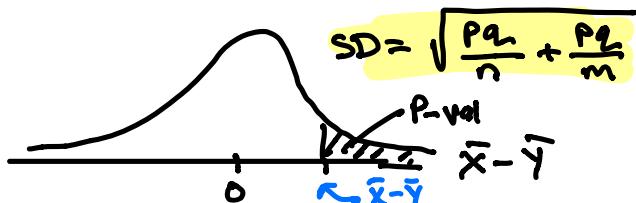
$$(37 - 28) \pm 2 \sqrt{\frac{.37(.63)}{400} + \frac{.28(.72)}{600}} = (.029, .15)$$

Test for equality of two proportions:

$$H_0: p_x = p_y = p$$

$$H_1: p_x > p_y \quad (\text{or } p_x \neq p_y \quad \text{2-sided alternative})$$

T.S. $\bar{X} - \bar{Y}$



Assuming the null, we reject null if

$$P\text{-val} = 1 - \Phi\left(\frac{(\bar{X} - \bar{Y}) - 0}{\sqrt{\frac{pq}{n} + \frac{pq}{m}}}\right) < .05$$

example 10.5.7

7. A simple random sample of 200 students is taken at University A. Independently, a simple random sample of 300 students is taken at University B.

- $\hat{P}_A = 20\%$ of the students in Sample A are football fans, whereas $\hat{P}_B = 50\%$ of the students in Sample B are football fans.

Test the hypothesis

$$H_0: P_A = P_B = 30\%$$

$$H_1: P_B > P_A$$

it is just an accident that this is $50\% - 20\% = \bar{B} - \bar{A}$

T.S. $\bar{B} - \bar{A}$

$$SD(\bar{B} - \bar{A}) = \sqrt{\frac{(3)(.7)}{200} + \frac{(3)(.7)}{300}} = .04$$

$$P\text{-val} = 1 - \Phi\left(\frac{-3 - 0}{.04}\right) = 1 - \Phi(-7.5) \approx 0 < .05$$

\Rightarrow reject null.

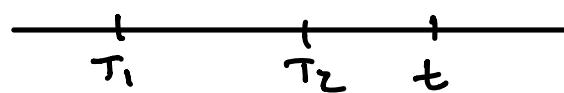
② Extra problems

\leq

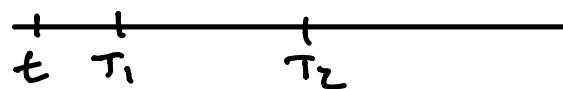
Suppose the time until the next earthquake in a particular place is exponentially distributed with rate 1 per year. Find the probability that the next earthquake happens within x days.

$$\begin{aligned}
 T &\sim \text{Exp}(1) \quad \text{rate in years} \\
 P(T < \frac{x}{365}) &= \int_0^{\frac{x}{365}} 1 \cdot e^{-1t} dt \\
 &= -e^{-\frac{x}{365}} + e^0 \\
 &= \boxed{-e^{-\frac{x}{365}} + 1} \leftarrow F\left(\frac{x}{365}\right)
 \end{aligned}$$

Trick to find cdf of min and max :



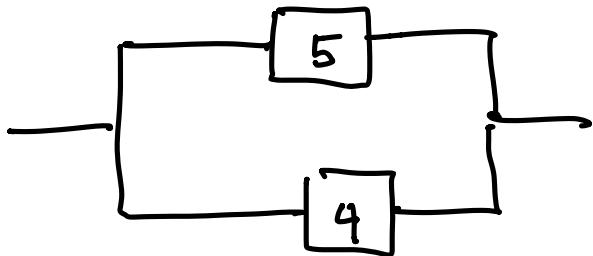
$$P(\max(T_1, T_2) < t) = P(T_1 < t, T_2 < t)$$



$$P(\min(T_1, T_2) > t) = P(T_1 > t, T_2 > t)$$

ex

An electrical circuit consists of 2 components in the following diagram.



The lifetimes of the components, measured in days, have independent exponential distributions with means given in the diagram

Let T be the lifetime of the circuit.
Find the cdf of T .

$$\left. \begin{array}{l} T_1 \sim \text{Exp}(\frac{1}{5}) \\ T_2 \sim \text{Exp}(\frac{1}{4}) \end{array} \right\} \text{indep}$$

$$T = \max(T_1, T_2) \quad \text{since circuit works until both components die.}$$

$$\begin{aligned} F(t) &= P(T < t) = P(T_1 < t, T_2 < t) \\ &= P(T_1 < t) P(T_2 < t) \\ &= \boxed{\left(1 - e^{-\frac{1}{5}t}\right) \left(1 - e^{-\frac{1}{4}t}\right)} \end{aligned}$$

(3) Sec 11.1 Bias and Variance,

Suppose we are trying to estimate a constant numerical parameter, Θ , and our estimator is the statistic T .

Below Θ is red

T is blue for different samples

what are the two best estimators?



Let's make a quantitative analysis.

Mean square error

$$MSE_{\Theta}(T) = E_{\Theta}((T - \Theta)^2)$$

$\nwarrow \Theta \text{ is true Parameter}$

Think of this as the average distance squared of T from θ , we want $MSE_{\theta}(T)$ to be as small as possible.

Decomposition of Error

Bias:

$$\text{let } B_{\theta}(T) = E_{\theta}(T) - \theta \quad \left(\begin{array}{l} \text{unbiased means} \\ E_{\theta}(T) = \theta \end{array} \right)$$

This is a constant !! bias

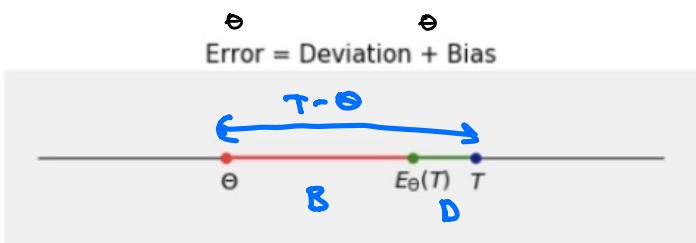
Deviation:

$$\text{let } D_{\theta}(T) = T - E_{\theta}(T) \quad \left(\begin{array}{l} \text{deviations of } T \\ \text{from the mean} \end{array} \right)$$

This is a RV !! deviation

$$\text{What is } E_{\theta}(D_{\theta}(T)) ? \quad \begin{aligned} E(D(T)) &= E(T - E(T)) \\ &= E(T) - E(T) = 0 \end{aligned}$$

$$\text{What is } E_{\theta}(D_{\theta}^2(T)) ? \quad \begin{aligned} E(D^2(T)) &= E((T - E(T))^2) \\ &= \text{Var}(T) \end{aligned}$$



Bias-variance Decomposition

$$\begin{aligned}
 MSE_{\theta}(T) &= E_{\theta}((T - \theta)^2) \\
 &= E_{\theta}((D_{\theta}(T) + B_{\theta}(T))^2) \\
 &= E_{\theta}(D_{\theta}^2(T) + 2D_{\theta}(T)B_{\theta}(T) + B_{\theta}^2(T))
 \end{aligned}$$

Simplify

$$= E(D^2(\tau)) + 2B(\tau)E(D(\tau)) + B^2(\tau)$$

$$= \boxed{v_{av}(T) + B^2(T)}$$

A small, dark blue ink smudge or mark is located near the top center of the page.