

⑤ RANDOM FUNCTIONS ASSOCIATED w THE NORMAL DISTRIBUTION

x_1, x_2, \dots, x_n are independent random variables with $\{\mu_1, \mu_2, \dots, \mu_n\}$ means and $y = \sum_{i=1}^n c_i x_i \Rightarrow y \sim N\left(\sum_{i=1}^n c_i \mu_i, \sum_{i=1}^n c_i^2 \sigma_i^2\right)$ $\left(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2\right)$ variances

Ex: $x_1 \sim N(\mu=2, \sigma^2=3)$, $x_2 \sim N(\mu=1, \sigma^2=4)$ x_1, x_2 are not.

Distribution of y ?

- $y = 2x_1 + 3x_2 \Rightarrow y \sim N(2 \cdot 2 + 3 \cdot 1, 2^2 \cdot 3 + 3^2 \cdot 4) = N \sim (7, 48)$
 - $y = x_1 - x_2 \Rightarrow y \sim N(2 - 1, (1)^2 \cdot (3) + (-1)^2 \cdot (4)) = N \sim (1, 7)$
 - $y = x_1 + x_2 \Rightarrow y \sim N(2 + 1, (1)^2 \cdot (3) + (1)^2 \cdot (4)) = N \sim (3, 7)$
 - $y = x_1 \cdot x_2 \Rightarrow y \sim N(\mu_1 - \mu_2, \sigma_1^2 + \sigma_2^2)$
- σ^2 is the same!

SAMPLING DISTRIBUTION:

when x_1, x_2, \dots, x_n are observations from a random sample of n obs. from a normal dist. w mean μ and variance $\sigma^2 \sim N(\mu, \sigma^2)$, then

• SAMPLE MEAN \bar{x}

$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ is normally distributed w $N\left(\mu, \frac{\sigma^2}{n}\right)$

$$\bar{x} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

the larger n , the smaller the variance $\frac{\sigma^2}{n}$

• SAMPLE VARIANCE s^2

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

- \bar{x} and s^2 are independent

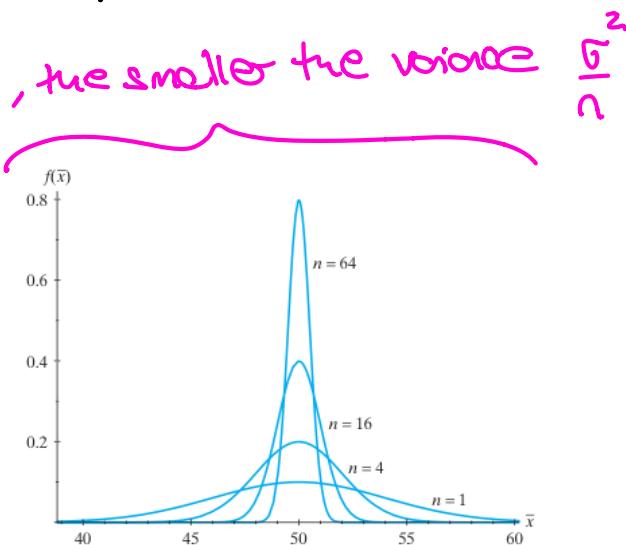


Figure 5.5-1 pdfs of means of samples from $N(50, 16)$

when sampling from a normal distribution :

sum of squared differences of:

• population mean μ

$$\frac{\sum_{i=1}^n (x_i - \mu)^2}{\sigma^2} \sim \chi^2_{(n)}$$

$$\frac{n(\bar{x} - \mu)}{\sigma^2} \sim \chi^2_{(1)}$$

• sample mean \bar{X}

we lose 1 df
when we estimate the μ
with \bar{X}

$$\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sigma^2} = \frac{(n-1) s^2}{\sigma^2} \sim \chi^2_{(n-1)}$$

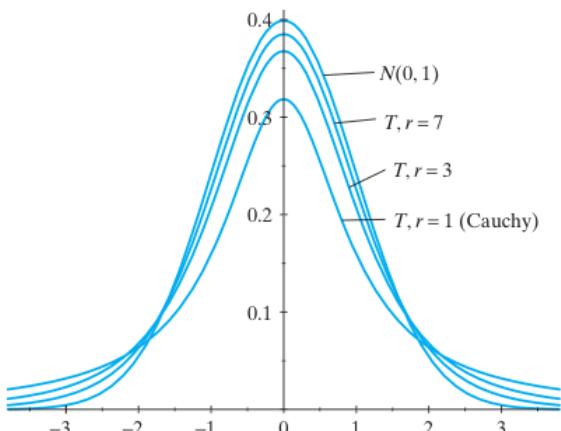
$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{(n-1)} \Rightarrow \sum_{i=1}^n (x_i - \bar{x})^2 = s^2 (n-1)$$

STUDENT'S t DISTRIBUTION:

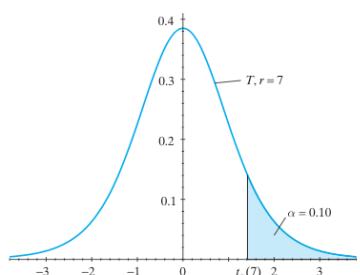
$$T = \frac{Z}{\sqrt{U/r}},$$

where Z is a random variable that is $N(0, 1)$, U is a random variable that is $\chi^2(r)$, and Z and U are independent. Then T has a t distribution with pdf

$$f(t) = \frac{\Gamma((r+1)/2)}{\sqrt{\pi r} \Gamma(r/2)} \frac{1}{(1+t^2/r)^{(r+1)/2}}, \quad -\infty < t < \infty.$$



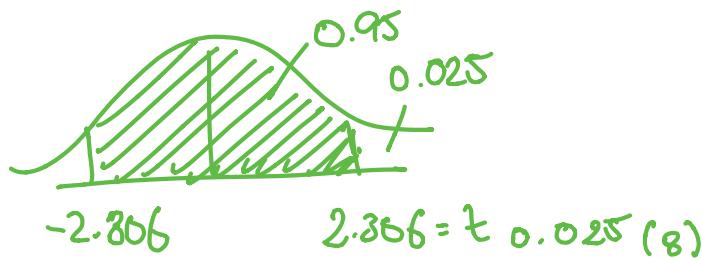
- Support δ is $-\infty < t < \infty$
- $\text{pdf symmetric around } 0$
- bell-shaped
- density curve \sim standard normal curve but with "heavier" tails : it is more likely to get extreme values in the t dist.
- as the df r increase $\rightarrow t$ -dist approaches the $N(0, 1)$ dist



Ex.

$$t(8) \quad P(|T| < 2.306) = P(-2.306 < T < 2.306) =$$

$$P(T < 2.306) - P(T > 2.306) = 0.975 - 0.025 = 0.95$$



RANDOM VARIABLE T

Given a random sample x_1, x_2, \dots, x_n from a normal dist. $N(\mu, \sigma^2)$:

- $Z = \frac{\bar{X} - \mu}{\sqrt{\frac{\sigma^2}{n}}} \sim N(0, 1)$

- $V = \frac{(n-1)s^2}{\sigma^2} \sim \chi^2_{(n-1)}$

- Z and V are independent

$$T = \frac{\bar{X} - \mu}{\sqrt{\frac{s^2}{n}}} = \frac{\bar{X} - \mu}{\sqrt{\frac{(n-1)s^2}{n^2}}} = \frac{\bar{X} - \mu}{\sqrt{\frac{n-1}{n}} s} = \frac{\bar{X} - \mu}{\sqrt{\frac{s}{\sqrt{n}}}} \sim t_{(n-1)}$$

T has a Student's t dist with $(n-1)$ df.

OTHER DISTRIBUTIONS :

- GAMMA: $\gamma(\alpha, \theta)$

$$Y \sim \gamma(\alpha = \frac{r}{2}, \theta = 2) \rightarrow Y \sim \chi^2_{(r)}$$

If y follows a γ w $\alpha = \frac{r}{2}$ ad $\theta = 2$ it follows a chi-squared r df.

- SUM OF INDEPENDENT χ^2 DISTRIBUTIONS $\sim \chi^2$ w sum of n df

$$Y_i \sim \chi^2_{(r_i)} \text{ w } i=1, 2, \dots, n \rightarrow Y = \sum_{i=1}^n Y_i \sim \chi^2_{\left(\sum_{i=1}^n r_i\right)}$$

- STANDARD NORMAL SQUARED $\sim \chi^2$ w 1 df

$$Z \sim N(0,1) \rightarrow Z^2 \sim \chi^2_{(1)}$$

- BETA DISTRIBUTION :

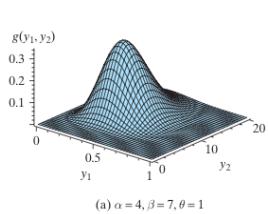
$$X_1 \sim \gamma(\alpha, \theta, \beta) \quad X_2 \sim \gamma(\alpha, \theta, \beta)$$

$$\begin{aligned} X_1 &= Y_1, Y_2 \\ X_2 &= Y_2 - Y_1, Y_2 \\ Y_1 &= \frac{X_1}{X_1 + X_2} \\ Y_2 &= X_1 + X_2 \end{aligned}$$

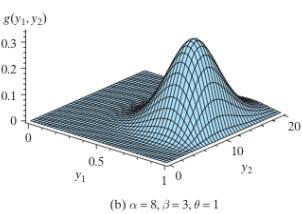
$$\text{Marginal pdf of } Y_1 = f(y_1, \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \cdot y_1^{\alpha-1} \cdot (1-y_1)^{\beta-1}$$

$$Y_1 \sim \beta(\alpha, \beta)$$

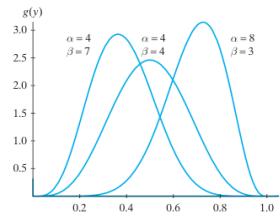
$$0 < y_1 < 1$$



(a) $\alpha = 4, \beta = 7, \theta = 1$



(b) $\alpha = 8, \beta = 3, \theta = 1$

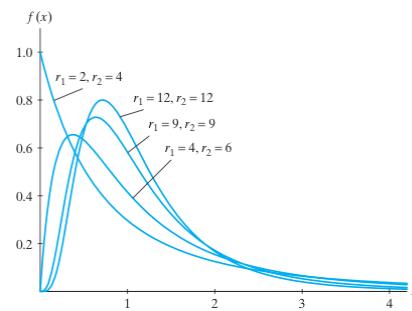


- F DISTRIBUTION: w r_1 and r_2 df.

$$Y_1 \sim \chi^2_{(r_1)} \quad Y_2 \sim \chi^2_{(r_2)}$$

$$\omega = \frac{Y_1}{Y_2} \quad \omega \sim F(r_1, r_2)$$

$$\text{pdf } f(\omega) = \frac{(\frac{r_1}{r_2})^{r_1/2} \Gamma[(r_1+r_2)/2]}{\Gamma(r_1/2) \Gamma(r_2/2) [1 + (\frac{r_1}{r_2}) \omega]^{(r_1+r_2)/2}} \omega^{r_1/2 - 1}$$



CENTRAL LIMIT THEOREM

The sampling distribution is approximately normally distributed regardless of the distribution of the underlying random sample when n is sufficiently large.

Random sample $x_1, x_2 \dots x_n$ from any distribution w/ μ and σ^2 w/ sufficiently large n :

- sample mean $\bar{X} \sim N(E(\bar{X}) = \mu_{\bar{X}} = \mu, \text{Var}(\bar{X}) = \sigma_{\bar{X}}^2 = \frac{\sigma^2}{n})$

$$\bar{X} \xrightarrow[n \rightarrow \infty]{\sim} N(\mu, \frac{\sigma^2}{n})$$

$$Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{\sum_{i=1}^n x_i - n\mu}{\sqrt{n}\sigma} \xrightarrow[n \rightarrow \infty]{\sim} N(0, 1)$$

APPROXIMATIONS FOR DISCRETE DISTRIBUTIONS

NORMAL APPROXIMATION TO BINOMIAL

$Y \sim b(n, p)$ is the sum of n $x_i \sim \text{Bennelli}(p)$

$y = \sum_{i=1}^n x_i$ random variables, $x_1, x_2 \dots x_n$

	mean	variance
$x_i \sim \text{Bennelli}(p)$ with	$\mu = E(x) = p$	$\sigma^2 = \text{Var}(x) = p(1-p) = pq$
$y \sim \text{binomial}(np)$ with	$\mu = np$	$\sigma^2 = np(1-p) = npq$

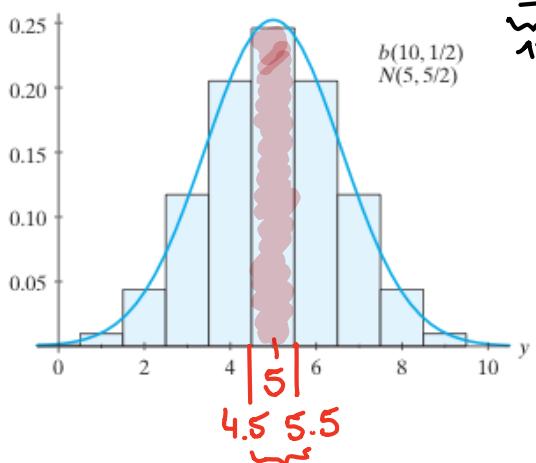
CLT: $Y \sim N(np, npq)$ if n is sufficiently large $\begin{cases} np \geq 5 \\ nq \geq 5 \end{cases}$

$$w = \frac{\bar{X} - \mu}{\sqrt{\frac{\sigma^2}{n}}} = \frac{\bar{X} - p}{\sqrt{\frac{pq}{n}}} \xrightarrow{n \rightarrow \infty} N(0, 1)$$

$$w = \frac{\sum_{i=1}^n x_i - np}{\sqrt{npq}} = \frac{Y - np}{\sqrt{n \cdot npq}} = \frac{Y - np}{\sqrt{npq}} \xrightarrow{n \rightarrow \infty} N(0, 1)$$

CONTINUITY CORRECTION

We approximate probabilities for discrete distributions with a continuous distribution. $P(Y=k)$ is represented by the area of the rectangle



height = $P(Y=k)$
base = length = 1 centered around k

$$P(Y=k) = P(k - \frac{1}{2} < Y < k + \frac{1}{2})$$

$$P(Y=5) = P(4.5 < Y < 5.5)$$

$$= P\left(\frac{4.5-5}{\sqrt{2.5}} < Z < \frac{5.5-5}{\sqrt{2.5}}\right) = P(-0.32 < Z < 0.32)$$

$$= P(Z < 0.32) - P(Z > 0.32) = 0.625 - 0.375 =$$

$$= 0.251$$

SAMPLE PROPORTION : $\hat{p} = \frac{\sum_{i=1}^n X_i}{n}$ is the proportion of the sample that meets condition of interest.

It can be estimated by :

$$Z = \frac{\sum_{i=1}^n X_i - np}{\sqrt{np(1-p)}} = \frac{\sum_{i=1}^n X_i - np}{\sqrt{\frac{np(1-p)}{n}}} = \frac{\frac{1}{\sqrt{n}}(\hat{p} - p)}{\frac{p(1-p)}{\sqrt{n}}} \xrightarrow{n \rightarrow \infty} N(0, 1)$$

NORMAL APPROXIMATION TO POISSON

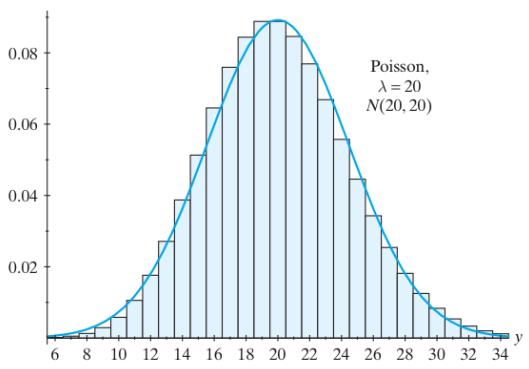
X_1, X_2, \dots, X_n are independent Poisson random variables

$X_i \sim \text{Poisson}(\lambda)$ with mean = 1

$Y = \sum_{i=1}^n X_i$ $Y \sim \text{Poisson}(\lambda)$ with mean = λ and variance = λ

CLT : $Y \sim \text{Poisson}(\lambda)$

$W = \frac{Y - \lambda}{\sqrt{\lambda}} \xrightarrow{n \rightarrow \infty} N(0, 1)$ when λ is sufficiently large



ESTIMATION

Random sample :

observed values of a specific random sample: $x_1, x_2 \dots x_n$
lowercase

Random variables arising from a random sample: $X_1, X_2 \dots X_n$
uppercase

- PARAMETER: θ
- PARAMETER SPACE: Ω range of possible values of parameter θ
Ex: proportion of students who smoke $\Omega = \{p: 0 \leq p \leq 1\}$
- STATISTIC: the function of $x_1, x_2 \dots x_n$ used to estimate parameter θ
Ex: the statistic of relative frequency $\frac{\text{no of success}}{n}$ estimates p .
- POINT ESTIMATOR: the statistic $\hat{u}(x_1, x_2 \dots x_n)$ used to estimate θ

* Diff between point est and statistic

Ex: $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ is a point estimator of population μ

$\hat{p} = \frac{1}{n} \sum_{i=1}^n x_i$ " population proportion p

- POINT ESTIMATE: the function of $u(x_1, x_2 \dots x_n)$ computed from a set of data is an observed point estimate of θ

Ex: x_i = observed grade point avg. of $n = 88$ students

$\bar{x} = \frac{1}{88} \sum_{i=1}^{88} x_i = 3.12$ is a point estimate of the pop. μ .

- ESTIMAND: what is being estimated, often a parameter θ

Ex: μ .

MAXIMUM LIKELIHOOD ESTIMATION

what for?

Random sample $X_1, X_2 \dots X_n$

with a prob. distribution that depends on parameter $\Theta \rightarrow X_i \sim N(\mu, \sigma^2)$

Goal: find a point estimator $u(X_1, X_2 \dots X_n)$

$$\text{ex} \\ X_1, X_2 \dots X_n \\ \bar{X} = \frac{\sum X_i}{n} \\ \text{to estimate the } \mu$$

such that $u(X_1, X_2 \dots X_n)$ is a "good" point estimate

the observed values of a random sample

LIKELIHOOD FUNCTION: $L(\Theta)$ (for 1 parameter)

Random sample: $X_1, X_2 \dots X_n$. X_i is independent bc it's random sample

the pdf for each X_i is $f(x_i; \Theta)$

Joint probability mass function of X_i independent vbles

$$L(\Theta) = P(X_1=x_1, X_2=x_2 \dots X_n=x_n) = f(x_1; \Theta) \cdot f(x_2; \Theta) \dots f(x_n; \Theta) = \prod_{i=1}^n f(x_i; \Theta)$$

Ex: Random sample $X_1, X_2 \dots X_n$. $X_i = 0 \rightarrow \text{no cor}$
 $X_i = 1 \rightarrow \text{cor}$

$X_i \sim \text{Bernoulli}(\rho) \rightarrow X_i$ are independent Bernoulli random vbles

pmf of each X_i : $f(x_i; \rho) = \rho^{x_i} (1-\rho)^{1-x_i}$ for $x_i=0 \text{ or } 1$
 $0 < \rho < 1$

$$L(\rho) = \prod_{i=1}^n f(x_i; \rho) = \rho^{\sum x_i} (1-\rho)^{n - \sum x_i} = \rho^{\sum x_i} (1-\rho)^{n - \sum x_i}$$

Find the ρ that maximizes the likelihood of $L(\rho)$

↪ differentiate the $L(\rho)$ (to find where the function has a maximum → review calculus)

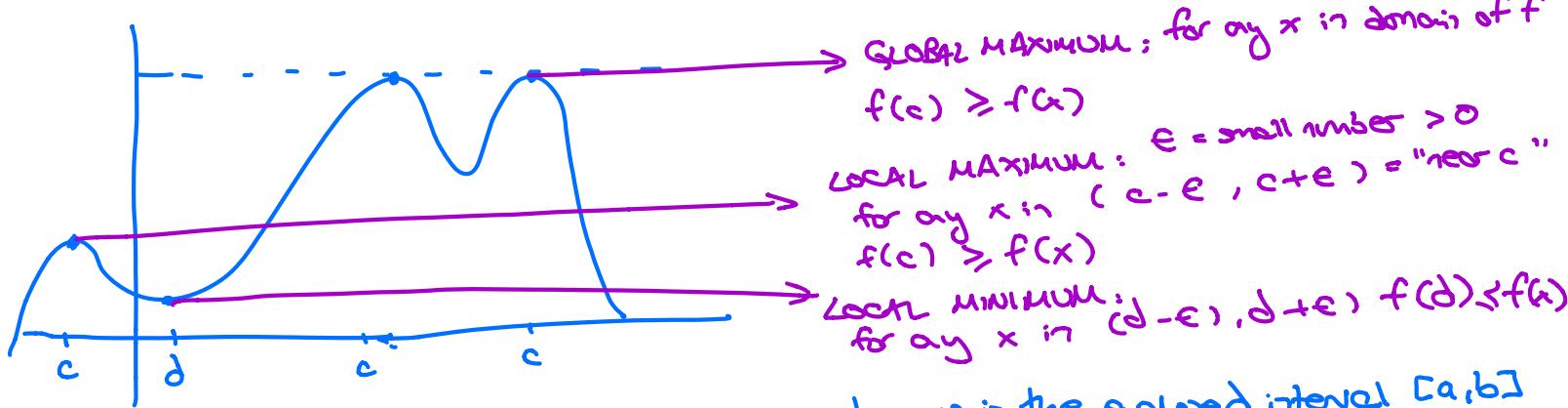
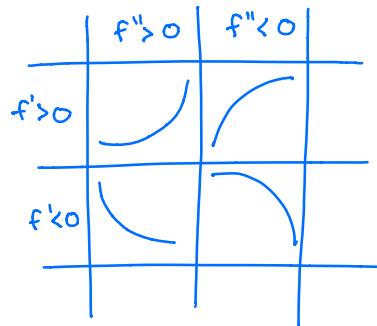
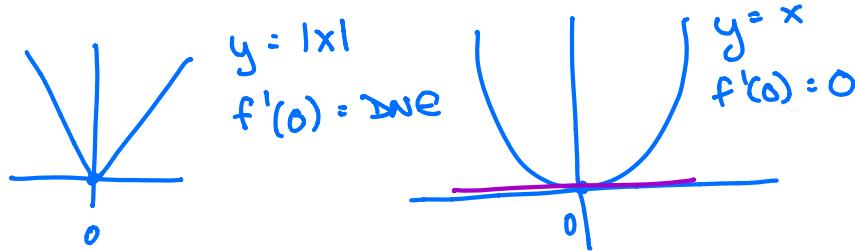
$$\frac{d}{d\rho} L(\rho) = L'(\rho) = 0$$

CALCULUS REVIEW :

EXTREUMS:

FERMAT'S THEOREM: Suppose f is a function defined in the interval $(a,b) \ni c \Rightarrow c$ is in (a,b) . If $f(c)$ is differentiable at c , $f'(c) \neq 0 \Rightarrow c$ is not an extreme value of f .

- $f'(c) = 0 \Rightarrow c$ is an extreme value of f = local extremum.
- $f'(c) = \text{DNE} \Rightarrow c$ is a critical point where the tangent is horizontal
- $f'(c) = \text{DNE} \Rightarrow c$ is an extreme value of f = local extremum.



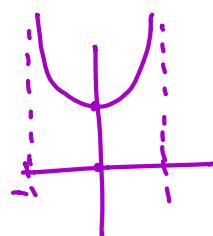
Extreme Value Theorem: if $f(x)$ is continuous in the closed interval $[a,b]$ then there are c and d in $[a,b]$ so that for all x in $[a,b]$

$f(x) \leq f(c) \Rightarrow f(x)$ attains a MAXIMUM VALUE
 $f(x) \geq f(d) \Rightarrow f(x)$ attains a MINIMUM VALUE

HOW TO FIND MAX MIN VALUES: ex: $f(x) = \frac{1}{(x^2-1)^2}$ on $(-1,1)$

- ① Differentiate
- ② List critical points
- ③ Check them
- ④ Check limiting behaviour

- ① $f'(x) = -2(x^2-1)^{-3} \cdot 2x = \frac{-4x}{(x^2-1)^3} = \text{DNE}$
- ② $-4x = 0 \Rightarrow x = 0$
- ③ $\begin{array}{c} f'(x) \\ \hline -1 & 0 & 1 \end{array}$
- ④ $\lim_{x \rightarrow -1^+} f(x) = \infty$
 $\lim_{x \rightarrow 1^+} f(x) = \infty$



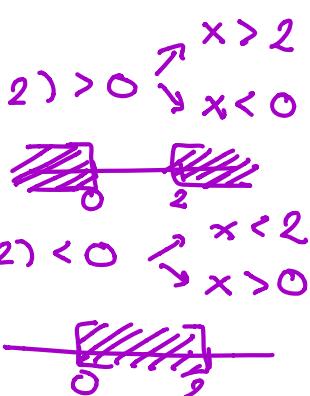
CALCULUS REVIEW CONT'

Ex: $f(x) = x - |x^2 - 2x|$ in $[0, 3]$

$$f(x) = \begin{cases} x - (x^2 - 2x) & \text{if } x^2 - 2x \geq 0 \\ x + (x^2 - 2x) & \text{if } x^2 - 2x < 0 \end{cases}$$

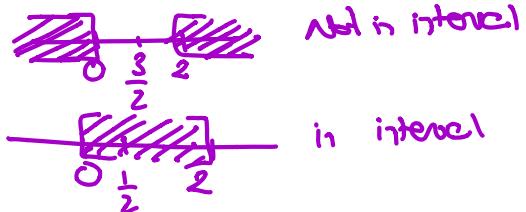
$$\textcircled{1} \quad f'(x) \begin{cases} \textcircled{1} \quad f'(x) = 1 - (2x-2) = 3-2x \text{ if } x^2-2x \geq 0 \Rightarrow x(x-2) \geq 0 \rightarrow \begin{array}{l} x > 2 \\ x < 0 \end{array} \\ \textcircled{2} \quad f'(x) = 1 + (2x-2) = -1+2x \text{ if } x^2-2x < 0 \Rightarrow x(x-2) < 0 \rightarrow \begin{array}{l} x < 2 \\ x > 0 \end{array} \end{cases}$$

$$\left. \begin{array}{ll} x=2 & f'(2)^{\textcircled{1}} = 3-4 = -1 \neq \\ & f'(2)^{\textcircled{2}} = 3 \\ x=0 & f'(0)^{\textcircled{1}} = 3 \neq \\ & f'(0)^{\textcircled{2}} = -1 \end{array} \right\} f'(x) = \text{DNE} \quad \text{not differentiable}$$



$$\textcircled{2} \quad f'(x) = 0 \quad \textcircled{1} \quad 3-2x = 0 \Rightarrow x = \frac{3}{-2} = \frac{3}{2}$$

$$\boxed{x = \frac{1}{2}} \quad \textcircled{2} \quad -1+2x = 0 \Rightarrow x = \frac{1}{2}$$



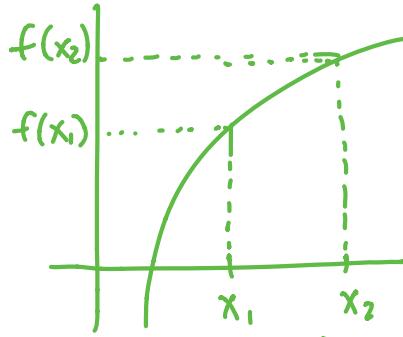
Do derivative without taking log:

$$\begin{aligned}
 L'(p) &= \frac{\partial}{\partial p} p^{\sum_{i=1}^n x_i} \cdot (1-p)^{n-\sum_{i=1}^n x_i} = \text{product rule } \frac{\partial}{\partial x} f(x) \cdot g(x) = f'(x) \cdot g(x) + g'(x) f(x) \\
 &= \frac{\partial}{\partial p} p^{\sum_{i=1}^n x_i} \cdot (1-p)^{n-\sum_{i=1}^n x_i} + p^{\sum_{i=1}^n x_i} \cdot \boxed{\frac{\partial}{\partial p} (1-p)^{n-\sum_{i=1}^n x_i}} \\
 &\quad \downarrow \\
 &= \sum_{i=1}^n x_i p^{\sum_{i=1}^n x_i - 1} \cdot (1-p)^{n-\sum_{i=1}^n x_i} + p^{\sum_{i=1}^n x_i} \cdot (n - \sum_{i=1}^n x_i)(1-p)^{n-\sum_{i=1}^n x_i - 1} \cdot (-1) = \\
 &= \sum_{i=1}^n x_i p^{\sum_{i=1}^n x_i - 1} \cdot (1-p)^{n-\sum_{i=1}^n x_i} - p^{\sum_{i=1}^n x_i} \cdot (n - \sum_{i=1}^n x_i)(1-p)^{n-\sum_{i=1}^n x_i - 1} = 0 \quad 0 < p < 1 \\
 &= \sum_{i=1}^n x_i \underbrace{p^{\sum_{i=1}^n x_i - 1}}_{\text{f'(g(x))}} \cdot \underbrace{p}_{\text{g(x)}} \cdot \underbrace{(1-p)^{n-\sum_{i=1}^n x_i}}_{\text{g'(x)}} = \\
 &\therefore p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i} \left[\frac{\sum_{i=1}^n x_i}{p} - \frac{n - \sum_{i=1}^n x_i}{1-p} \right] = 0 \quad \text{thus } = 0 \text{ when} \\
 &\frac{\sum_{i=1}^n x_i}{p} - \frac{n - \sum_{i=1}^n x_i}{1-p} = 0 \Rightarrow \frac{\sum_{i=1}^n x_i}{p} = \frac{n - \sum_{i=1}^n x_i}{1-p} \Rightarrow \sum_{i=1}^n x_i (1-p) = (n - \sum_{i=1}^n x_i)p \\
 &\sum_{i=1}^n x_i - np = 0 \Rightarrow p = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}
 \end{aligned}$$

$L''(\bar{x}) < 0$ so $L(\bar{x})$ is a maximum

The statistic $\hat{p} = \frac{\sum_{i=1}^n x_i}{n} = \bar{X}$ is the **MAXIMUM LIKELIHOOD ESTIMATOR**.

To make differentiation easier take the logarithm of the $f(x_i, p)$



$y = \ln(x)$ the natural logarithm is an increasing function of x
 \rightarrow if $x_1 < x_2 \Rightarrow f(x_1) < f(x_2)$

so the value of p that maximizes
 the logarithm of the likelihood function $\ln(L(p))$
 also maximizes the likelihood function $L(p)$

$$L(p) = p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i}$$

CHAIN RULE

$$\frac{\partial}{\partial p} \ln(1-p) = \frac{1}{1-p} \cdot (-1)$$

$$\ln(L(p)) = \sum_{i=1}^n x_i \cdot \ln(p) + (n - \sum_{i=1}^n x_i) \cdot \boxed{\ln(1-p)}$$

$$\frac{\partial \ln(L(p))}{\partial p} = \frac{\sum x_i}{p} - \frac{(n - \sum x_i)}{1-p} = 0 \Rightarrow \frac{\sum x_i}{p} - \frac{(n - \sum x_i)}{1-p} \cdot p(1-p) = 0$$

$$\sum x_i (1-p) - (n - \sum x_i) \cdot p = 0 \Rightarrow \sum x_i - \cancel{\sum x_i p} - np + \cancel{\sum x_i p} = 0$$

$$\sum x_i - np = 0 \Rightarrow \hat{p} = \frac{-\sum x_i}{-n} = \frac{\sum x_i}{n} \rightarrow \text{estimate (lower case X)}$$

$$\hat{p} = \frac{\sum x_i}{n} \rightarrow \text{estimator (upper case X)}$$

To verify we obtained a maximum $\rightarrow f''(\ln(p)) < 0$

Given a random sample $x_1, x_2 \dots x_n$ from a distribution that depends on $\theta_1, \theta_2 \dots \theta_m$
 unknown parameters w/ $\left\{ \text{pmf } f(x_i; \theta_1, \theta_2 \dots \theta_m) \right.$ then:
 restricted to parameter space \curvearrowright

• LIKELIHOOD FUNCTION: $L(\theta_1, \theta_2 \dots \theta_m) = \prod_{i=1}^n f(x_i; \theta_1, \theta_2 \dots \theta_m)$

• MAXIMUM LIKELIHOOD ESTIMATOR: $\hat{\theta}_i = u_i(x_1, x_2 \dots x_n)$ of θ_i for $i=1, 2 \dots m$
 if $[u_1(x_1, x_2 \dots x_n), u_2(x_1, x_2 \dots x_n) \dots u_m(x_1, x_2 \dots x_n)]$ is the m -tuple
 that maximizes $L(\theta_1, \theta_2 \dots \theta_m)$ (the likelihood function)

• MAXIMUM LIKELIHOOD ESTIMATES : the corresponding observed values of the statistics
 $[u_1(x_1, x_2 \dots x_n), u_2(x_1, x_2 \dots x_n) \dots u_m(x_1, x_2 \dots x_n)]$ of θ_i for $i=1, 2 \dots m$

For a given MLE / statistic $u(x_1, x_2, \dots, x_n)$,

$$u(x_1, x_2, \dots, x_n) \text{ is } \Theta$$

if $E(u(x_1, x_2, \dots, x_n)) = \Theta \rightarrow$ UNBIASED ESTIMATOR of parameter Θ

if $E(u(x_1, x_2, \dots, x_n)) \neq \Theta \rightarrow$ BIASED ESTIMATOR of parameter Θ

Ex. $x_i \sim \text{Bernoulli}(\rho)$. Is the MLE of $\rho \rightarrow \hat{\rho} = \frac{1}{n} \sum_{i=1}^n x_i$ $\begin{cases} \text{biased} \\ \text{unbiased} \end{cases}$
estimator?

$$E(\hat{\rho}) = E\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{1}{n} \sum_{i=1}^n E(x_i) = \frac{1}{n} \sum_{i=1}^n \rho = \frac{1}{n} \cdot np = \rho \rightarrow \text{UNBIASED!}$$

Ex. $x_i \sim N(\mu, \sigma^2)$ MLE $\hat{\mu} = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}$ or they $\begin{cases} \text{biased} \\ \text{unbiased} \end{math}$ estimators?

$$E(\hat{\mu}) = E\left(\frac{\sum_{i=1}^n x_i}{n}\right) = \frac{1}{n} E(\sum_{i=1}^n x_i) = \frac{1}{n} \sum_{i=1}^n E(x_i) = \frac{1}{n} \sum_{i=1}^n \mu = \frac{1}{n} n\mu = \mu \text{ UNBIASED!}$$

$$E(\hat{\sigma}^2) = ?$$

$$\begin{aligned} \hat{\sigma}^2 &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n (x_i^2 - 2x_i\bar{x} + \bar{x}^2) = \\ &= \frac{1}{n} \sum_{i=1}^n x_i^2 - 2\bar{x} \cdot \frac{1}{n} \sum_{i=1}^n x_i + \frac{1}{n} \sum_{i=1}^n \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - 2\bar{x}^2 + \frac{1}{n} \sum_{i=1}^n \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 \Rightarrow \end{aligned}$$

$$E(\hat{\sigma}^2) = E\left[\left(\frac{1}{n} \sum_{i=1}^n x_i^2\right) - \bar{x}^2\right] = E\left(\frac{1}{n} \sum_{i=1}^n x_i^2\right) - E(\bar{x}^2) = \frac{1}{n} \sum_{i=1}^n E(x_i^2) - E(\bar{x}^2) =$$

$$\text{var}(x) = \sigma^2 = E(x^2) - E(x)^2 = E(x^2) - \mu^2 \Rightarrow E(x^2) = \sigma^2 + \mu^2$$

$$\begin{aligned} \text{var}(\bar{x}) &= \frac{\sigma^2}{n} = E(\bar{x}^2) - \mu^2 \\ &= \frac{1}{n} \sum_{i=1}^n E(x_i^2) - \frac{1}{n} \sum_{i=1}^n \mu^2 = \frac{1}{n} (n\sigma^2 + n\mu^2) - \frac{\sigma^2}{n} - \mu^2 = \sigma^2 - \frac{\sigma^2}{n} + \mu^2 - \mu^2 = \\ &= \sigma^2 - \frac{\sigma^2}{n} = \frac{n\sigma^2 - \sigma^2}{n} = \frac{(n-1)\sigma^2}{n} \neq \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n} \text{ BIASED ESTIMATOR!} \end{aligned}$$

Is S^2 an unbiased estimator?

Recall: Sample variance of $X_i \sim N(\mu, \sigma^2)$: $\delta^2 = \frac{1}{n-1} \cdot \sum_{i=1}^n (x_i - \bar{x})^2 \Rightarrow \sum_{i=1}^n (x_i - \bar{x})^2 = S^2(n-1)$

$\left\{ \begin{array}{l} \frac{(n-1)S^2}{\sigma^2} \sim \chi^2_{(n-1)} \\ X_i \sim \chi^2_{(\sigma^2)} \rightarrow E(X) = \sigma^2 \Rightarrow \text{true df.} \end{array} \right.$

$$E(\delta^2) = E\left(\frac{\sum(x_i - \bar{x})^2}{n-1}\right) = E\left(\frac{S^2 \cdot (n-1)}{(n-1)} \cdot \frac{\sigma^2}{\sigma^2}\right) = E\left(\frac{S^2 \cdot (n-1)}{(n-1)} \cdot \frac{\sigma^2}{\sigma^2}\right) =$$

$$\frac{\sigma^2}{(n-1)} \cdot E\left(\frac{S^2 \cdot (n-1)}{\sigma^2}\right) = \frac{\sigma^2}{(n-1)} \cdot (n-1) = \sigma^2 \quad \delta^2 \text{ is UNBIASED of } \sigma^2 \text{ ESTIMATOR !!}$$

$\rightarrow \left(\begin{array}{l} X_i \sim \chi^2_{(n-1)} \\ E(X_i) = (n-1) \end{array} \right)$

S^2 is ALWAYS an unbiased estimator of σ^2 , but not of σ .

METHOD OF MOMENTS

for $K = 1, 2, \dots$

k^{th}	THEORETICAL MOMENT	SAMPLE MOMENTS
About the origin	$E(X^K)$	$M_K = \frac{1}{n} \sum_{i=1}^n X_i^K$
About the mean	$E[(X-\mu)^K]$	$M_K = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^K$

• ABOUT THE ORIGIN:

- ① Equate theoretical moments with sample moments
- 1^{st} $E(X) = M_1 = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}$
 - 2^{nd} $E(X^2) = M_2 = \frac{1}{n} \sum_{i=1}^n X_i^2$
- ② solve for the parameters (MME)
- ③ Results = METHOD OF MOMENTS ESTIMATORS! The empirical distribution converges to the probability distribution so the corresponding moments should be about equal when $n \rightarrow \infty$

* FROM THE ORIGIN : THEORETICAL MOMENT SAMPLE MOMENT solve for

Ex: $X_i \sim \text{Benoilli}(\rho)$

• 1st

$$E(X_i) = \rho = \frac{1}{n} \sum_{i=1}^n x_i$$

DONE!

$$\hat{\rho}_{\text{MM}} = \frac{1}{n} \sum_{i=1}^n x_i$$

Ex: $X_i \sim N(\mu, \sigma^2)$

• 1st

$$E(X_i) = \mu = \frac{1}{n} \sum_{i=1}^n x_i \Rightarrow \hat{\mu}_{\text{MM}} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$$

• 2nd

$$E(X_i^2) = \sigma^2 + \mu^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 \Rightarrow \hat{\sigma}_{\text{MM}}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \mu^2$$

$$\text{Var}(X) = \sigma^2 = E(X^2) - E(X)^2 = E(X^2) - \mu^2 \Rightarrow E(X^2) = \sigma^2 + \mu^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

The method of moments estimators = MLE estimators

$$\hat{\sigma}_{\text{MM}} = \hat{\sigma}; \hat{\mu}_{\text{MM}} = \hat{\mu}$$

* ABOUT THE MEAN

① Equate 1st theoretical moment about the origin with $M_1 = E(X) = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$
1st sample moment about the origin

② Equate 2nd theoretical moment about the mean with $M_2 = E(X - \mu)^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$
2nd sample moment about the mean

③ $K = 3, 4, \dots$

$$M_K = E(X - \mu)^K = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^K$$

④ Solve for parameters

⑤ Result = METHOD OF MOMENTS ESTIMATORS

Ex: $X_i \sim \text{Exp}(\alpha, \theta)$.

pdf of Exp dist: $f(x_i) = \frac{1}{\Gamma(\alpha)} \theta^\alpha \cdot x^{\alpha-1} \cdot e^{-x/\theta}$ for $x > 0$

Likelihood function: $L(\alpha, \theta) = \left(\frac{1}{\Gamma(\alpha)} \theta^\alpha \right)^n \cdot (x_1 x_2 \dots x_n)^{\alpha-1} \cdot e^{-\frac{1}{\theta} \sum x_i}$

$L'(\alpha, \theta)$ is hard to do bc of $\Gamma(\alpha)$ so MLE is hard.

$$\hat{\alpha}_{\text{MM}} = ? \quad \hat{\theta}_{\text{MM}} = ?$$

THEORETICAL	SAMPLE	SOLVE FOR PARAMETERS
$\cdot 1^{\text{st}} (\text{mean}) \quad E(X_i) = \alpha \theta = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$	$\alpha = \frac{\bar{x}}{\theta}$	$\text{Var}(x_i) = \frac{\bar{x} \cdot \theta^2 - \bar{x}^2}{\theta} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$
$\cdot 2^{\text{nd}} (\text{mean}) \quad \text{Var}(x_i) = E((X_i - \mu)^2) = \alpha \theta^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$	$\hat{\theta}_{\text{MM}} = \frac{1}{\bar{x}} \cdot \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$	
	$\hat{\alpha}_{\text{MM}} = \frac{\bar{x}}{\theta} = \frac{\bar{x}}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} = \frac{n \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2}$	

ASYMPTOTIC DISTRIBUTIONS OF MAXIMUM LIKELIHOOD ESTIMATORS :

continuous distribution \Rightarrow pdf $f(x; \theta)$.

we want to find $\hat{\theta}$ by solving $\frac{\partial(\ln L(\theta))}{\partial \theta} = 0$

If the approximation $\hat{\theta}$ is close enough to θ , by a series of proofs:

- CLT
- Taylor series.

$$\hat{\theta} \sim N \left(\theta, \frac{1}{-n E \left[\frac{\partial^2 \ln f(x; \theta)}{\partial \theta^2} \right]} \right)$$

This can be used for continuous and discrete dist

Ex: continuous: x_1, x_2, \dots, x_n random sample from exponential dist.

$$x_i \sim \text{Exp}(\theta). \quad \text{pdf: } f(x; \theta) = \frac{1}{\theta} e^{-\frac{x}{\theta}} \quad 0 < x < \infty \quad \theta \in \Omega (0 < \theta < \infty)$$

$$L(\theta) = L(\theta; x_1, x_2, \dots, x_n) = \left(\frac{1}{\theta} e^{-\frac{x_1}{\theta}} \right) \cdot \left(\frac{1}{\theta} e^{-\frac{x_2}{\theta}} \right) \cdot \left(\frac{1}{\theta} e^{-\frac{x_n}{\theta}} \right) = \frac{1}{\theta^n} \cdot e^{-\frac{\sum x_i}{\theta}} \quad 0 < \theta < \infty$$

$$\ln(L(\theta)) = -n \cdot \ln(\theta) - \frac{1}{\theta} \cdot \sum x_i$$

$$\frac{\partial \ln(L(\theta))}{\partial \theta} = \frac{-n}{\theta} + (-1) \cdot \frac{1}{\theta^2} \cdot \sum x_i = -\frac{n}{\theta} + \frac{\sum x_i}{\theta^2} = \frac{-\theta n + \sum x_i}{\theta^2} = 0$$

$$\theta = -\frac{\sum x_i}{n} = \frac{\sum x_i}{n} = \bar{x} \text{ is the MLE.}$$

$$\ln f(x; \theta) = \ln \left(\frac{1}{\theta} e^{-\frac{x}{\theta}} \right) = (-1) \cdot \ln(\theta) \quad \frac{-x}{\theta} \cdot \ln(e) = -\ln(\theta) - \frac{x}{\theta}$$

$$\frac{\partial(\ln f(x; \theta))}{\partial \theta} = -\frac{1}{\theta} - x \cdot (-1) \cdot \frac{1}{\theta^2} = -\frac{1}{\theta} + \frac{x}{\theta^2}$$

$$\frac{\partial^2 (\ln f(x; \theta))}{\partial \theta} = (-1) \cdot (-1) \cdot \frac{1}{\theta^2} + x \cdot (-2) \cdot \frac{1}{\theta^3} = \frac{1}{\theta^2} - \frac{2x}{\theta^3}$$

$$-E\left(\frac{\partial^2 (\ln f(x; \theta))}{\partial \theta}\right) = E\left(\frac{1}{\theta^2} - \frac{2x}{\theta^3}\right) = -\frac{1}{\theta^2} + \frac{2}{\theta^3} \cdot E(x) = -\frac{1}{\theta^2} + \frac{2}{\theta^3} \cdot \bar{x} = \frac{1}{\theta^2}$$

$\bar{x} = \theta$

VARIANCE AND ALSO THE CRLB

$$\hat{\theta} \sim N\left(\theta, \frac{1}{-n E\left[\frac{\partial^2 \ln f(x; \theta)}{\partial \theta^2}\right]}\right) \Rightarrow \hat{\theta} \sim N\left(\theta, \underbrace{\frac{1}{n \cdot \frac{1}{\theta^2}}}_{\text{mean}} = \frac{\theta^2}{n}\right)$$

standard deviation $\approx \sqrt{\frac{\theta}{n}}$

The variance of \bar{x} = the CRLB (lower bound)

$\hookrightarrow \bar{x}$ is an unbiased minimum variance estimator (UMVUE)
 $\hat{\theta}$

The random interval $\bar{x} \pm 1.96 \cdot \frac{\theta}{\sqrt{n}}$ has a prob of 95% that it covers θ .

$\bar{x} \pm 1.96 \cdot \frac{\theta}{\sqrt{n}}$ is an approximate 95% CI for θ

Ex: Discrete: $X_i \sim \text{Poisson } (\lambda) \rightarrow \text{pmf: } f(x; \lambda) = \frac{\lambda^x e^{-\lambda}}{x!} \quad x=0, 1, 2, \dots \quad \lambda \in \mathbb{R}^+ \text{ (OK)} \quad$

MLE for $\hat{\lambda} = \bar{x}$

$$\ln f(x; \lambda) = x \cdot \ln(\lambda) - \lambda - \ln(x!)$$

$$\frac{\partial (\ln f(x; \lambda))}{\partial \lambda} = \frac{x}{\lambda} - 1 \quad ; \quad \frac{\partial^2 (\ln f(x; \lambda))}{\partial \lambda^2} = -\frac{x}{\lambda^2}$$

$$-E\left(-\frac{x}{\lambda^2}\right) = \frac{1}{\lambda^2} \cdot E(x) = \frac{1}{\lambda^2} \cdot \lambda = \frac{1}{\lambda}$$

$$\hat{\lambda} \sim N\left(\lambda, \frac{1}{-n E\left[\frac{\partial^2 \ln f(x; \lambda)}{\partial \lambda^2}\right]}\right) \Rightarrow \hat{\lambda} \sim N\left(\lambda, \underbrace{\frac{1}{n \cdot \frac{1}{\lambda}}}_{\text{mean}} = \frac{\lambda}{n}\right)$$

standard dev = $\sqrt{\frac{\lambda}{n}} = \sqrt{\bar{x}}$

VARIANCE AND ALSO THE CRLB (lower bound)

The variance of \bar{x} = the CRLB (lower bound)

$\hookrightarrow \bar{x}$ is an unbiased minimum variance estimator (UMVUE)

$\hat{\lambda}$

* see next page.

UNIFORMLY MINIMUM VARIANCE UNBIASED ESTIMATOR (UMVUE)

The variance of $\hat{\theta}$ serves as the lower bound for the variance of any unbiased estimator of $\theta \Rightarrow$ therefore, if we find an unbiased estimator that has a variance = to the lower bound of $\hat{\theta} \Rightarrow$ best estimator or UMVUE.

because we cannot find a better estimator.
 so for $x_1, x_2 \dots x_n$ random sample of size n from dist $f(x; \theta)$
 the estimator U^* is UMVUE if $\left\{ \begin{array}{l} U^* \text{ is an UNBIASED estimator of } \theta \\ \text{for any other unbiased estimator of } \theta \text{ } U \\ \text{Var}(U^*) \leq \text{Var}(U) \text{ for all } \theta \in \Omega \end{array} \right.$

CRAIG - RAO LOWER BOUND (CRLB)

or INEQUALITY.

For $x_1, x_2 \dots x_n$ random sample of size n from continuous dist $f(x; \theta)$
 The support of x does not depend on θ . $\Omega \subset \mathbb{R} = \{\theta; c < \theta < d\}$

$y = u(x_1, x_2 \dots x_n)$ is an unbiased estimator of θ .

$$\text{Var}(y) \geq \frac{1}{n E \left[\frac{\partial \ln f(x; \theta)}{\partial \theta} \right]^2} = \frac{-1}{n E \left[\frac{\partial^2 \ln f(x; \theta)}{\partial \theta^2} \right]} = \frac{1}{I}$$

$$I = \text{FISHER'S INFORMATION} = n E \left[\frac{\partial \ln f(x; \theta)}{\partial \theta} \right]^2 = -n E \left[\frac{\partial^2 \ln f(x; \theta)}{\partial \theta^2} \right]$$

$$\text{for continuous dist: } n \int_{-\infty}^{\infty} \left[\frac{\partial \ln f(x; \theta)}{\partial \theta} \right]^2 f(x; \theta) dx \stackrel{\uparrow}{=} \int_{-\infty}^{\infty} \left[\frac{\partial^2 \ln f(x; \theta)}{\partial \theta^2} \right] f(x; \theta) dx$$

$$\text{for discrete dist: } n \sum \left[\frac{\partial \ln f(x; \theta)}{\partial \theta} \right]^2 f(x; \theta) \quad n \sum \left[\frac{\partial^2 \ln f(x; \theta)}{\partial \theta^2} \right] f(x; \theta)$$

EFFICIENCY of an estimator: ratio of the Rao-Cramer lower bound to the variance of any unbiased estimator.

e.g.: an estimator with efficiency of 50% means that $\frac{1}{0.5} = 2$ times as many samples observations are needed to do as well in an estimation as can be done with the MVUE (that is a 100% efficient estimator).

ex: pdf $x = f(x; \theta) = \theta x^{\theta-1} \quad 0 < x < 1 \quad \theta \in \Omega = \{ \theta : 0 < \theta < \infty \}$

$$\ln f(x; \theta) = \ln(\theta) + (\theta-1) \cdot \ln(x)$$

$$\frac{\partial(\ln f(x; \theta))}{\partial \theta} = \frac{1}{\theta} + \ln(x) \quad ; \quad \frac{\partial^2(\ln f(x; \theta))}{\partial \theta^2} = (-1) \cdot \frac{1}{\theta^2} = -\frac{1}{\theta^2}$$

$$E\left(\frac{\partial^2(\ln f(x; \theta))}{\partial \theta^2}\right) = E\left(-\frac{1}{\theta^2}\right) = -\frac{1}{\theta^2}$$

$$\text{Var}(y) \geq \frac{-1}{n E\left[\frac{\partial^2 \ln f(x; \theta)}{\partial \theta^2}\right]} = -\frac{1}{n \cdot \left(-\frac{1}{\theta^2}\right)} = \frac{\theta^2}{n} \quad \text{is the greatest lower bound of the variance of every unbiased estimator of } \theta$$

MLE of $\theta = \hat{\theta}$

$$\mathcal{L}(\theta) = \theta x_1^{\theta-1} \cdot \theta x_2^{\theta-1} \cdots \theta x_n^{\theta-1} = n \theta \left(\prod_{i=1}^n x_i \right)^{\theta-1}$$

$$\ln(\mathcal{L}(\theta)) = n \ln(\theta) + (\theta-1) \cdot \ln\left(\prod_{i=1}^n x_i\right)$$

$$\frac{\partial \ln(\mathcal{L}(\theta))}{\partial \theta} = \frac{n}{\theta} + \ln\left(\prod_{i=1}^n x_i\right) \equiv 0 \Rightarrow \frac{n + \ln\left(\prod_{i=1}^n x_i\right) \cdot \theta}{\theta} = 0$$

$\hat{\theta} = \frac{n}{\ln\left(\prod_{i=1}^n x_i\right)}$ is the MLE

$$\hat{\theta} \sim N\left(\theta, \frac{\theta^2}{n}\right)$$

$\underbrace{\frac{\theta^2}{n}}$ = variance of $\hat{\theta}$

so $\hat{\theta}$ is the MVUE.

Properties of MLE (Summary) Under certain regularity conditions, the MLE

1. exists and is unique
2. is a consistent (\approx asymptotically unbiased) estimator
3. is asymptotically normal
4. is asymptotically efficient.

$$\hat{\theta}_{MLE} \sim N\left(\theta, -\frac{1}{n E\left[\frac{\partial^2 \ln f}{\partial \theta^2}\right]}\right)$$

Problems: you need to know the pdf. ($f(x; \theta)$)

Definition (Mean Squared Error) If Y is an estimator of θ , then the bias is given by

$$\text{Bias}(Y) = E(Y) - \theta$$

and the mean squared error (MSE) of Y is given by

$$\text{MSE}(Y) = E[(Y - \theta)^2].$$

Theorem If Y is an estimator of θ , then

$$\text{MSE}(Y) = \text{Var}(Y) + [\text{Bias}(Y)]^2.$$

DELEM METHOD.

Estimator \hat{Y}_n of a parameter $\gamma \rightarrow \hat{Y}_n \sim N\left(\gamma, \frac{\sigma_\gamma^2}{n}\right)$

$$E(\hat{Y}_n) \approx \gamma \quad \text{Var}(\hat{Y}_n) \approx \frac{\sigma_\gamma^2}{n}$$

Estimator $f(\hat{Y}_n)$ is an estimator of $f(\gamma)$

↳ Taylor series expansion of $f(\hat{Y}_n)$ about $f(\gamma) \Rightarrow$

$$f(\hat{Y}_n) = f(\gamma) + f'(\gamma)(\hat{Y}_n - \gamma) + (\text{small remainder term})$$

$$E(f(\hat{Y}_n)) = E\left[f(\gamma) + f'(\gamma)(\hat{Y}_n - \gamma)\right] = f(\gamma) + f'(\gamma)E(\hat{Y}_n - \gamma)$$

$$\text{Var}(f(\hat{Y}_n)) = \text{Var}[f(\gamma) + f'(\gamma)(\hat{Y}_n - \gamma)] = f'(\gamma)^2 \cdot \text{Var}(\hat{Y}_n) = f'(\gamma)^2 \left(\frac{\sigma_\gamma^2}{n}\right)$$

$$\text{Var}(a+bX) = \text{Var}(a) + \text{Var}(bX) + 2\text{Cov}(a, bX)$$

$$\underbrace{\text{Var}(f(\gamma))}_{0} + \underbrace{\text{Var}(f'(\gamma)(\hat{Y}_n - \gamma))}_{0} + \underbrace{2\text{Cov}(a, bX)}_{0}$$

$$f(\hat{Y}_n) \sim N\left(f(\gamma), f(\gamma)^2 \left(\frac{\sigma_\gamma^2}{n}\right)\right)$$

DELTA METHOD :

standard error for odds / risk ratio. You have to take the log of the odds.

Approximation by Taylor series

$$\hat{\theta} \sim N(\theta, \frac{\sigma_\theta^2}{n})$$

$$f(\theta) \sim N(f(\theta), f'(\theta) \text{SE}_\theta) = (f(\theta), f'(\theta)^2 \text{Var}(\theta))$$

TAYLOR SERIES: $\frac{f^k(\theta)(x-\theta)^k}{k!}$ you can use it to approximate $f(x) = y$

$$f(x) \approx f(\theta) + f'(\theta)(x-\theta) + \frac{f''(\theta)(x-\theta)^2}{2!} \text{ residual term is even smaller}$$

$$E(y) = E[f(x)] + E[f'(\theta)(x-\theta)] = f(x) + \underbrace{f'(\theta)(\theta - \theta)}_0$$

$$E(y) = f(x)$$

$$\text{Var}(x) \approx \text{Var}[f(\theta) + f'(\theta)(x-\theta)] = \text{Var}[f'(\theta)(x-\theta)] = \\ f'(\theta)^2 \text{Var}x - \text{Var}(f'(\theta)) = f'(\theta)^2 \cdot \text{Var}(x)$$

ANOTHER METHOD:

$$\text{If } \hat{\theta} \rightarrow \theta \quad \frac{dx}{\hat{\theta}} \quad \frac{f(\hat{\theta}) - f(\theta)}{\hat{\theta} - \theta} \approx f'(\hat{\theta}) \Rightarrow \frac{f(\hat{\theta}) - f(\theta)}{f'(\theta)} \approx \hat{\theta} - \theta \left(\frac{\text{SE}_\theta}{\text{SE}_\theta} \right)$$

$$\Rightarrow \frac{f(\hat{\theta}) - f(\theta)}{f'(\theta) \text{SE}_\theta} \approx \frac{\hat{\theta} - \theta}{\text{SE}_\theta}$$

CONFIDENCE INTERVALS FOR ONE MEAN

Problem?

- When we use the sample $\{\text{mean } \bar{x}$
 $\text{proportion } \hat{p}\}$
- To estimate the population $\{\text{mean } \mu$
 $\text{proportion } p\}$

C.I. =	$\begin{cases} 1 \text{ sample: } ?(Y=1) \\ \text{of the mean: } 2 \text{ samples, } z, t, \text{ bin } \end{cases}$
• of the variance	$\begin{cases} 1 \text{ sample: } \chi^2 \\ 2 \text{ samples: } F \end{cases}$

Can we be confident that \bar{x} is close to μ ?

- And how close is the sample statistic to the pop. parameter?

Instead of using a point estimate, we can use an interval of values that we are confident contains the actual unknown pop. parameter.

CONFIDENCE INTERVAL: the interval of lower (L) and upper (U) values between which we can be really confident the population mean falls. $L < \mu < U$
 $L < p < U$

CONFIDENCE COEFFICIENT: reported as a proportion: $1 - \alpha : 0.9, 0.95, 0.99$

CONFIDENCE LEVEL: reported as a percentage: $(1 - \alpha) \cdot 100 : 90\%, 95\%, 99\%$

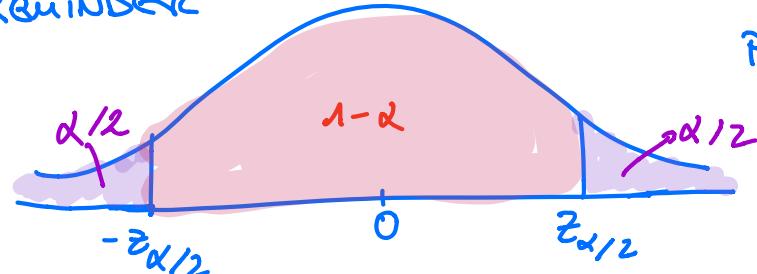
"we can be 95% confident that the population falls between L and U .

CONFIDENCE LIMITS: endpoints a and b of a CI

CRITICAL VALUES: Cut-off values. ex: $z_{\alpha/2}$

CONFIDENCE PROBABILITY (CP): CP of the random interval $[L, U]$ is the probability that $[L, U]$ covers the true θ and is denoted by $P(\theta \in [L, U])$

REMINDER



$$P(z > z_{\alpha/2}) = \alpha/2$$

$$P(-z < z_{\alpha/2}) = 1 - P(z > z_{\alpha/2}) = \alpha/2$$

CONFIDENCE INTERVAL FOR POPULATION MEAN.

- ASSUME σ IS KNOWN:

- Generic confidence interval for the mean: $P(a < \frac{\bar{X}-\mu}{\sigma/\sqrt{n}} < b) = 1-\alpha$

- Z-interval for a mean:

$$P(-z_{\alpha/2} < \frac{\bar{X}-\mu}{\sigma/\sqrt{n}} < z_{\alpha/2}) = 1-\alpha$$

Assume:

- X_1, X_2, \dots, X_n $X_i \sim N(\mu, \sigma^2)$ and $Z = \frac{\bar{X}-\mu}{\sigma/\sqrt{n}} \sim N(0,1)$

- The population variance σ^2 is known.

$(1-\alpha) \cdot 100\%$ confidence interval for the mean μ is $\bar{x} \pm z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right)$

Because μ is the parameter of interest, we isolate μ in the middle.

$$\begin{aligned} P(-z_{\alpha/2} < \frac{\bar{X}-\mu}{\sigma/\sqrt{n}} < z_{\alpha/2}) &= P(-z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} - \bar{x} < -\mu < z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} - \bar{x}) \\ &= P(\bar{x} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}) = 1-\alpha. \end{aligned}$$

A common choice of $\alpha = 0.05 \rightarrow z_{\alpha/2} = 1.96$.

We can be $(1-\alpha) \cdot 100 = (1-0.05) \cdot 100 = 95\%$ confident that the μ is in the interval $\left[\bar{x} - 1.96 \frac{\sigma}{\sqrt{n}}, \bar{x} + 1.96 \frac{\sigma}{\sqrt{n}} \right]$

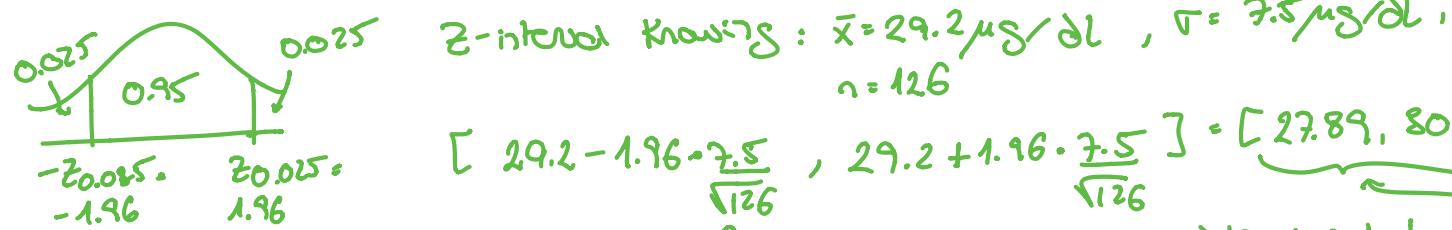
Ex: sample $n=126$ police officers whose average $29.2 \mu\text{g/dL}$ blood lead level

X = blood lead level. $X \sim N(\mu, \sigma^2 = 7.5 \mu\text{g/dL})$

Historically avg blood level of humans w/o no exposure = $18.2 \mu\text{g/dL}$.

Do the police officers have elevated blood concentrations?

95% CI $\rightarrow 1-\alpha = 0.95 \rightarrow \alpha = 0.05 \rightarrow \alpha/2 \rightarrow 0.025 \rightarrow z_{0.025} = 1.96$.



$$\left[29.2 - 1.96 \cdot \frac{7.5}{\sqrt{126}}, 29.2 + 1.96 \cdot \frac{7.5}{\sqrt{126}} \right] = \left[\underbrace{27.89}_{\text{lower bound}}, \underbrace{30.51}_{\text{upper bound}} \right]$$

We are 95% confident that the mean blood level is between $[27.89, 30.51]$ because $[27.89, 30.51]$ does not include the historical avg $18.2 \mu\text{g/dL} \rightarrow$ police officer is elevated.

INTERPRETATION:

correct: we are $(1-\alpha)\%$ confident that the interval includes the true pop mean μ

Ex: if we took 1000 samples \rightarrow calculated 95% CI for each \rightarrow

 950 of the 1000 intervals would actually contain the unknown value μ .
 if we take only 1 sample, the interval will either contain or not the μ .
 so we can only say we are 95% confident the interval contains μ .

incorrect: the probability that the pop mean μ falls between L and U is $1-\alpha$

because: the pop. μ is a constant, not a random variable.
 you can't make a probability statement about a constant that

doesn't change, only about random variables.

Ex: so if you say $P[30 < \mu < 90] = 95\%$ that's not correct. if $\begin{cases} \mu = 40 \rightarrow P[30 < \mu < 90] = 1 \\ \mu = 100 \rightarrow P[30 < \mu < 90] = 0 \end{cases}$

LENTH OF AN INTERVAL: we want the interval to be as narrow as possible.

if the interval is $L < \bar{x} < U \rightarrow$ Length $U-L$.

$$\text{for the Z-interval: Length } [\bar{x} + z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right)] - [\bar{x} - z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right)] = 2 \cdot z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right)$$

Length of the interval is influenced by:

- standard deviation σ : as $\sigma \downarrow \rightarrow$ length \downarrow
 - sample size n : as $n \uparrow \rightarrow$ length \downarrow
 - confidence level (α): as confidence level $\downarrow \rightarrow$ length \downarrow
- Ex.: 95% $\rightarrow z = 1.96$. 90% $\rightarrow z = 1.645$

- ASSUME σ IS UNKNOWN.

We estimate the population st.dev. σ from the sample standard deviation s .

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

For $x_1, x_2, \dots, x_n \rightarrow x_i \sim N(\mu, \sigma^2)$ then a $(1-\alpha) \cdot 100\%$ confidence interval for μ is
 - t-interval for the mean: $\bar{x} \pm t_{\alpha/2, n-1} \left(\frac{s}{\sqrt{n}} \right)$

To derive the confidence interval we use $\frac{\bar{x}-\mu}{s/\sqrt{n}}$ instead of $\frac{\bar{x}-\mu}{\sigma/\sqrt{n}} \sim N(0, 1)$

For $x_1, x_2, \dots, x_n \rightarrow x_i \sim N(\mu, \sigma^2)$ then $T = \frac{\bar{x}-\mu}{s/\sqrt{n}}$ follows t dist w/ $(n-1)$ df

PROOF:

a T random variable $T = \frac{\bar{Z}}{\sqrt{U/n}}$ if independent $\Rightarrow T \sim t(r)$

For $X_1, X_2, \dots, X_n \Rightarrow X_i \sim N(\mu, \sigma^2) \Rightarrow \left. \begin{array}{l} Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \\ \frac{(n-1)\sigma^2}{\sigma^2} \sim \chi^2_{(n-1)} \end{array} \right\} Z \text{ and } S^2 \text{ are independent}$

$$T = \frac{\bar{Z}}{\sqrt{U/n}} = \frac{\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}}{\sqrt{\frac{(n-1)\sigma^2}{\sigma^2}/(n-1)}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \cdot \frac{1}{\sqrt{\frac{(n-1)}{n}}} = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$$

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$$

$$P[-t_{\alpha/2, n-1} < T = \frac{\bar{X} - \mu}{S/\sqrt{n}} < t_{\alpha/2, n-1}] = 1 - \alpha$$

$$\bar{X} - t_{\alpha/2, n-1} \cdot \frac{S}{\sqrt{n}} < \mu < \bar{X} + t_{\alpha/2, n-1} \cdot \frac{S}{\sqrt{n}}$$

- POINT ESTIMATE: \bar{x} is a point estimate of μ
- INTERVAL ESTIMATE: $\bar{x} \pm t_{\alpha/2, n-1} \left(\frac{S}{\sqrt{n}} \right)$ → MARGIN OF ERROR: $t_{\alpha/2, n-1} \left(\frac{S}{\sqrt{n}} \right)$
- STANDARD ERROR OF THE MEAN: $\frac{S}{\sqrt{n}}$

ONE SIDED CONFIDENCE INTERVAL:

Replace $\frac{z}{t}_{\alpha/2}$ with $\frac{z}{t}_{\alpha}$ $[\bar{x} - \frac{z_{\alpha} \frac{S}{\sqrt{n}}}{t}, \infty)$ or $(-\infty, \bar{x} + \frac{z_{\alpha} \frac{S}{\sqrt{n}}}{t}]$

for exercises, 1st you have to make sure your data is normal
2nd find your t value: ex $n=16$ 75% CI $\Rightarrow t_{0.025, 15} = 2.1314$
3rd Apply formula.

NON-NORMALLY DISTRIBUTED DATA

- with sample > 30 (although possibly more depending on skewness of dist.) $\bar{x} \pm t_{\alpha/2, n-1} \left(\frac{S}{\sqrt{n}} \right)$ and $\bar{x} \pm z_{\alpha/2} \left(\frac{S}{\sqrt{n}} \right)$ give similar results.
- with sample < 30 :
 - use $\bar{x} \pm t_{\alpha/2, n-1} \left(\frac{S}{\sqrt{n}} \right)$ with caution or
 - use non-parametric CI

BAYESIAN INTERPRETATION of the interval estimator:

	FREQUENTIST	BAYESIAN
95% coverage	in 95% of repeated experiments the interval will cover the true Θ	the probability that the Θ is in the interval is 95%
Randomness	repetition of experiments	uncertainty about the value of the Θ "CREDIBLE INTERVAL"

CONFIDENCE INTERVALS FOR 2 MEANS

KNOWN VARIANCES:

$x_i = x_1, x_2, \dots, x_n \sim N(\mu_x, \sigma_x^2)$. $y_i = y_1, y_2, \dots, y_m \sim N(\mu_y, \sigma_y^2)$

σ_x^2 and σ_y^2 are known. we want to compare μ_x and μ_y

$$\bar{x} - \bar{y} \text{ is } N\left(\mu_x - \mu_y, \frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}\right) \rightarrow P\left(\bar{x} - \bar{y} - z_{\alpha/2} \sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}} < \mu_x - \mu_y < \bar{x} - \bar{y} + z_{\alpha/2} \sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}}\right)$$

$$P\left(-z_{\alpha/2} < \frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_y)}{\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}}} < z_{\alpha/2}\right) = 1 - \alpha$$

UNKNOWN BUT EQUAL VARIANCES:

$x_i = x_1, x_2, \dots, x_n \sim N(\mu_x, \sigma_x^2)$. $y_i = y_1, y_2, \dots, y_m \sim N(\mu_y, \sigma_y^2)$

x_i and y_i are independent random variables.

(1- α) 100% confidence interval for $\mu_x - \mu_y$ the difference is the population mean if:

$$(\bar{x} - \bar{y}) \pm t_{\alpha/2, (n+m-2)} s_p \sqrt{\frac{1}{n} + \frac{1}{m}}$$

where s_p^2 = POOLED SAMPLE VARIANCE is an unbiased estimator of the common variance σ^2 .
it's an average of the sample variances weighted by the sample sizes.

$$s_p^2 = \frac{(n-1)s_x^2 + (m-1)s_y^2}{n+m-2}$$

ASSUMPTIONS:

- x_i and y_i are independent
- the measurements of each pop are normally distributed.
- the measurements of each pop have the same variance σ^2
- the measurements of each pop have the same variance σ^2

INTERPRETATION:
if the difference of the sample means is 0, we cannot conclude that the pop. means differ.

PROOF:

$\bar{x} \sim N\left(\mu_x, \frac{\sigma^2}{n}\right)$; $\bar{y} \sim N\left(\mu_y, \frac{\sigma^2}{m}\right) \Rightarrow \bar{x} - \bar{y} \sim N\left(\mu_x - \mu_y, \frac{\sigma^2}{n} + \frac{\sigma^2}{m}\right)$ because x, y independent

$$z = \frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_y)}{\sqrt{\frac{\sigma^2}{n} + \frac{\sigma^2}{m}}} \sim N(0, 1)$$

bc x, y are normal

$$U = \frac{(n-1)s_x^2}{\sigma^2} + \frac{(m-1)s_y^2}{\sigma^2} \sim \chi^2_{n+m-2}$$

$$T = \frac{z}{\sqrt{\frac{U}{n+m-2}}} = \frac{1}{\sigma} \frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_y)}{\sqrt{\frac{1}{n} + \frac{1}{m}}} = \frac{\bar{x} - \bar{y} - (\mu_x - \mu_y)}{s_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{n+m-2}$$

$$\frac{1}{\sigma} \sqrt{\frac{(n-1)s_x^2 + (m-1)s_y^2}{(n+m-2)}} = s_p$$

• UNKNOWN AND UNEQUAL VARIANCES:

- LARGE SAMPLE SIZE

$$(\bar{x} - \bar{y}) \pm z_{\alpha/2} \cdot \sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}$$

PROOF:

$$\frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_y)}{\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}} \rightarrow z \sim N(0, 1)$$

ASSUMPTIONS:

x and y are normally dist or
 n and $m \rightarrow \infty$

- SMALL SAMPLE SIZE

WELCH'S t-interval: $\bar{x} - \bar{y} \pm t_{\alpha/2, r} \cdot \sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}$

uses a modified t-test w/ df
Protects against smaller
sample size associated w/ larger
variance.

$$df = \frac{\left(\frac{s_x^2}{n} + \frac{s_y^2}{m}\right)^2}{\frac{\left(\frac{s_x^2}{n}\right)^2}{n-1} + \frac{\left(\frac{s_y^2}{m}\right)^2}{m-1}}$$

use the integer $[r]$

PROOF:

$$\frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_y)}{\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}} \sim t_{(r)} \text{ with small sample sizes}$$

WHEN TO POOL VARIANCE?

$\frac{s_x^2}{s_y^2} > 4$ or $\frac{s_y^2}{s_x^2} > 4 \rightarrow$ use Welch's t-interval.
Otherwise use 2-sample t-interval.

Ex: $X \sim N(0, 1)$ $Y \sim N(0, 36)$ with Student's t-test superimposed

$$n=6 \quad m=18 \quad n < m \\ s_x^2=1 \quad s_y^2=36 \quad \sigma_x^2 < \sigma_y^2$$

$$n=6 \quad m=18 \quad n > m \\ s_x^2=1 \quad s_y^2=36 \quad \sigma_x^2 < \sigma_y^2$$

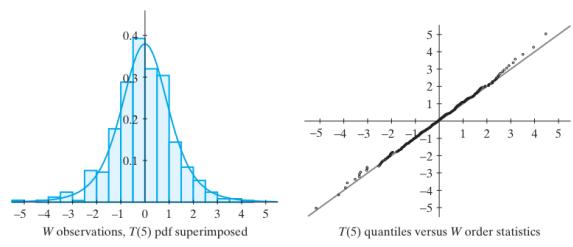
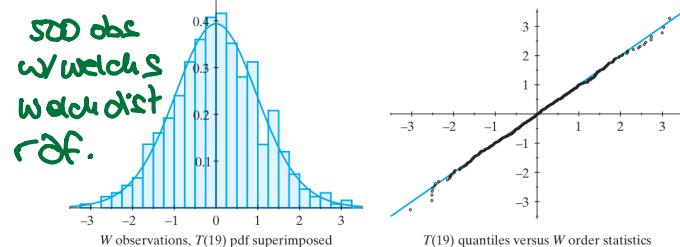
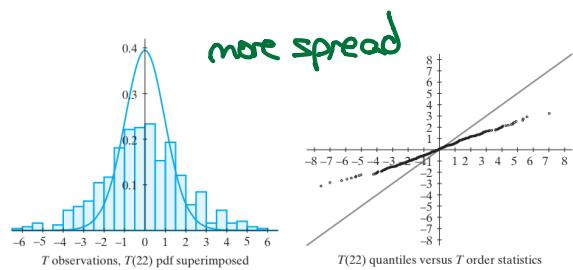
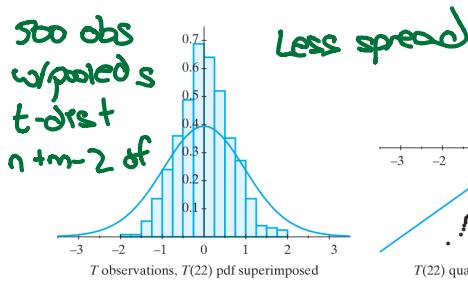


Figure 7.2-1 Observations of T and of W , $n = 6, m = 18, \sigma_x^2 = 1, \sigma_y^2 = 36$

Figure 7.2-2 Observations of T and of W , $n = 18, m = 6, \sigma_x^2 = 1, \sigma_y^2 = 36$

- PAIRED SAMPLES :

we cannot use t-statistics bc x and y are not independent

$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ pairs of measurements (before, after etc)

$D_i = x_i - y_i$ $D_i \sim N(\mu_D, \sigma_D^2)$ is a random sample of pairs.

$(1-\alpha) \cdot 10\%$ CI of the sample mean difference to estimate the pop mean difference is

$$\bar{d} \pm t_{\alpha/2, n-1} \left(\frac{s_d}{\sqrt{n}} \right)$$

Subject	Red (x)	Green (y)	$d = x - y$	$\rightarrow \bar{d} = -0.0625$
1	0.30	0.43	-0.13	$s_d = 0.0765$
2	0.23	0.32	-0.09	
3	0.41	0.58	-0.17	
4	0.53	0.46	0.07	
5	0.24	0.27	-0.03	
6	0.36	0.41	-0.05	
7	0.38	0.38	0.00	
8	0.51	0.61	-0.10	

CONFIDENCE INTERVALS OF VARIANCES

ONE VARIANCE:

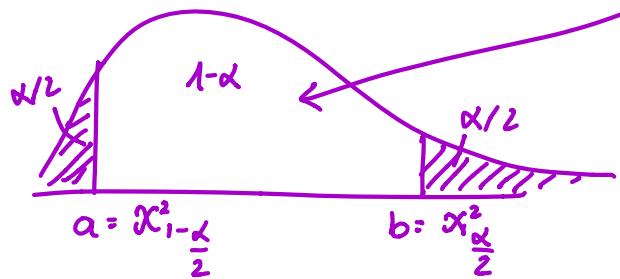
for $X_i = x_1, x_2, \dots, x_n$ $X_i \sim N(\mu, \sigma^2)$ with $a = \chi^2_{1-\alpha/2, n-1}$ $b = \chi^2_{\alpha/2, n-1}$

$(1-\alpha)\%$ confidence interval for

- the pop variance σ^2 $\left(\frac{(n-1)s^2}{b} \leq \sigma^2 \leq \frac{(n-1)s^2}{a} \right)$
- the standard deviation σ $\left(\sqrt{\frac{(n-1)s^2}{b}} \leq \sigma \leq \sqrt{\frac{(n-1)s^2}{a}} \right)$

Proof

We know that if $X_i \sim N(\mu, \sigma^2)$ $\frac{(n-1)s^2}{\sigma^2} \sim \chi^2_{n-1}$



$$\frac{(n-1)s^2}{\sigma^2} \sim \chi^2_{n-1}$$

$$P\left[a \leq \frac{(n-1)s^2}{\sigma^2} \leq b\right] = 1-\alpha$$

$$\frac{1}{a} > \frac{\sigma^2}{(n-1)s^2} > \frac{1}{b} \Rightarrow \frac{(n-1)s^2}{b} \leq \sigma^2 \leq \frac{(n-1)s^2}{a} \Rightarrow \sqrt{\frac{(n-1)s^2}{b}} \leq \sigma \leq \sqrt{\frac{(n-1)s^2}{a}}$$

Ex: Packs of candy with target weight 52g. Estimate σ .
sample of candy $n=10$, sample variance $s=4.2$. 95% CI?

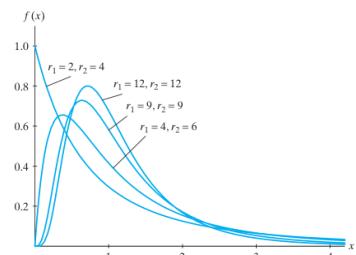
$$a = \chi^2_{1-\alpha/2, n-1} = \chi_{0.975, 9} = 2.7 \quad b = \chi^2_{\alpha/2, n-1} = \chi_{0.025, 9} = 19.02$$

$$\left(\frac{\sqrt{9 \cdot 4.2}}{\sqrt{19.02}} < \sigma < \frac{\sqrt{9 \cdot 4.2}}{\sqrt{2.7}} \right) \Rightarrow (1.41 < \sigma < 3.74)$$

F-DISTRIBUTION:

$U \sim \chi^2_{(r_1)}$ $V \sim \chi^2_{(r_2)}$. U and V independent. $F = \frac{U}{V} \sim F(r_1, r_2)$

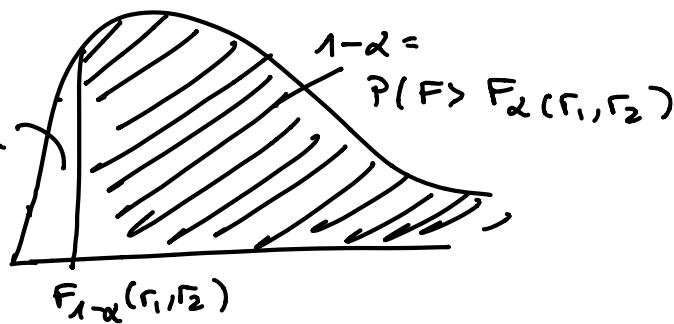
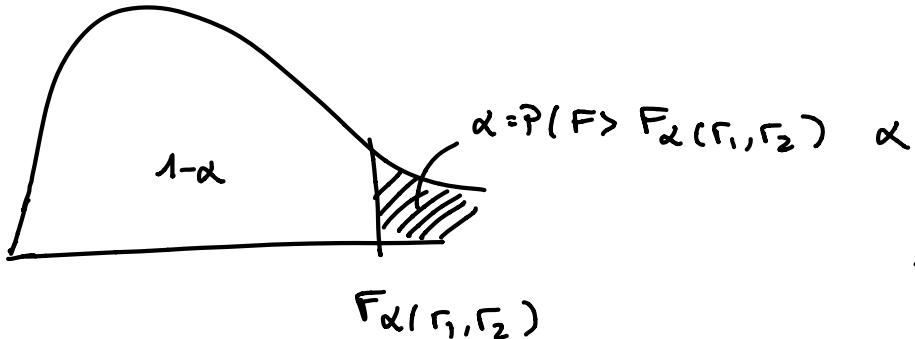
$$P[W \geq F_\alpha(r_1, r_2)] = \alpha.$$



$$\text{if } \omega \sim F(r_1, r_2) \Rightarrow \frac{1}{\omega} \sim F(r_2, r_1)$$

$$\alpha = P[\omega \leq F_{1-\alpha}(r_1, r_2)] = P\left[\frac{1}{\omega} \geq \frac{1}{F_{1-\alpha}(r_1, r_2)}\right]$$

$$\frac{1}{F_{1-\alpha}(r_1, r_2)} = F_\alpha(r_2, r_1) \text{ or } F_{1-\alpha}(r_1, r_2) = \frac{1}{F_\alpha(r_2, r_1)}$$



Upper $100 \cdot \alpha^{\text{th}}$ percentile of a F dist is $F_{\alpha}(r_1, r_2)$

RATIO OF 2 VARIANCES :

x_i and y_i are independent, normally distributed random samples

$$x_i = x_1, x_2, \dots, x_n \sim N(\mu_x, \sigma_x^2) \quad y_i = y_1, y_2, \dots, y_m \sim N(\mu_y, \sigma_y^2)$$

$$c = F_{1-\alpha/2}(m-1, n-1) = \frac{1}{F_{\alpha/2}(n-1, m-1)} \quad ; \quad d = F_{\alpha/2}(m-1, n-1)$$

The $(1-\alpha) \cdot 100\%$ CI of $\frac{s_x^2}{s_y^2}$ is $\frac{1}{F_{\alpha/2}(n-1, m-1)} \cdot \frac{s_y^2}{s_x^2} \leq \frac{s_x^2}{s_y^2} \leq F_{\alpha/2}(m-1, n-1) \cdot \frac{s_x^2}{s_y^2}$

PROOF:

$$x_i \sim N(\mu_x, \sigma_x^2) \Rightarrow \frac{(n-1)s_x^2}{\sigma^2} \sim \chi_{n-1}^2 \quad y_i \sim N(\mu_y, \sigma_y^2) \Rightarrow \frac{(m-1)s_y^2}{\sigma^2} \sim \chi_{m-1}^2$$

$$F = \frac{\frac{(m-1)s_y^2}{\sigma^2}}{\frac{(n-1)s_x^2}{\sigma^2}} = \frac{\frac{s_y^2}{\sigma^2}}{\frac{s_x^2}{\sigma^2}} \sim F(m-1, n-1)$$

$$\mathbb{P}\left[F_{1-\frac{\alpha}{2}}(m-1, n-1) \leq \frac{s_y^2}{s_x^2} \leq F_{\frac{\alpha}{2}}(m-1, n-1)\right] = 1-\alpha$$

$$\mathbb{P}\left[\frac{1}{F_{\frac{\alpha}{2}}(n-1, m-1)} \cdot \frac{s_x^2}{s_y^2} \leq \frac{s_y^2}{s_x^2} \leq F_{\frac{\alpha}{2}}(m-1, n-1)\right]$$

CIs for variances are "sensitive" to the normality assumption, are not "robust", are not accurate when data is not normally distributed. However it can be important when we are testing if variances are equal, like in - t-test, or checking the assumptions of ANOVA.

REMINDER

NORMAL APPROXIMATION TO BINOMIAL

$Y \sim b(n, p)$ is the sum of n $X_i \sim \text{Bernoulli}(p)$

$y = \sum_{i=1}^n X_i$ random variables, X_1, X_2, \dots, X_n

	mean	variance
$X_i \sim \text{Bernoulli}(p)$ with	$\mu = E[X] = p$	$\sigma^2 = \text{Var}(X) = p(1-p) = pq$
$Y \sim \text{binomial}(np)$ with	$\mu = np$	$\sigma^2 = np(1-p) = npq$

CLT : $Y \sim N(np, npq)$ if n is sufficiently large $\left\{ \begin{array}{l} np \geq 5 \\ npq \geq 5 \end{array} \right.$

$$W = \frac{\bar{X} - \mu}{\sqrt{\frac{\sigma^2}{n}}} = \frac{Y - np}{\sqrt{np(1-p)}} = \frac{\frac{Y - np}{\sqrt{np}}}{\sqrt{\frac{p(1-p)}{n}}} \xrightarrow{n \rightarrow \infty} N(0, 1)$$

$$W = \frac{\sum_{i=1}^n X_i - np}{\sqrt{n \cdot pq}} = \frac{Y - np}{\sqrt{n \cdot pq}} = \frac{Y - np}{\sqrt{npq}} \xrightarrow{n \rightarrow \infty} N(0, 1)$$

SAMPLE PROPORTION : $\hat{p} = \frac{\sum_{i=1}^n X_i}{n}$ is the proportion of the sample that meets condition of interest.

It can be estimated by :

$$Z = \frac{\sum_{i=1}^n X_i - np}{\sqrt{np(1-p)}} = \frac{\sum_{i=1}^n X_i - np}{\sqrt{\frac{np(1-p)}{n}}} = \frac{\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}}}{\sqrt{\frac{p(1-p)}{n}}} \xrightarrow{n \rightarrow \infty} N(0, 1)$$

CONFIDENCE INTERVALS FOR PROPORTIONS

ONE PROPOSITION:

For large random sample $\rightarrow (1-\alpha) \cdot 100\%$ confidence interval for a population proportion p is:

$$\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

ONE SIDED: $[0, \hat{p} + z_{\alpha} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}]$ or $[\hat{p} - z_{\alpha} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, 1]$

CONDITIONS:

- $np = \text{no of expected successes} \geq 5$

AND

- $n(1-p) = \text{no of expected failures} \geq 5$

PROOF:

CLT \Rightarrow for a large n $Z = \frac{\bar{X} - \mu}{\sigma}$ follows $\mathcal{N}(0, 1)$

$$\mu = E(X_i) = p$$

$$\sigma^2 = \text{Var}(X_i) = p(1-p)$$

$$\bar{X} = \hat{p}$$

$$\text{for a large } n \quad Z = \frac{\bar{X} - p}{\sqrt{\frac{p(1-p)}{n}}} \Rightarrow P\left[-z_{\alpha/2} \leq \frac{\bar{X} - p}{\sqrt{\frac{p(1-p)}{n}}} \leq z_{\alpha/2}\right] \approx 1-\alpha$$

$$P\left[-\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right] =$$

$$P\left[\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right]$$

we need to know the population proportion p to estimate population proportion P !!
 Instead of population proportion p we will use sample proportion \hat{p} to get the approximate.

$$\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

Ex: $n=418$ Pennsylvanians surveyed on opinion. $4=280$ blaue doctors.

$$\hat{p} = \frac{280}{418} = 0.67 \text{ is the sample proportion}$$

Estimate proportion of all pennsylvanians who blaue doctors with a 95% CI

$x_i = 1 \rightarrow$ randomly selected Pennsylvania does home doctors

$x_i = 0 \rightarrow$ randomly selected Pennsylvania does not home doctors.

Number of P. who home docs = $\sum_{i=1}^{418} x_i = 280 \Rightarrow \hat{p} = \frac{\sum_{i=1}^{418} x_i}{n} = \frac{280}{418} = 0.67$.

$$z_{0.025} = 1.96 \Rightarrow 0.67 \pm 1.96 \sqrt{\frac{0.67(1-0.67)}{418}} = 0.67 \pm 0.045$$

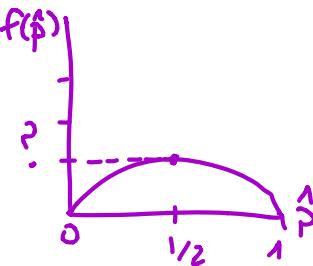
We are 95% confident that (0.625, 0.715) Pennsylvania home doctors

Margin of error: $z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$; Margin: 0.045 = 4.5%

Maximum of $\hat{p}(1-\hat{p}) = \frac{1}{4}$ \Rightarrow PROOF. $f(\hat{p})$

The largest margin of error

is based on n : $z_{\alpha/2} \sqrt{\frac{1}{4}} \cdot \sqrt{\frac{1}{n}}$



$$f(\hat{p}) = \hat{p}(1-\hat{p})$$

peak value: derivative = 0

$$\frac{d(\hat{p}-\hat{p}^2)}{d\hat{p}} = 1 - 2\hat{p} = 0$$

$$\hat{p} = \frac{1}{2}$$

$$f(\hat{p}=\frac{1}{2}) = \frac{1}{2}(1-\frac{1}{2}) = \frac{1}{4}$$

for a 95% CI: $1.96 \sqrt{\frac{1}{4}} \cdot \sqrt{\frac{1}{n}} \approx 2 \sqrt{\frac{1}{4}} \cdot \sqrt{\frac{1}{n}} \approx \frac{1}{\sqrt{n}}$

$\frac{1}{\sqrt{n}}$ approximates the 95% CI margin of error.

Ex. $\frac{1}{\sqrt{418}} = 0.048$: is the approximate margin of error.

When n is small, y or $n-y$ is close to 0:

$$\tilde{p} \pm z_{\alpha/2} \sqrt{\tilde{p}(1-\tilde{p})/n+4} \Rightarrow \tilde{p} \approx \frac{y+2}{n+4} = \frac{\tilde{p}+2}{n+4}$$

Ex: $n = 40$ $y = 8$ $\frac{y}{n} = \hat{p} = \frac{8}{40} = 0.20 \Rightarrow 0.20 \pm 1.645 \sqrt{\frac{(0.2)(0.8)}{40}}$
 $1-\alpha = 0.90$

$$\tilde{p} = \frac{8+2}{40+4} = \frac{10}{44} = 0.227 \Rightarrow 0.227 \pm 1.645 \sqrt{\frac{(0.227)(0.773)}{44}}$$

Two proportions:

for large random samples or approximate $(1-\alpha) \cdot 100\%$ confidence interval for $p_1 - p_2$

$$\hat{p}_1 - \hat{p}_2 \pm z_{\alpha/2} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

Proof: $\hat{p}_1 = \frac{y_1}{n_1} \sim N(p_1, \frac{p_1(1-p_1)}{n_1})$; $\hat{p}_2 = \frac{y_2}{n_2} \sim N(p_2, \frac{p_2(1-p_2)}{n_2})$

$$(\hat{p}_1 - \hat{p}_2) \sim N(p_1 - p_2, \frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2})$$

$$P[-z_{\alpha/2} \leq \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}} \leq z_{\alpha/2}] \approx 1-\alpha$$

etc --- same as
w/ 1 proportion —
you replace p with
 \hat{p} .

Ex:

SAMPLE SIZE

ESTIMATING A POPULATION MEAN μ (NORMAL / CONTINUOUS DATA)

Sample size necessary for estimating μ with $(1-\alpha)100\%$ CI and error $\leq \varepsilon$

$$n = \frac{(z_{\alpha/2})^2 s^2}{\varepsilon^2}$$

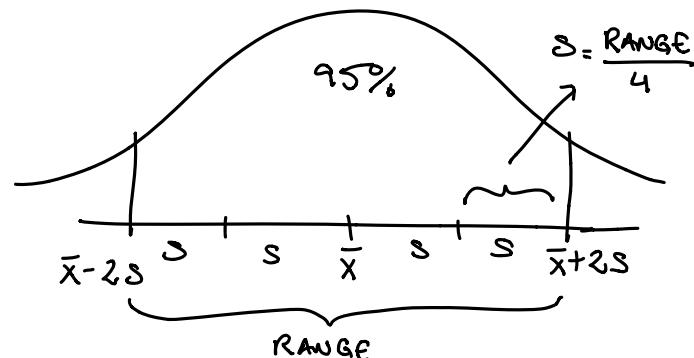
How to determine s^2 :

- scientific lit

- pilot study

- Empirical Rule:

$$\begin{cases} 95\% \text{ of the observations to fall in interval } \bar{x} \pm 2s \Rightarrow s = \frac{\text{MAX-MIN}}{4} \\ 99.7\% \text{ of the observations to fall in interval } \bar{x} \pm 3s \Rightarrow s = \frac{\text{MAX-MIN}}{6} \end{cases}$$



$$\text{MAXIMUM ERROR OF THE ESTIMATE} = \varepsilon = \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}} \Rightarrow n = \frac{z_{\alpha/2}^2 \sigma^2}{\varepsilon^2}$$

$$P \left[\underbrace{\bar{x} - z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right)}_{\varepsilon} \leq \mu \leq \underbrace{\bar{x} + z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right)}_{\varepsilon} \right] = 1-\alpha$$

OK:

$s^2 = 10^2$	$\varepsilon = 1$	$\varepsilon = 3$	$\varepsilon = 5$
90% ($z_{0.05} = 1.645$)	271	31	11
95% ($z_{0.025} = 1.96$)	385	43	16
99% ($z_{0.005} = 2.576$)	664	74	27

We can also change the estimate of the variance. For example, if we change the sample variance to $s^2 = 8^2$, then the necessary sample sizes for various errors ε and confidence levels $(1-\alpha)$ become:

$s^2 = 8^2$	$\varepsilon = 1$	$\varepsilon = 3$	$\varepsilon = 5$
90% ($z_{0.05} = 1.645$)	174	20	7
95% ($z_{0.025} = 1.96$)	246	28	10
99% ($z_{0.005} = 2.576$)	425	48	17

FACTORS AFFECTING SAMPLE SIZE

- STANDARD DEVIATION (s) = $s \uparrow \rightarrow n \uparrow$ bc it is in the numerator
- ERROR (ε) = $\varepsilon \downarrow \rightarrow n \uparrow$ bc it is in the denominator
- CONFIDENCE INTERVAL CI = $CI \uparrow \rightarrow n \uparrow$ bc it's $z_{\alpha/2}$ that is in the numerator

Ex: estimate μ the mean systolic blood pressure of adult Americans w 95% CI
 ε no larger than 3 mm Hg . How many adults should be sampled $\rightarrow n$?
 say sample variance $s^2 = 10^2$

$$\bar{x} \pm \varepsilon = \bar{x} \pm 3$$

$$CI \text{ for a pop. mean } \mu = \bar{x} + t_{\alpha/2, n-1} \left(\frac{s}{\sqrt{n}} \right) \Rightarrow \varepsilon = t_{\alpha/2, n-1} \left(\frac{s}{\sqrt{n}} \right) \Rightarrow \text{solve for } n$$

$$n = \frac{(t_{\alpha/2, n-1})^2 \cdot s^2}{\varepsilon^2} \Rightarrow n \approx \frac{(z_{\alpha/2}^2) s^2}{\varepsilon^2} = \frac{(1.96)^2 (10)^2}{3^2} = 42.7 \approx 43 \text{ people}$$

AS $n \uparrow \rightarrow t$ distribution approaches the $N(0,1)$ dist.

• can't use t value dependent on n if we're trying to discover n .

ESTIMATING A PROPORTION (BINARY DATA)

- LARGE POPULATION

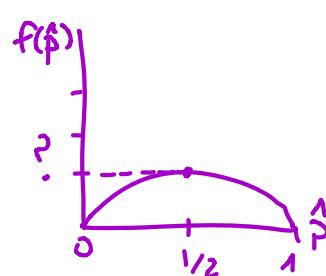
SAMPLE size needed to estimate a pop. proportion \hat{p} with $\{(1-\alpha) \cdot 100\% \text{ CI}$
 $\text{error} \leq \varepsilon$

$$n = \frac{z_{\alpha/2}^2 \hat{p}(1-\hat{p})}{\varepsilon^2}$$

$$\varepsilon = z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

How to determine $\hat{p}(1-\hat{p})$

- use prior knowledge
- set $\hat{p} = \frac{1}{2} = 0.5 \Rightarrow \hat{p}(1-\hat{p}) = \frac{1}{4}$



$$f(\hat{p}) = \hat{p}(1-\hat{p})$$

peak value: derivative = 0

$$\frac{d(\hat{p}(1-\hat{p}))}{d\hat{p}} = 1 - 2\hat{p} = 0$$

$$\hat{p} = \frac{1}{2}$$

$$f(\hat{p}=\frac{1}{2}) = \frac{1}{2}(1-\frac{1}{2}) = \frac{1}{4}$$

Ex: Pollster wants to estimate p = the proportion of Americans who are Dem's.
 95% C $\varepsilon \leq 0.03$. $n = ?$

$$\hat{p} \pm \varepsilon = \hat{p} \pm 0.03$$

$$n = \frac{(1.96)^2 \cdot (0.8)(0.2)}{0.03^2}$$

$$n = \frac{(1.96)^2 (0.5)(0.5)}{0.03^2}$$

$\hat{p} = 0.8$	$\varepsilon = 0.01$	$\varepsilon = 0.03$	$\varepsilon = 0.05$
90% ($z_{0.05} = 1.645$)	4330	482	174
95% ($z_{0.025} = 1.96$)	6147	683	246
99% ($z_{0.005} = 2.576$)	10618	1180	425

$\hat{p} = 0.5$	$\varepsilon = 0.01$	$\varepsilon = 0.03$	$\varepsilon = 0.05$
90% ($z_{0.05} = 1.645$)	6766	752	271
95% ($z_{0.025} = 1.96$)	9604	1068	385
99% ($z_{0.005} = 2.576$)	16590	1844	664

- SMALL POPULATION:

$$n = \frac{m}{1 + \frac{m-1}{N}}$$

$$m = \frac{z_{\alpha/2}^2 \hat{p}(1-\hat{p})}{\varepsilon^2} = \text{sample size needed for large pop.}$$

An approximate $(1-\alpha) \cdot 100\%$ CI for small population proportion is

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n} \cdot \frac{N-n}{N-1}}$$

Proof: we take random sample $X_i = x_1, x_2, \dots, x_n$ of size n from a pop of size N .

$$\text{Ex. } N=2000 \quad . \quad N_1 = \text{yes} \quad N-N_1 = \text{no} \Rightarrow p = \frac{N_1}{N} \quad 1-p = \frac{N-N_1}{N}$$

$x = \text{no. of respondents}$ $x_i=1 = \text{yes}$ $x_i=0 = \text{no}$
 who say yes = $\sum_i x_i$

$$\hat{p} = \frac{\sum_i x_i}{n} = \hat{p} = \frac{x}{n}$$

$$Ex: n \cdot \frac{N_1}{N} = np \quad \text{Var}(X) = np(1-p) \left(\frac{N-n}{N-1} \right)$$

$$E\hat{p} = p \quad \text{Var}(\hat{p}) = \frac{p(1-p)}{n} \cdot \left(\frac{N-n}{N-1} \right)$$

CLT $\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n} \left(\frac{N-n}{N-1} \right)}} \sim N(0, 1) \Rightarrow \hat{p} \pm z_{\alpha/2} \underbrace{\sqrt{\frac{\hat{p}(1-\hat{p})}{n} \cdot \frac{N-n}{N-1}}}_{\text{as small pop.}}$

if $n \ll N$
 $\frac{N-n}{N-1} \approx 1$

$$\Downarrow = \hat{p} \pm z_{\alpha/2} \underbrace{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}}_{\text{as large pop}}$$

As $N \downarrow \rightarrow n \downarrow$ As the pop size \downarrow so does the size of the sample.

SIMPLE LINEAR REGRESSION

Evaluates the relationship between 2 continuous variables:

X : independent / predictor / explanatory variable.

weight
↓
weight

Y : dependent / response / outcome variable

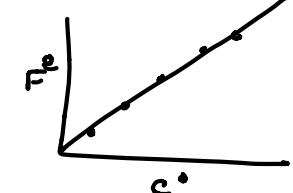
speed
+
gas mileage

$T =$
 \downarrow

gas mileage
gas consumption

DETERMINISTIC RELATIONSHIP: exact relationship $y = \alpha + \beta X$

$$F^0: \frac{9}{5}C^0 + 32$$



SIMPLISTIC RELATIONSHIP: not exact relationship: $y = \alpha + \beta X + \varepsilon$

"Trend" between X and Y , but there is also some "scatter"

$$y_i \sim N(\alpha + \beta(x_i - \bar{x}), \sigma^2)$$

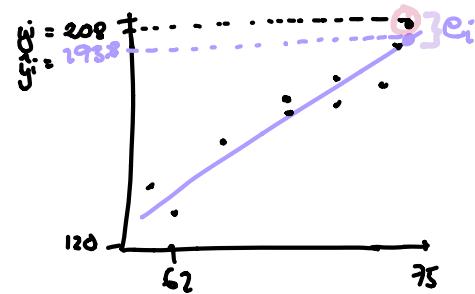
$$y_i = \alpha + \beta(x_i - \bar{x}) + \varepsilon_i$$

Deterministic + Random

which is the best line?

$$\text{weight} = -266.5 + 6.1 \cdot \text{height}$$

$$\text{weight} = -331.2 + 7.1 \cdot \text{height}$$



$$x_i = 75 \\ y_i = 208 \\ \hat{y}_i = -266.5 + 6.1 \cdot 75 = 193.8$$

y_i : observed response of i^{th} experimental unit

x_i : predictor value of i^{th} experimental unit

\hat{y}_i : predicted response of i^{th} experimental unit
fitted value

$e_i = y_i - \hat{y}_i$: residual error of i^{th} experimental unit $\Rightarrow \sum_{i=1}^n (y_i - \hat{y}_i) = 0$
or
predicted error

(LEAST SQUARES CRITERION): you want to find a line that fits the data where e_i is as small as possible for all the observations i : $\hat{y}_i = a_1 + b x_i$; find a and b !

LEAST SQUARES REGRESSION LINE:

minimize the sum of the squared predicted errors

$$Q: \sum_{i=1}^n (\hat{y}_i - y_i)^2 = \sum_{i=1}^n (y_i - (a_1 + b x_i))^2$$

LEAST SQUARES REGRESSION LINE: $\hat{y}_i = a_1 + b x_i = a + b(x_i - \bar{x})$

$x_i = 0$ can be not meaningful, ex:
 $x = \text{height}$ $y = \text{weight}$,
 $x = 0$ makes ... doesn't work

$$\hat{y}_i = a_1 \text{ when } x_i = 0 \quad \hat{y}_i = a \text{ when } x_i = \bar{x}$$

NOT CENTERED? CENTERED MIGHT WORK?

LEAST SQUARES ESTIMATES

$$a = \bar{y}$$

$$b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

does not change with centering

i	x_i	y_i	\hat{y}_i	$y_i - \hat{y}_i$	$(y_i - \hat{y}_i)^2$
1	64	121	123.2	-2.2	4.84
2	73	181	187.1	-6.1	37.21
3	71	156	172.9	-16.9	285.61
4	69	162	158.7	3.3	10.89
5	66	142	137.4	4.6	21.16
6	69	157	158.7	-1.7	2.89
7	75	208	201.3	6.7	44.89
8	71	169	172.9	-3.9	15.21
9	63	127	116.1	10.9	118.81
10	72	165	180.0	-15.0	225.00
<hr/>					
766.51					

i	x_i	y_i	\hat{y}_i	$y_i - \hat{y}_i$	$(y_i - \hat{y}_i)^2$
1	64	121	126.271	-5.3	28.09
2	73	181	181.509	-0.5	0.25
3	71	156	169.234	-13.2	174.24
4	69	162	156.959	5.0	25.00
5	66	142	138.546	3.5	12.25
6	69	157	156.959	0.0	0.00
7	75	208	193.784	14.2	201.64
8	71	169	169.234	-0.2	0.04
9	63	127	120.133	6.9	47.61
10	72	165	175.371	-10.4	108.16
<hr/>					
597.29					



This line would be better.
But there are a 1000 lines!

PROOF:

$$\text{minimize } Q = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (a + b(x_i - \bar{x})))^2$$

$$\frac{\partial Q}{\partial a} : \underbrace{2 \sum_{i=1}^n (y_i - (a + b(x_i - \bar{x})))}_{2} \cdot (-1) \stackrel{\text{set}}{=} 0 \Rightarrow \text{Solve for } a$$

$$-\sum_{i=1}^n y_i + \sum_{i=1}^n a + b \sum_{i=1}^n (x_i - \bar{x}) = -\sum_{i=1}^n y_i + n a = 0 \Rightarrow a = \frac{\sum_{i=1}^n y_i}{n} = \bar{y}$$

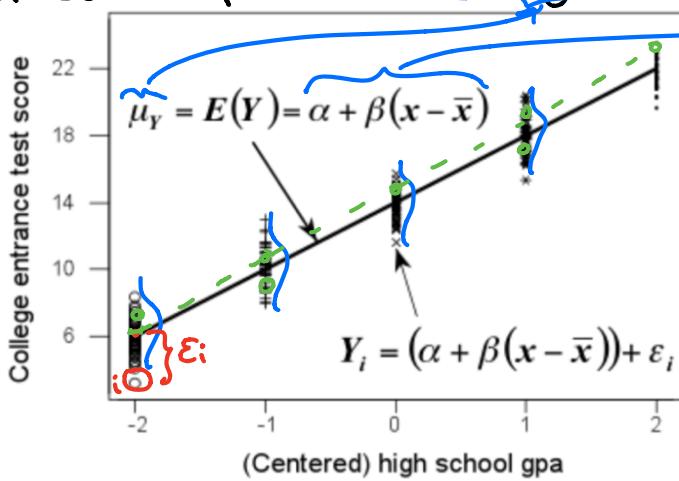
Replace a for \bar{y} in Q formula: $Q = \sum_{i=1}^n (y_i - (\bar{y} + b(x_i - \bar{x})))^2$

$$\frac{\partial Q}{\partial b} = 2 \cdot \sum_{i=1}^n (y_i - (\bar{y} + b(x_i - \bar{x}))) \cdot -(x_i - \bar{x}) \stackrel{\text{set}}{=} 0$$

$$\sum_{i=1}^n (y_i - \bar{y} - b(x_i - \bar{x})) \cdot -(x_i - \bar{x}) = -\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) + b \sum_{i=1}^n (x_i - \bar{x})^2 = 0$$

$$b = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Linear relationship between college entrance score and centered HS GPA w $\begin{cases} \text{intercept } \alpha \\ \text{slope } \beta \end{cases}$



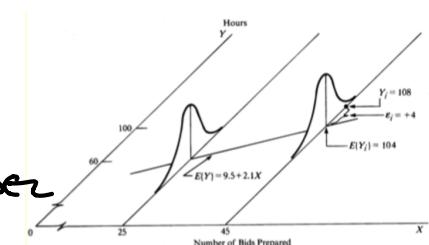
x = high school GPA

$x - \bar{x}$ = centered high school GPA

(± 2 points from the mean)

ASSUMPTIONS OF THE SIMPLE LINEAR REGRESSION MODEL

- The mean of the responses $E(y_i)$ is a linear function of x ;
- The errors ϵ_i , and the responses y_i are independent
- The errors ϵ_i , and the responses y_i are Normally dist. \Rightarrow residual plot $\begin{cases} \text{Residual plot} \\ \text{Assess model fit} \end{cases}$
- The errors ϵ_i , and the responses y_i have Equal variances for all x values LINE!



MAXIMUM LIKELIHOOD ESTIMATORS OF THE SIMPLE LINEAR REGRESSION PARAMETERS

if the 4 assumptions hold true:

	$\hat{\alpha}$	$\hat{\beta}$	$\hat{\sigma}^2$
CENTERED	$a = \bar{y}$	$b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{SS_{xy}}{SS_{xx}}$	$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{y})^2}{n} = \frac{\sum_{i=1}^n (y_i - \hat{y})^2}{n}$
NOT CENTERED	$a = \bar{y} - \hat{\beta}\bar{x}$	SAME	$\hat{\sigma}^2 = \sum_{i=1}^n (y_i - \hat{y})^2$

Proof:

Simple linear regression model states that $E_i \xrightarrow{\text{independent}} \sim N(0, 1)$

$$y_i = \alpha + \beta(x_i - \bar{x}) + \varepsilon_i \Rightarrow y_i \sim N(\alpha + \beta(x_i - \bar{x}), \sigma^2)$$

x_1, x_2, \dots, x_n random sample from normal dist $x_i \sim N(\mu, \sigma^2)$

MLE of μ and σ^2 ?

$$\left[-\frac{(x_i - \theta_1)^2}{2\theta_2} \right] \quad -\infty < \theta_1 < \infty$$

$$\text{pdf of Normal Dist } \left. \begin{array}{l} f(x_i; \theta_1, \theta_2) = \frac{1}{\sqrt{\theta_2} \sqrt{2\pi}} \cdot e^{-\frac{(x_i - \theta_1)^2}{2\theta_2}} \\ \theta_1 = \mu ; \theta_2 = \sigma^2 \end{array} \right\}$$

$$0 < \theta_2 < \infty \quad \left[-\frac{1}{2\theta_2} \sum_{i=1}^n (x_i - \theta_1)^2 \right]$$

$$\text{LIKELIHOOD FUNCTION : } L(\theta_1, \theta_2) = \prod_{i=1}^n f(x_i; \theta_1, \theta_2) = \theta_2^{-n/2} \cdot (2\pi)^{-n/2} \cdot e^{-\frac{1}{2\theta_2} \sum_{i=1}^n (y_i - \alpha - \beta(x_i - \bar{x}))^2}$$

$$\text{Likelihood function : } L(y_i | \alpha, \beta, \sigma^2) = \prod_{i=1}^n L(y_i | \alpha, \beta, \sigma^2) = (\sigma^2)^{-n/2} \cdot (2\pi)^{-n/2} \cdot e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \alpha - \beta(x_i - \bar{x}))^2}$$

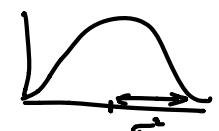
$$\ln(L(y_i)) = -\frac{n}{2} \sigma^2 - \frac{1}{2} \cdot 2\pi \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \alpha - \beta(x_i - \bar{x}))^2 \right)$$

$$\frac{\partial \ln(L(y_i))}{\partial \alpha} = 0 \Rightarrow \alpha$$

to maximize $\ln(L(y_i))$, we need to minimize $\sum_{i=1}^n (y_i - \alpha - \beta(x_i - \bar{x}))^2$ with respect to α and β with respect to α and β

$$\frac{\partial \ln(L(y_i))}{\partial \beta} = 0 \Rightarrow b$$

so the Maximum Likelihood Estimators (MLE) of α and β are the least squares estimators (LSE) of α and β under the simple linear regression assumptions



MAXIMUM LIKELIHOOD ESTIMATOR OF THE VARIANCE σ^2
UNBIASED BIASED.

- ONE POPULATION

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}$$

- MANY POPULATIONS

$$\text{MEAN SQUARE ESTIMATOR} \quad \text{MSE} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2}$$

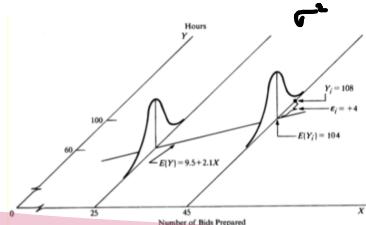
MAX LIKELIHOOD ESTIMATOR

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{y})^2}{n}$$

Proof:

$$\frac{\partial \ln(L(y_i))}{\partial \sigma^2} = 0 = \left[-\frac{n}{2\sigma^2} - \frac{1}{2} \sum_{i=1}^n (y_i - \alpha - \beta(x_i - \bar{x}))^2 \cdot \left(\frac{1}{\sigma^2} \right)^2 \right] = 0$$

$$\Rightarrow -n\sigma^2 + \sum_{i=1}^n (y_i - \alpha - \beta(x_i - \bar{x}))^2 = 0$$



BLUE : Best Linear Unbiased Estimator
 estimator with the smallest variance among all the unbiased estimators

If $\begin{cases} E(y_i) = \alpha + \beta x_i \\ \text{Var}(y_i) = \sigma^2 \\ \text{Cov}(y_i, y_j) = 0 \text{ for } i \neq j \end{cases}$ then the LSE $\begin{cases} \hat{\alpha} \\ \hat{\beta} \end{cases}$ are the BLUE

$\hat{\alpha}$ and $\hat{\beta}$ are unbiased estimators with finite variances

$$E(\hat{\alpha}) = \alpha \quad \text{Var}(\hat{\alpha}) = \frac{\sigma^2}{n}$$

$$E(\hat{\beta}) = \beta \quad \text{Var}(\hat{\beta}) = \frac{\sigma^2}{S_{xx}}$$

$$y_i \sim N(\alpha + \beta(x_i - \bar{x}), \sigma^2)$$

μ_y and σ_y^2 exist and are finite
 y_i 's are random but x_i 's are fixed
 $\hat{\alpha}, \hat{\beta}$ are linear functions of y_1, y_2, \dots, y_n

SAMPLING DISTRIBUTIONS OF THE ESTIMATORS

	$\hat{\alpha}$	$\hat{\beta}$		
CENTERED	$\hat{\alpha} \sim N\left(\alpha, \frac{\sigma^2}{n}\right)$	$\hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{S_{xx}}\right)$	$\text{Cov}(\hat{\alpha}, \hat{\beta}) = -\frac{\sigma^2}{S_{xx}} \bar{x}$	$\frac{n\sigma^2}{\sigma^2} \sim \chi^2_{(n-2)}$
NOT CENTERED	$\hat{\alpha} \sim N\left(\alpha, \frac{\sigma^2 \sum x_i^2}{n S_{xx}}\right)$	$\hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{S_{xx}}\right)$	$\text{Cov}(\hat{\alpha}, \hat{\beta}) = 0$	$\frac{n\sigma^2}{\sigma^2} \sim \chi^2_{(n-2)}$
CENTERED				

CONFIDENCE INTERVALS FOR REGRESSION PARAMETERS

Under the assumptions of the linear regression model (LINE):

$$\hat{\alpha} \sim N(\alpha, \frac{\sigma^2}{n})$$

Proof: Maximum likelihood estimator of $\alpha = \hat{\alpha} = \bar{y}$
least squares estimator

$$y_i \sim N(\alpha + \beta(x_i - \bar{x}), \sigma^2)$$

$$E(\hat{\alpha}) = E(\bar{y}) = \frac{1}{n} \cdot \sum_{i=1}^n E(y_i) = \frac{1}{n} \cdot \sum_{i=1}^n E(\alpha + \beta(x_i - \bar{x})) = \frac{1}{n} \left[n\alpha + \beta \cdot \sum_{i=1}^n (x_i - \bar{x}) \right] = \frac{1}{n} \cdot n\alpha = \alpha$$

$$\text{Var}(\hat{\alpha}) = \text{Var}(\bar{y}) = \frac{\sigma^2}{n}$$

$$\hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)$$

Proof: MLE of β : $\hat{\beta} = b = \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}$

$$y_i \sim N(\alpha + \beta(x_i - \bar{x}), \sigma^2) \quad E(y_i)$$

$$E(\hat{\beta}) = \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sum_{i=1}^n E[(x_i - \bar{x}) y_i] = \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sum_{i=1}^n (x_i - \bar{x}) \cdot (\alpha + \beta(x_i - \bar{x})) =$$

$$\frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \left[\underbrace{\sum_{i=1}^n \alpha}_{0} + \beta \sum_{i=1}^n (x_i - \bar{x})^2 \right] = \beta$$

$$\text{Var}(\hat{\beta}) = \left[\frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]^2 \sum_{i=1}^n (x_i - \bar{x})^2 \text{Var}(y_i) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sigma^2}{s_{xx}}$$

$$\bullet \quad n \cdot \frac{\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{(n-2)} \quad \therefore \alpha = \hat{\alpha}, \beta = \hat{\beta}, \frac{\hat{\sigma}^2}{\sigma^2} \Rightarrow \text{mutually independent.}$$

• a $(1-\alpha) \cdot 100\%$ confidence interval for the slope parameter β is

$$b \pm t_{\alpha/2, n-2} \left(\frac{\sqrt{n} \hat{\beta}}{\sqrt{n-2} \cdot \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \right) \approx \hat{\beta} \pm t_{\alpha/2, n-2} \cdot \sqrt{\frac{\text{MSE}}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

PROOF : $T = \frac{\hat{\beta} - \beta}{\sqrt{\frac{\text{MSE}}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \rightarrow T \sim N(0, 1)$ independent

$$\frac{\hat{\beta} - \beta}{\sqrt{\frac{\text{MSE}}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \sim N(0, 1) \quad \hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)$$

$$T \sim t_{(n)}$$

$$\frac{n \hat{\sigma}^2}{\sigma^2} \sim \chi^2_{(n-2)}$$

$$\begin{aligned}
 T_1 &= \frac{\frac{\hat{\beta} - \beta}{\sigma}}{\sqrt{\frac{\sum \varepsilon(x_i - \bar{x})^2}{n-2}}} = \frac{\hat{\beta} - \beta \cdot \sqrt{\sum \varepsilon(x_i - \bar{x})^2}}{\sqrt{\frac{n \hat{\sigma}^2}{n-2}} \sqrt{\frac{1}{n-2}}} = \frac{\hat{\beta} - \beta}{\frac{n \cdot \hat{\sigma}^2}{(n-2)} \sqrt{\sum \varepsilon(x_i - \bar{x})^2}} \\
 &= \frac{\hat{\beta} - \beta}{\frac{MSE}{\sum \varepsilon(x_i - \bar{x})^2}} \sim t(n-2)
 \end{aligned}$$

• a $(1-\alpha) \cdot 100\%$ confidence interval for the slope parameter α is

$$\alpha \pm t_{\alpha/2, n-2} \cdot \sqrt{\frac{\hat{\sigma}^2}{n-2}} \Rightarrow \alpha \pm t_{\alpha/2, n-2} \cdot \sqrt{\frac{MSE}{n}}$$

$$T_2 = \frac{\hat{\alpha} - \alpha}{\sqrt{\frac{\hat{\sigma}^2}{n-2}}} \sim t(n-2)$$

CORRELATION \neq CAUSATION
 \neq REGRESSION: in correlation x, y are interchangeable

MORE REGRESSION :

ESTIMATION : To know the value of the mean response $E(y) = \mu_y$ for a given value x , calculate the CONFIDENCE INTERVAL for the MEAN $E(y) = \mu_y$

PREDICTION : To know the value of a new observation y_{n+1} for a given value x , calculate the PREDICTION INTERVAL for the NEW OBSERVATION y_{n+1}

CONFIDENCE INTERVAL FOR THE MEAN $E(y) = \mu_y$

$$\hat{y} \pm t_{\alpha/2, n-2} \sqrt{MSE} \sqrt{\frac{1}{n} + \frac{(x-\bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

Proof: point estimate for μ_y is $\hat{y} = \hat{\alpha} + \hat{\beta}(x - \bar{x})$ and $\begin{cases} \hat{\alpha} \sim N(\alpha, \frac{\sigma^2}{n}) \\ \hat{\beta} \sim N(\beta + \frac{\sigma^2}{n} \sum_{i=1}^n (x_i - \bar{x})^2) \\ \frac{n\hat{\beta}^2}{\sigma^2} \sim \chi^2_{(n-2)} \end{cases}$ dependent

$$\begin{aligned} \text{Var}(\hat{y}) &= \text{Var}[\hat{\alpha} + \hat{\beta}(x - \bar{x})] = \text{Var}(\hat{\alpha}) + (x - \bar{x})^2 \text{Var}(\hat{\beta}) = \frac{\sigma^2}{n} + (x - \bar{x})^2 \cdot \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \\ &= \sigma^2 \left[\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \end{aligned}$$

$$\hat{y} \sim N(\mu_y, \sigma^2 \left[\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right])$$

$$T = \frac{\frac{\hat{y} - \mu_y}{\sigma \cdot \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}}}{\sqrt{\frac{n\hat{\beta}^2}{\sigma^2} / (n-2)}} = \frac{\hat{y} - \mu_y}{\sqrt{MSE} \cdot \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \sim t_{(n-2)}$$

$$MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2}$$

$$P \left[-t_{\alpha/2, (n-2)} \leq \frac{\hat{y} - \mu_y}{\sqrt{MSE} \cdot \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \leq t_{\alpha/2, (n-2)} \right] = 1 - \alpha$$

How to make the CI for the mean now?

$$\hat{y} \pm t_{\alpha/2, n-2} \sqrt{MSE} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum(x_i - \bar{x})^2}}$$

- ① estimate μ_y at \bar{x} = the mean of the predictor values
shortest CI for μ_y when $\bar{x} = x \Rightarrow (x - \bar{x})$ is smallest
- ② \downarrow confidence level $\rightarrow \downarrow t_{\alpha/2, n-2}$
- ③ \uparrow sample size $\rightarrow \uparrow n \rightarrow \downarrow \frac{1}{n}$
- ④ choose spread out predictor values $x_i \rightarrow \uparrow$ spread $\rightarrow \uparrow \sum(x_i - \bar{x})^2$

Ex: x = duration of previous eruptions y = time until next eruption
mean time till next eruption = $\bar{y} = 8.5$ min. $\rightarrow 95\%$ confidence that the mean time until next eruption will be between $(70.14, 72.703)$ min

PREDICTION INTERVAL for the new observation y_{n+1}

A $(1-\alpha) \cdot 100\%$ prediction interval for a new observation y_{n+1} when the predictor $x = x_{n+1}$

$$\hat{y}_{n+1} \pm t_{\alpha/2, n-2} \sqrt{MSE} \sqrt{1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum(x_i - \bar{x})^2}}$$

PROOF: $y_{n+1} \sim N(\underbrace{\alpha + \beta(x_{n+1} - \bar{x})}_{\text{and}}, \sigma^2)$

$$w = y_{n+1} - \hat{y}_{n+1} = y_{n+1} - \hat{\alpha} - \hat{\beta}(x_{n+1} - \bar{x})$$

$$\left. \begin{aligned} \hat{\alpha} &\sim N\left(\alpha, \frac{\sigma^2}{n}\right) \\ \hat{\beta} &\sim N\left(\beta, \frac{\sigma^2}{\sum(x_i - \bar{x})^2}\right) \\ \frac{n\hat{\beta}^2}{\sigma^2} &= \frac{(n-2)MSE}{\sigma^2} \sim \chi^2_{(n-2)} \end{aligned} \right\} \text{SPEZ. SCHEIT}$$

$$\begin{aligned} E(w) &= E[y_{n+1} - \hat{\alpha} - \hat{\beta}(x_{n+1} - \bar{x})] = \underbrace{E(y_{n+1}) - E(\hat{\alpha}) - (x_{n+1} - \bar{x}) E(\hat{\beta})}_{= 0} = \\ &= \cancel{\alpha + \beta(x_{n+1} - \bar{x})} - \cancel{\alpha} - (x_{n+1} - \bar{x}) \beta = 0 \end{aligned}$$

$$\begin{aligned} \text{Var}(w) &= \text{Var}[y_{n+1} - \hat{\alpha} - \hat{\beta}(x_{n+1} - \bar{x})] = \text{Var}(y_{n+1}) - \text{Var}(\hat{\alpha}) - (x_{n+1} - \bar{x})^2 \text{Var}(\hat{\beta}) = \\ &= \sigma^2 + \frac{\sigma^2}{n} + \frac{(x_{n+1} - \bar{x})^2 \sigma^2}{\sum(x_i - \bar{x})^2} = \sigma^2 \left[1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum(x_i - \bar{x})^2} \right] \end{aligned}$$

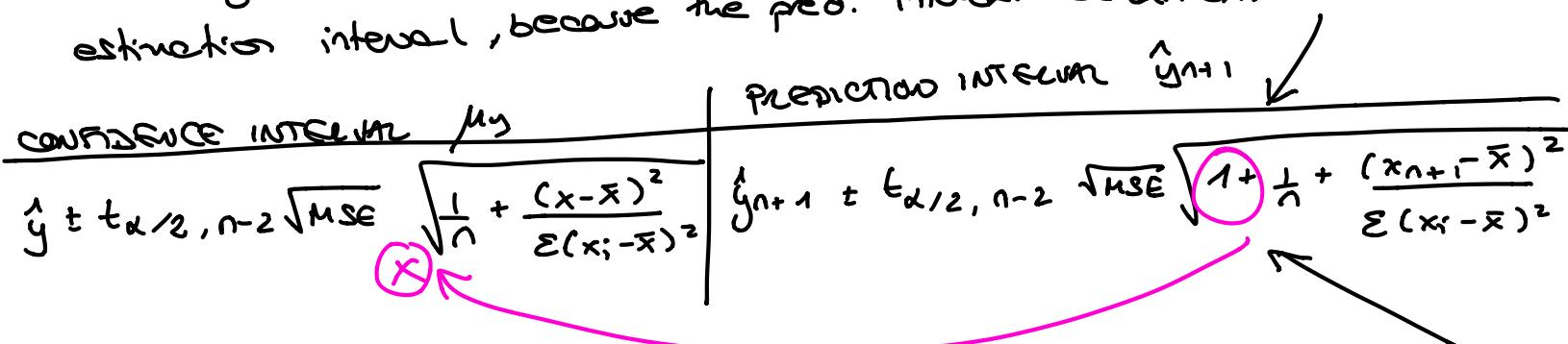
$$w \sim N(0, \sigma^2 \left[1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum(x_i - \bar{x})^2} \right])$$

$$T = \frac{\frac{(y_{n+1} - \hat{y}_{n+1}) - 0}{\sigma^2 \left[1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]}}{\sqrt{\frac{n \sigma^2}{n-2}} / (n-2)} = \frac{(y_{n+1} - \hat{y}_{n+1})}{\sqrt{\text{MSE}} \sqrt{1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum (x_i - \bar{x})^2}}} \sim t_{n-2}$$

$$P \left[-t_{\alpha/2, n-2} \leq \frac{(y_{n+1} - \hat{y}_{n+1})}{\sqrt{\text{MSE}} \sqrt{1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum (x_i - \bar{x})^2}}} \leq t_{\alpha/2, n-2} \right] = 1-\alpha$$

Ex: x = duration of previous eruptions y = time until next eruption
 previous eruption $\bar{y} = 8.5$ min. $\rightarrow 95\%$ confidence that the time until next eruption will be between $(58.109, 84.724)$ min

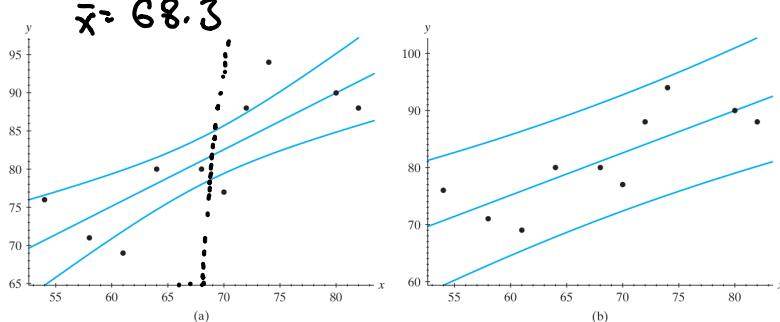
- for a given value x , the prediction interval is always larger than the estimation interval, because the pred. interval has an extra term in the SE.



- we can't make the standard error of the prediction for y_{n+1} approach 0

- make narrower \rightarrow same as with C.I. for μ_y .

CONFIDENCE BAND: collection of all $(1-\alpha) \cdot 100\%$ C.I. is for $\mu(x) : -\infty < x < \infty$
PREDICTION BAND: collection of all $(1-\alpha) \cdot 100\%$ pred.I. is for $y(x) = \alpha + \beta x + \varepsilon : -\infty < x < \infty$



The bands are narrowest at $x = \bar{x}$.
 Prediction band is always wider: the difference between 1 observation of y and its predictor varies more than the difference between the mean of the pop of y values and its estimator.

Figure 7.6-1 A pointwise 95% (a) confidence band for $\mu(x)$ and (b) prediction band for y

MULTIPLE REGRESSION

GENERAL LINEAR MODEL: $E(Y) = E(Y|x_1, \dots, x_K) = E(Y|x) = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_K x_K$

Ex.: x_1 : student's ACT score x_2 : student's US class rank y : 1st year GPA at college

To fit the linear model, $\beta_1 x_1 + \beta_2 x_2 + \dots + \beta_K x_K$. n observation points: $x_{ij}, x_{2j}, \dots, x_{Kj}$, y_j

① Minimize the function by least squares, $G = \sum_{j=1}^n (y_j - \beta_1 x_{1j} - \beta_2 x_{2j} - \dots - \beta_K x_{Kj})^2$

② K 1st partial derivatives and $\stackrel{\text{set}}{=} 0$.

$$\frac{\partial G}{\partial \beta_i} = \sum_{j=1}^n (-2)(y_j - \beta_1 x_{1j} - \beta_2 x_{2j} - \dots - \beta_K x_{Kj})(x_{ij}) \stackrel{\text{set}}{=} 0$$

K normal equations

$$\begin{aligned} \beta_1 \sum_{j=1}^n x_{1j}^2 + \beta_2 \sum_{j=1}^n x_{1j} x_{2j} + \dots + \beta_K \sum_{j=1}^n x_{1j} x_{Kj} &= \sum_{j=1}^n x_{1j} y_j, \\ \beta_1 \sum_{j=1}^n x_{2j} x_{1j} + \beta_2 \sum_{j=1}^n x_{2j}^2 + \dots + \beta_K \sum_{j=1}^n x_{2j} x_{Kj} &= \sum_{j=1}^n x_{2j} y_j, \\ \vdots &\quad \vdots \quad \ddots \quad \vdots \quad \vdots \\ \beta_1 \sum_{j=1}^n x_{Kj} x_{1j} + \beta_2 \sum_{j=1}^n x_{Kj} x_{2j} + \dots + \beta_K \sum_{j=1}^n x_{Kj}^2 &= \sum_{j=1}^n x_{Kj} y_j. \end{aligned}$$

③ Find solutions for $\beta_1, \dots, \beta_K \Rightarrow$ least squares estimates = MLE
 provided the random variables y_1, y_2, \dots, y_n are mutually independent and
 y_j is $N(\beta_1 x_{1j} + \beta_2 x_{2j} + \dots + \beta_K x_{Kj}, \sigma^2)$ $j=1 \dots n$

INTERPRETATION OF β : change in $E(Y)$ when x_i is increased by 1 unit,
 while each of the x_2, \dots, x_K are held constant (=adjusting / controlling
 other x).
 Caution! when data are transformed ex $\log x$, $\log y$.

HYPOTHESIS TESTING.

HYPOTHESIS: proposed explanation for a phenomenon

NUL HYPOTHESIS: H_0 : initial assumption about the population parameter

ALTERNATIVE HYP. = H_1 :

SIMPLE HYP.: each hypothesis has a single value

Ex: we assume $X \sim N(\mu, \sigma^2)$, $\mu = 50$ or $\mu = 55$. $H_0: \mu = 50$ then $H_1: \mu = 55$

COMPOSITE HYP.: A hypothesis can be a composite of different values / distributions.

Ex: we assume $X \sim N(\mu, \sigma^2)$. $H_0: \mu = 50$; $H_1: \mu > 50 \Rightarrow \mu = 51, \mu = 52 \dots$

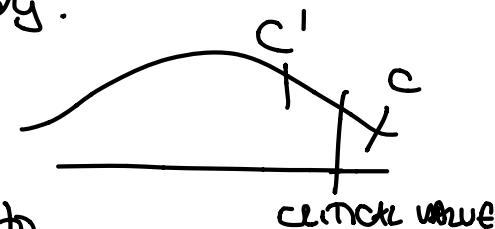
TEST STATISTIC: function of the sample, numerical summary.

CRITICAL REGION: Rejection region for H_0 . = C

ACCEPTANCE REGION: Acceptance region for H_0 . = C'

CRITICAL VALUE: threshold / cut-off value, corresponds to

SIGNIFICANCE LEVEL: probability of a Type-I error.



P-VALUE VS. CI: 95% CI ≈ interval containing all values where $p > 0.05$.
if $p < 0.05$ or C contains 0 → reject H_0 .

ERRORS

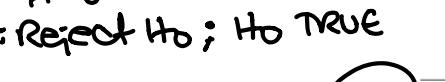


Table 2: True state vs. Decision

TYPE I ERROR: α . FALSE \ominus : Reject H_0 ; H_0 TRUE

$$P[(x_1, x_2, \dots, x_n) \in C; H_0] = c/a+c$$

= SIGNIFICANCE LEVEL.

"Fail to reject"

	True state	
	H_0 is true	H_0 is false
Decision	Accept H_0	a
Reject H_0	c	b

TYPE II ERROR: β . FALSE \oplus : Accept H_0 ; H_0 FALSE

$$P[(x_1, x_2, \dots, x_n) \in C'; H_1] = b/b+\delta$$

POSTERIOR PROB. TYPE I ERROR: $c/c+\delta$ } BAYESIAN

POSTERIOR PROB. TYPE II ERROR: $b/a+b$

POWER FUNCTION: $1 - \beta$: Reject H_0 ; H_1 TRUE

Ex: None to place right side. $H_0: p=0.5$ $H_1: p < 1/2$ $n=20$

$y = \sum_{i=1}^{20} x_i \sim \text{bin}(20, p)$ CRITICAL REGION, $C = [y \leq 6] = \sum_{i=1}^{20} y \leq 6$

$$\alpha = P(y \leq 6; p = \frac{1}{2}) = \sum_{y=0}^6 \binom{20}{y} \left(\frac{1}{2}\right)^y \left(\frac{1}{2}\right)^{20-y} \quad \text{example}$$

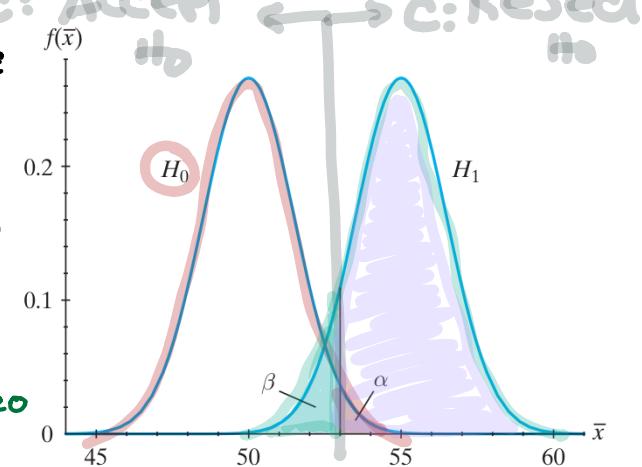


Figure 8.1-1 pdf of \bar{X} under H_0 and H_1

$$K_p = P[y \leq 6; p \leq \frac{1}{2}] = \sum_{y=0}^6 \binom{20}{y} \cdot (p)^y \cdot (1-p)^{20-y}$$

Power is a function because p can be any value and the power changes depending on the value of p . In this example p can be any value $0 < p < 1/2$ and the power changes accordingly.

• HYPOTHESIS TESTING: OK: 4-sided die tossed 1000 times
= CLINICAL VALUE APPROACH 290 # 4 observed. → is this more than expected?

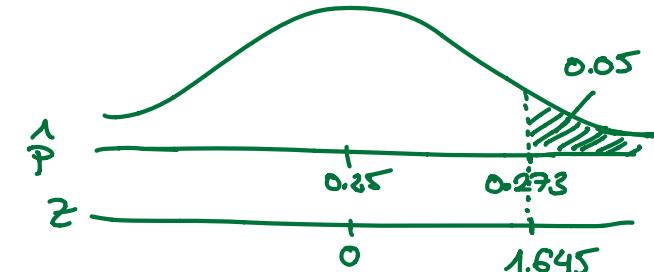
$$\textcircled{1} \text{ STATE } \left\{ \begin{array}{l} H_0 \\ H_1 \end{array} \right.$$

H_0 : unbiased die ⇒ each side = likely $\Rightarrow \frac{1}{4} = 0.25 = p$.
 H_1 : biased die \Rightarrow sides \neq likelihood $\Rightarrow p > 0.25$

\textcircled{2} Decide appropriate test statistic

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$

one-tailed (right-tailed) \hat{p}
0.05



$$p = 0.273 \Rightarrow z = 1.645$$

if our statistic is more extreme than this value \Rightarrow reject the H_0 . = threshold value.

CRITICAL REGION:

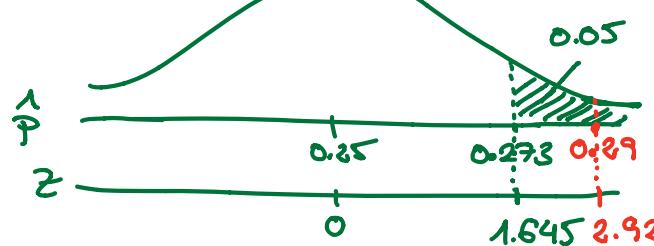


sample proportion $\hat{p} = \frac{290}{1000} = 0.29$.

sampling distribution $p \sim N(p=0.25, \sigma = \sqrt{\frac{0.25 \cdot 0.75}{1000}} = 0.01363)$

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \sim N(0,1) \Rightarrow z = \frac{0.29 - 0.25}{0.01363} = 2.92$$

Reject H_0
Fail to reject H_1



Errors

Type I error: Reject H_0 ; H_0 true

α

Type II error: fail to reject H_0 ; H_0 false

β

⇒ ex: if pop. proportion $p = 0.27$.

$$\Pr(\beta) = P(\hat{p} < 0.273; \text{if } p = 0.27) =$$

$$P\left[z < \frac{0.273 - 0.27}{\sqrt{\frac{0.27 \cdot 0.73}{1000}}}\right] = P(z < 0.214) = 0.5847$$

• SIGNIFICANCE TESTING = ex: cancer patients: 90% die within 3 years.
 P-VALUE APPROACH.

new trt \rightarrow reduced rate? $y=128$ die
 SIGNIFICANCE LEVEL: $\alpha = 0.05$ $n = 150 \rightarrow \hat{p} = \frac{128}{150} = 0.853$

① STATE H_0
 H_1

ONE TAILED TEST: $H_0: p=0.90$ $H_1: p < 0.90$

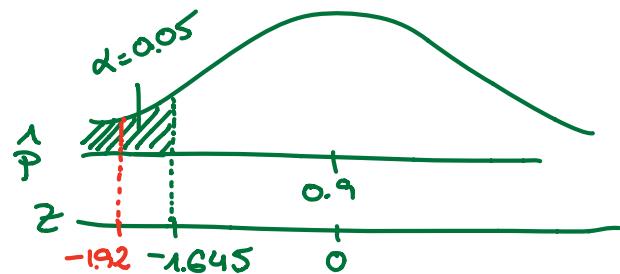
ONE TAIL!

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{0.853 - 0.90}{\sqrt{\frac{0.90 \cdot 0.10}{150}}} = -1.92$$

$\alpha = 0.05$

$$-1.92 < -1.645$$

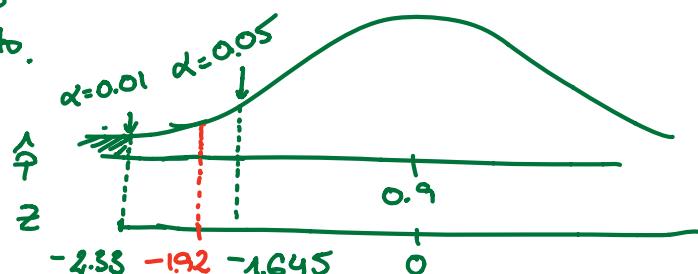
Reject the H_0 .



$\alpha = 0.01$

$$-1.92 > -2.33$$

Fail to reject the H_0 .



P-VALUE:

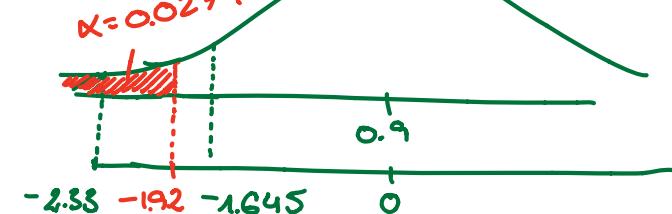
smallest significance level α that leads to us rejecting the H_0
 or

the probability that we'd observe a more extreme statistic than we did if the H_0 were true

$$0.0274 < 0.05 \Rightarrow \text{reject } H_0$$

if p-value $< \alpha \Rightarrow$ Reject H_0

$P(\text{p-value} \leq \alpha | H_0) \leq \alpha \Rightarrow$ random vde



p-value $\left\{ \begin{array}{l} \text{The smallest } \alpha\text{-level that would lead to rejection of } H_0 \\ \alpha\text{-level associated with the test statistic -1.92} \\ P(z < -1.92) = 0.0274 = \text{probability that we'd observe a } \hat{p} = 0.853 \text{ if the true proportion were } 0.9 (\text{=} H_0 \text{ TRUE}) \end{array} \right.$

Ex: TWO-TAILED TEST : ① $H_0: p=0.9$
 SIGNIFICANCE LEVEL $\frac{\alpha}{2} = 0.025$ ② $H_1: p \neq 0.9$

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{0.853 - 0.90}{\sqrt{\frac{0.90 \cdot 0.10}{150}}} = -1.92$$

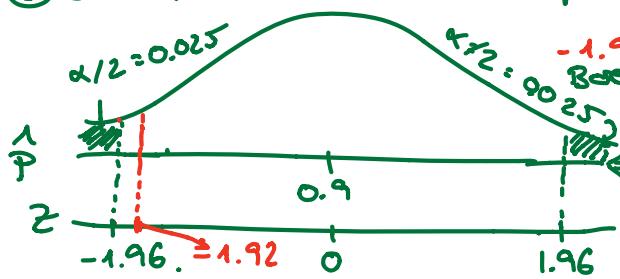
HYPOTHESIS TESTING :

Reject H_0 when:

$$|z| > 1.96 \rightarrow z > 1.96 \quad \downarrow z < -1.96$$

-1.92 > -1.96
 Barely fail to reject H_0 .

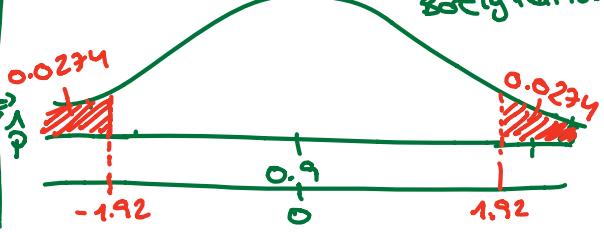
Inufficient evidence
 at $\alpha = 0.05$ level
 to conclude that the sample p differs
 significantly from $p=0.9$



SIGNIFICANCE TESTING.

$$\begin{aligned} ③ \text{P-value} &= P(|z| > 1.92) = P(z > 1.92) + P(z < -1.92) \\ &= 2 \times 0.0274 = 0.055 > 0.05 \end{aligned}$$

Barely fail to reject H_0



TESTS ABOUT PROPORTIONS:

• ONE SAMPLE PROPORTION:

$Y \sim \text{binomial}(n, p)$ $y = \text{no. of successes}$

- ① $H_0: p = p_0$ $H_1: p \neq p_0 \Rightarrow \text{2-sided test}$ \hat{p} is a probability of success = $\frac{y}{n}$
 $H_1: p > p_0$ or $H_1: p < p_0 \Rightarrow \text{1-sided test}$.

- ② STATISTIC: $Z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \sim N(0,1)$ when $n \rightarrow \infty$ due to the CLT (normal approximation of binomial)
 Advantage: better approximation to α -level significance tests (p-value approach)

- ③ a) TWO-TAIL:
 Reject H_0 if $|Z| > z_{\alpha/2}$. Under H_0 $P(|Z| > z_{\alpha/2}) \approx \alpha$.

ONE-TAIL $Z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} > z_\alpha$ or $Z = \frac{y - np_0}{\sqrt{\frac{np_0(1-p_0)}{n}}} \leq z_\alpha$

Advantage: the interpretation of the confidence interval is consistent with the hypothesis test decision
 Also better estimate of the SD if the H_0 is clearly false.

② b) ALTERNATIVE STATISTIC $Z = \frac{\hat{p} - p_0}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}}$

Two-tail $H_1: p \neq p_0$
 = if p_0 is not in the $(1-\alpha) \cdot 100\%$ confidence interval

$$\begin{cases} Z = \frac{\hat{p} - p_0}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} > z_{\alpha/2} \\ \hat{p} - p_0 > z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \\ p_0 < \hat{p} - z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \end{cases}$$

$$\begin{cases} Z = \frac{\hat{p} - p_0}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} < -z_{\alpha/2} \\ \hat{p} - p_0 < -z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \\ p_0 > \hat{p} + z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \end{cases}$$

Reject H_0 if

One-tail $H_1: p < p_0$
 = if p_0 is not in the upper $(1-\alpha) \cdot 100\%$ confidence interval

$$\begin{cases} Z = \frac{\hat{p} - p_0}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} < -z_\alpha \\ \hat{p} - p_0 > -z_\alpha \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \\ p_0 > \hat{p} + z_\alpha \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \end{cases} = [0, \hat{p} + z_\alpha \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}]$$

Ex: Manufacturer observes approx $p = 0.06$ 6 circuits fail.
 New procedure → trials 4; $\sim \text{bin}(200, p) \Rightarrow \hat{p} = \frac{y}{200}$ is the new procedure better?

- C.I. approach: $[0, \hat{p} + 1.645 \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}]$ → if it contains 0.06 → new procedure not necessarily better.
 ↓ if doesn't contain 0.06 → 95% confidence that the true p is less than 0.06. → new p. better.

• HYPOTHESIS approach

$$\textcircled{1} \quad H_0: p = 0.06 = \hat{p} \quad H_1: \hat{p} < 0.06 \quad \text{if } p = 0.06, y = 12 \text{ failures}$$

$$\textcircled{2} \quad \text{CRITICAL VALUE: New p. is better if } p = 0.035, y = 7 \text{ failures}$$

Probability of $p = 0.06$ and observing 7 failures \Rightarrow Type I error $= \alpha$

$$\text{Binomial} \Rightarrow \alpha = P(Y \leq 7; p = 0.06) = \sum_{y=0}^7 \binom{200}{y} (0.06)^y (0.94)^{200-y} = 0.0829$$

$$\text{Poisson approximation} \Rightarrow \alpha \approx \sum_{y=0}^7 \frac{12^y e^{-12}}{y!} = 0.090 \quad \begin{matrix} \text{because } n \text{ small, } p \text{ large} \\ \text{low!} \end{matrix}$$

Prob. of no improvement to for ex. $p = 0.03$ and observing 7 failures \Rightarrow Type II error $= \beta$

$$\text{Binomial} \Rightarrow \beta = P(Y > 7; p = 0.03) = \sum_{y=8}^{200} \binom{200}{y} (0.03)^y (0.97)^{200-y}$$

$$\text{Poisson} \Rightarrow \beta \approx 1 - \sum_{y=0}^7 \frac{6^y e^{-6}}{y!} = 1 - 0.744 = 0.256 \rightarrow \text{high!} \quad \therefore$$

$\frac{Y}{n} \sim N(\mu = p_0, \sigma^2 = \frac{p_0(1-p_0)}{n})$ when $n \rightarrow \infty$ under $H_0: Z \sim N(0, 1)$

if $Z = \frac{\frac{Y}{n} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \geq z_\alpha \Rightarrow \frac{Y}{n}$ exceeds p_0 by z_α standard deviations of $\frac{Y}{n}$, we reject H_0 .

The prob of $Z = \frac{\frac{Y}{n} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \geq z_\alpha$ when H_0 is true = α

significance level = α

• TWO SAMPLE PROPORTION

$$\textcircled{1} \quad H_0: p_1 = p_2$$

$$\textcircled{2} \text{ STATISTIC: } z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\hat{p}(1-\hat{p}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

$$H_1: p_1 - p_2 = 0 \quad \text{or} \quad p_1 \neq p_2$$

$\hat{p} = \frac{y_1 + y_2}{n_1 + n_2}$: proportion of successes of the two samples combined.

ONE-TAIL

$$H_1: p_1 < p_2$$

$$\text{2 TAIL: } H_1: p_1 \neq p_2$$

$$\text{Reject } H_0 \quad | \quad | z | = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\hat{p}(1-\hat{p}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \geq z_{\alpha/2}$$

$$H_1: p_1 > p_2$$

$$z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\hat{p}(1-\hat{p}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} > z_\alpha$$

$$z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\hat{p}(1-\hat{p}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} < z_\alpha$$

$$\text{Proof: } \hat{p}_1 - \hat{p}_2 \sim N(\mu = p_1 - p_2, \sigma^2 = \frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2})$$

$$\text{if } H_0 \text{ is true} \rightarrow p_1 = p_2 = p \Rightarrow \sigma^2 = p(1-p) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$$

$$z = \frac{\hat{p}_1 - \hat{p}_2 - (p_1 - p_2)}{\sqrt{p(1-p) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \sim N(0,1)$$

$$\text{we don't know the assumed common } p = p_1 = p_2 \Rightarrow \text{estimate } \hat{p} = \frac{y_1 + y_2}{n_1 + n_2}$$

$$\text{STATISTIC: } z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\hat{p}(1-\hat{p}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

\textcircled{2B}

ALTERNATIVE:

$$z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\frac{p_1(1-\hat{p}_1)}{n_1} + \frac{p_2(1-\hat{p}_2)}{n_2}}}$$

Advantage: the interpretation of the confidence interval is consistent with the hypothesis test decision

Ex: Should cigarette tax be raised? $y_i = \text{"yes"}$. Ac the opinions between groups \neq ?

	No-smokers	smokers
n	$n_1 = 605$	$n_2 = 195$
y	$y_1 = 851$	$y_2 = 41$
\hat{p}	$\hat{p}_1 = \frac{851}{605} = 0.58$	$\hat{p}_2 = \frac{41}{195} = 0.21$

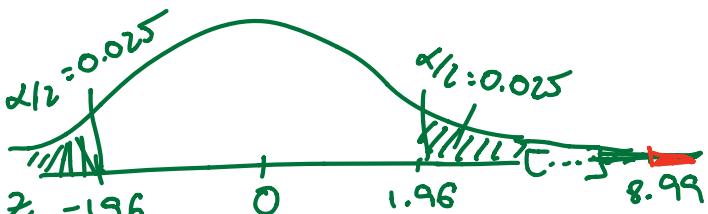
$$\hat{p} = \frac{41 + 851}{605 + 195} = 0.49 \quad \textcircled{1} \quad H_0: p_1 = p_2 \\ H_1: p_1 \neq p_2$$

$$z = \frac{(0.58 - 0.21) - 0}{\sqrt{0.49 \cdot 0.51 \left(\frac{1}{605} + \frac{1}{195} \right)}} = 8.99.$$

$z > z_{\alpha/2} \rightarrow \text{Reject } H_0$.

$$P = 2 \times P(z > 8.99) = 2(0,0000) = 0$$

$$P\text{-value} < 0.0001$$



TESTS ABOUT MEANS

ASSUMPTION: sample pop. is $N \sim (\mu, \sigma)$
random sample

• ONE SAMPLE MEAN

TWO-TAIL

$$\textcircled{1} \quad H_0: \mu = \mu_0 \quad H_1: \mu \neq \mu_0$$

ONE-TAIL

$$\text{or } H_1: \mu < \mu_0 \quad \text{or } H_1: \mu > \mu_0$$

② STATISTIC

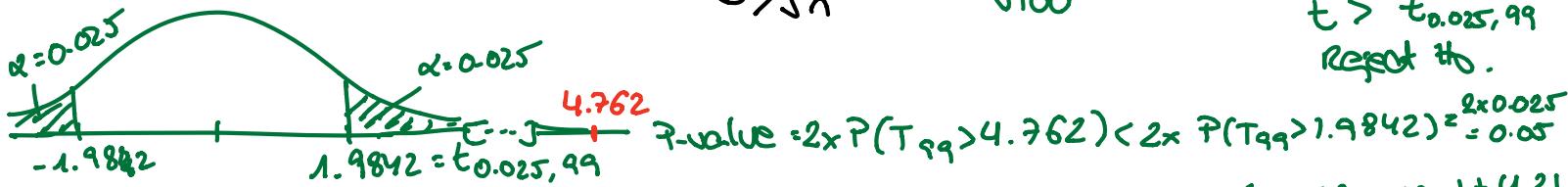
- KNOWN VARIANCES : $Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \sim N(0,1)$

- UNKNOWN VARIANCES: $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \sim t_{(n-1)}$ if the data are normally distributed.

Ex: mean blood pressure in honolulu : $\mu = 120$ mmHg. $n = 100$. $\bar{X} = 130.1$. $SD = 21.2$ mmHg

$$\textcircled{1} \quad H_0: \mu = 120 \quad H_1: \mu \neq 120 \quad \textcircled{2} \quad T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} = \frac{130.1 - 120}{21.2 \sqrt{100}} = \underline{4.762} \quad \textcircled{3} \quad \alpha = 0.05 \quad t_{0.025, 99} = 1.9842$$

$t > t_{0.025, 99}$
reject H_0 .



$$\text{P-value} = 2 \times P(T_{99} > 4.762) < 2 \times P(T_{99} > 1.9842) = \frac{\alpha}{2} = 0.025$$

Equivalent to C.I.: $\bar{X} \pm t_{0.025, 99} \left(\frac{s}{\sqrt{n}} \right) = 130.1 \pm 1.9842 \cdot \left(\frac{21.21}{\sqrt{100}} \right) = 130.1 \pm 4.21$

95% confident the honolulu's blood presue is between [125.89, 134.31]

- PAIRED T-TEST

$$\textcircled{1} \quad H_0: \mu_x = \mu_y \Rightarrow \mu_x - \mu_y = \mu_D = 0 \quad H_1: \mu_D \neq 0$$

$$\textcircled{2} \quad \text{STATISTIC: } \frac{\bar{D} - \mu_0}{SD/\sqrt{n}} \sim t_{(n-1)} \Rightarrow \frac{\bar{D} - 0}{\frac{SD}{\sqrt{n}}}$$

• TWO SAMPLE MEAN

ASSUMPTION: sample pop. of $N \sim (\mu, \sigma)$
random sample

- UNKNOWN EQUAL VARIANCES:

① $H_0: \mu_x = \mu_y$ or $\mu_x - \mu_y = 0$ $H_1: \mu_x - \mu_y \neq 0, \mu_x - \mu_y < 0, \mu_x - \mu_y > 0$

② STATISTIC: $T = \frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_y)}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{(n+m-2)}$

POOLED 2-SAMPLE
T-TEST

pooling sample variance.
 $S_p = \frac{(n-1)S_x^2 + (m-1)S_y^2}{n+m-2}$

PROOF:
 $\bar{x} \sim N(\mu_x, \frac{\sigma^2}{n}) ; \bar{y} \sim N(\mu_y, \frac{\sigma^2}{m}) \Rightarrow \bar{x} - \bar{y} \sim N(\mu_x - \mu_y, \frac{\sigma^2}{n} + \frac{\sigma^2}{m})$ because x, y independent

$Z = \frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_y)}{\sqrt{\frac{\sigma^2}{n} + \frac{\sigma^2}{m}}} \sim N(0, 1)$ bc x, y are normal $U = \frac{(n-1)S_x^2}{\sigma^2} + \frac{(m-1)S_y^2}{\sigma^2} \sim \chi^2_{n+m-2}$

$\frac{(n-1)S_x^2}{\sigma^2} \sim \chi^2_{n-1}; \frac{(m-1)S_y^2}{\sigma^2} \sim \chi^2_{m-1}$

$$T = \frac{Z}{\sqrt{\frac{U}{n+m-2}}} = \frac{1}{\sigma} \cdot \frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_y)}{\sqrt{\frac{1}{n} + \frac{1}{m}}} = \frac{\bar{x} - \bar{y} - (\mu_x - \mu_y)}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{n+m-2}$$

• UNKNOWN AND UNEQUAL VARIANCES:

① $H_0: \mu_x = \mu_y$ or $\mu_x - \mu_y = 0$ $H_1: \mu_x - \mu_y \neq 0, \mu_x - \mu_y < 0, \mu_x - \mu_y > 0$

② STATISTIC: $T = \frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_y)}{\sqrt{\frac{S_x^2}{n} + \frac{S_y^2}{m}}} \sim t_r$ $r = \frac{\left(\frac{S_x^2}{n} + \frac{S_y^2}{m} \right)^2}{\frac{\left(S_x^2 \right)^2}{n-1} + \frac{\left(S_y^2 \right)^2}{m-1}}$

some proof as ↓

WELCH'S t-interval: $\bar{x} - \bar{y} \pm t_{\alpha/2, r} \cdot \sqrt{\frac{S_x^2}{n} + \frac{S_y^2}{m}}$ use the integer [r]

if large sample size $Z = \frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_y)}{\sqrt{\frac{S_x^2}{n} + \frac{S_y^2}{m}}} \sim N(0, 1)$

Note: T and Z statistics do not do well if
The you need alternative methods.

samplesizes { severely
variances { different
raw data is highly skewed

Ex: Is the mean fastest speed driven by men \neq than by females?

	n	mean	variance
Male	n = 34	$\bar{x} = 105.5$	$s_x^2 = 20.1$
Female	m = 29	$\bar{y} = 90.9$	$s_y^2 = 12.2$

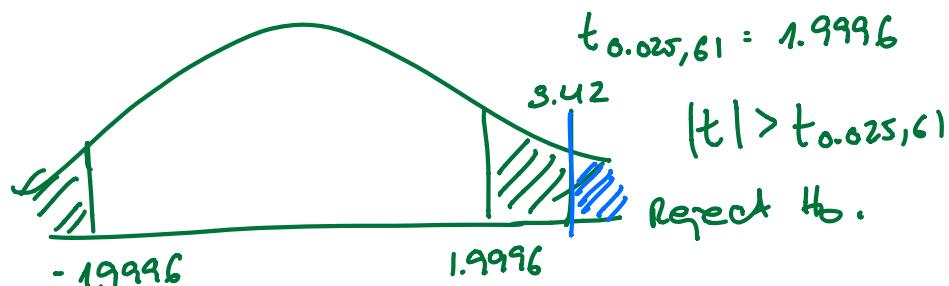
Assume the 2 pops. are normally dist.

$$H_0: \mu_y = \mu_x \quad H_1: \mu_y \neq \mu_x$$

EQUAL VARIANCES
pooled sample standard deviation

$$S_p = \sqrt{\frac{(n-1)s_x^2 + (m-1)s_y^2}{n+m-2}} = \sqrt{\frac{33 \cdot (20.1)^2 + 28 \cdot (12.2)^2}{61-2}} = 16.9$$

$$T = \frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_y)}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} = \frac{(105.5 - 90.9) - 0}{16.9 \sqrt{\frac{1}{34} + \frac{1}{29}}} = 3.42$$



$$P\text{-value} = 2 \times P(T_{61} > 3.42) = 2(0.0006) = 0.0012$$

TEST ABOUT VARIANCES:

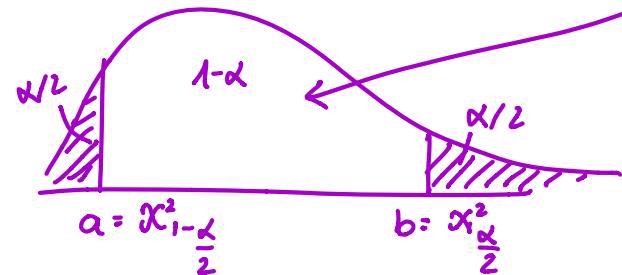
• ONE VARIANCE

$$① H_0: \sigma^2 = \sigma_0^2 \quad H_1: \sigma^2 \neq \sigma_0^2 \quad \sigma^2 > \sigma_0^2 \quad \sigma^2 < \sigma_0^2$$

$$② \text{STATISTIC: } \chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$$

Proof

We know that if $X_i \sim N(\mu, \sigma^2)$ $\frac{(n-1)s^2}{\sigma^2} \sim \chi^2_{n-1}$



NOT EQUAL VARIANCES

$$F = \frac{\left(\frac{s_x^2}{n} + \frac{s_y^2}{m}\right)^2}{\frac{\left(\frac{s_x^2}{n}\right)^2}{n-1} + \frac{\left(\frac{s_y^2}{m}\right)^2}{m-1}} = \frac{\left(\frac{(20.1)^2}{34} + \frac{(12.2)^2}{29}\right)^2}{\left(\frac{(20.1)^2}{34}\right)^2 / 33 + \left(\frac{(12.2)^2}{29}\right)^2 / 28} = 55.5$$

$$t = \dots$$

• 2 VARIANCES

$$\textcircled{1} \quad H_0: \sigma_x^2 = \sigma_y^2 \approx \frac{\sigma_x^2}{\sigma_y^2} = 1 \quad H_1: \sigma_x^2 \neq \sigma_y^2 \quad \sigma_x^2 > \sigma_y^2 \quad \sigma_x^2 < \sigma_y^2$$

$$\textcircled{2} \quad \text{STATISTIC: } F = \frac{s_x^2}{s_y^2}$$

\textcircled{3} Reject H_0 if F is too large

if F is too small

$$F > F_{\frac{\alpha}{2}, (n-1, m-1)}$$

$$F < F_{1-\frac{\alpha}{2}, (n-1, m-1)} = \frac{1}{F_{\frac{\alpha}{2}, (m-1, n-1)}}$$

PROOF:

$$x_i \sim N(\mu_x, \sigma_x^2) \Rightarrow \frac{(n-1)s_x^2}{\sigma_x^2} \sim \chi_{n-1}^2$$

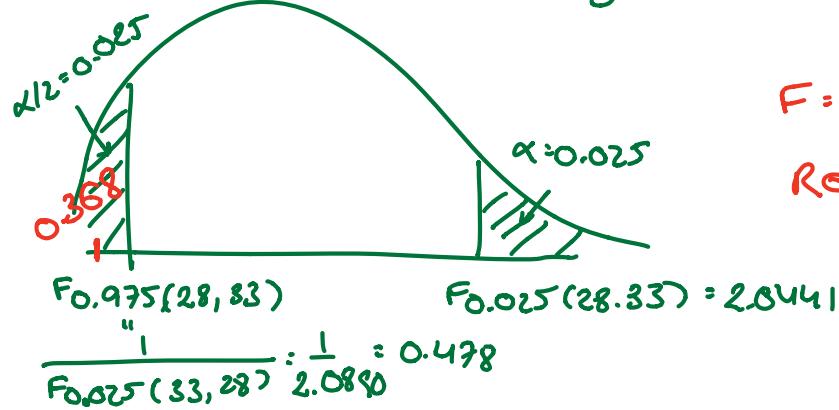
$$y_i \sim N(\mu_y, \sigma_y^2) \Rightarrow \frac{(m-1)s_y^2}{\sigma_y^2} \sim \chi_{m-1}^2$$

$$F = \frac{\frac{(m-1)s_y^2}{\sigma_y^2}}{\frac{n-1}{s_x^2}} = \frac{\sigma_x^2}{\sigma_y^2} \cdot \frac{s_y^2}{s_x^2} \sim F(m-1, n-1)$$

ASSUMPTION: sample pop. or $N \sim (\mu, \sigma)$
random sample

Ex: Is the variance of fastest speed driven by M \neq than by Fem?

$$\textcircled{1} \quad H_0: \sigma_x^2 = \sigma_y^2 \quad H_1: \sigma_x^2 \neq \sigma_y^2 \quad \textcircled{2} \quad F = \frac{(12.2)^2}{(20.1)^2} = 0.368.$$



$$F = 0.368 < 0.478$$

Reject H_0 : two variances are not equal. So to test the difference between means, we should use Welch's t-test not the pooled 2-sample t-test.

- ONE OR TWO-SIDED TEST? Two-sided is more conservative; harder to reject H_0 . Hard to know directionality a priori
- Fail to reject H_0 vs accept H_0 : we don't know the H_0 is true, just that we don't have enough data to reject it. also acknowledges 0.05 random off.
- Don't forget assumptions! They are actually implicitly included in the H_0 .

CATEGORICAL . = NON-PARAMETRIC METHODS.

CHI-SQUARE .

• ONE SAMPLE :

CHI-SQUARE GOODNESS OF FIT STATISTIC : $Q_1 = \sum_{i=1}^2 \frac{(Y_i - np_i)^2}{np_i} = \sum_{i=1}^2 \frac{(\text{OBSERVED} - \text{EXPECTED})^2}{\text{EXPECTED}}$

$$Z^2 = \frac{(Y_1 - np_1)^2}{np_1 p_1} + \frac{(Y_2 - np_2)^2}{np_2 p_2}$$

PROOF : Ex: $Y_1 + Y_2 = n$; $p_1 + p_2 = 1$

Y_1 = successes / females / ... $Y_1 \sim b(n, p_1) \rightarrow E(Y_1) = np_1$
 $\rightarrow \text{Var}(Y_1) = np_1(1-p_1)$

Y_2 = failures / males / ... $Y_2 = n - Y_1 \sim b(n, p_2) = b(n, 1-p_1) \rightarrow E(Y_2) = np_2 = n(1-p_1)$
 $\rightarrow \text{Var}(Y_2) = np_2(1-p_2) =$

when $n \rightarrow \infty$ $\begin{cases} np_1 > 5 \\ n(1-p_1) > 5 \end{cases} \xrightarrow{\text{CLT}} Z = \frac{(Y_1 - np_1)}{\sqrt{np_1(1-p_1)}} \sim N(0, 1) \quad \begin{aligned} &= n(1-p_1)(1-(1-p_1)) \\ &= np_1(1-p_1) \end{aligned}$

$$Z^2 = \frac{(Y_1 - np_1)^2}{np_1(1-p_1)} = Q_1 \sim \chi^2_{(1)}$$

$$Q_1 = \frac{(Y_1 - np_1)^2}{np_1(1-p_1)} = \frac{(Y_1 - np_1)^2}{np_1(1-p_1)} \cdot ((1-p_1) + p_1) = \frac{(Y_1 - np_1)^2 \cdot (1-p_1)}{np_1(1-p_1)} + \frac{(Y_1 - np_1)^2 p_1}{np_1(1-p_1)} =$$

$$\frac{(Y_1 - np_1)^2}{np_1} + \frac{(Y_1 - np_1)^2}{n(1-p_1)} = \frac{(Y_1 - np_1)^2}{np_1} + \frac{(Y_1 - np_1)^2}{np_2} = \frac{(Y_1 - np_1)^2}{np_1} + \frac{(n - Y_1 - n(1-p_2))^2}{np_2} = \frac{(Y_1 - np_1)^2}{np_1} + \frac{(-Y_2 - np_2)^2}{np_2}$$

$$= \frac{(Y_1 - np_1)^2}{np_1} + \frac{(Y_2 - np_2)^2}{np_2}$$

$\rightarrow Q_1 \sim \chi^2_{(1)}$: 1 df because if there is only 1 count Y_1 , \rightarrow if you know Y_1 , you know $Y_2 = n - Y_1$

- we need CLT for proof so $\rightarrow np_1 > 5$; $n(1-p_1) > 5$

- Q_1 is large if Observed \gg Expected. Repeat to when Q_1 is large \rightarrow depends on $\chi^2_{(1)}$

• for small samples ≤ 5 } use FISCHER'S EXACT TEST

• very unequally distributed values among cells } CHI-SQUARE is not adequate because it provides an approximation to the χ^2 that depends on a large sample size.

Ex: Penn state pop. 60% ♀, 40% ♂. Sample $n=100$. $Y_1 = \text{♀} = 53$ $Y_2 = \text{♂} = 47$.

Is the sample random and representative of the population?

$H_0: p_F = 0.60$ $H_1: p_F \neq 0.60$. $\alpha = 0.05$. $\chi^2_{0.05,1} = 3.84$ Reject H_0 if $Q_1 \geq \chi^2_{0.05,1}$

$$\chi^2 = \frac{(Y_1 - np_1)^2}{np_1 q_1} = \frac{(Y_1 - np_1)^2}{np_1} + \frac{(Y_2 - np_2)^2}{np_2} = Q_1 = \frac{(53-60)^2}{60} + \frac{(47-40)^2}{40} = 2.04.$$

$$Q_1 = 2.04 < \chi^2_{0.05,1} = 3.84 \Rightarrow \text{fail to reject } H_0.$$

Relationship χ^2 and z : you can do a z test and square and it's the same!

$$z = \frac{0.53 - 0.60}{\sqrt{\frac{(0.60)(0.40)}{100}}} = -1.428 \Rightarrow z^2 = (-1.428)^2 = 2.04. |z| > 1.96 \Rightarrow Q_1 > (1.96)^2 > 3.84$$

$$P = 2 \times P(z > 1.428) = 2(0.0766) = 0.1532 \Rightarrow P = P(\chi^2_{1,1} > 2.04) = 0.1532$$

- K SAMPLES : K CATEGORIES: $Q_{K-1} = \sum_{i=1}^K \frac{(Y_i - np_i)^2}{np_i} \sim \chi^2_{(K-1)}$

CATEGORIES	1	2	...	K-1	K
OBSERVED	Y_1	Y_2		Y_{K-1}	$n - Y_1 - Y_2 - \dots - Y_{K-1}$
EXPECTED	np_1	np_2		np_{K-1}	np_K

CONTINGENCY TABLES

i = n° of row categories $i = 1 \dots h$

j = n° of col. categories $j = 1 \dots K$

$\hat{p}_{ij} = Y_{ij}/n_i$ proportion of ith sample in jth category

$P(A_i \cap B_j)$ prob. that a randomly selected obs falls in $A_i \cap B_j$ cell

$$\hat{p}_j = \sum_{i=1}^h \hat{p}_{ij}$$

$(Y_{1j} + Y_{2j}) / n_1 + n_2$ proportion in jth category

$p_{.j} = P(B_j)$ prob that randomly selected obs falls in B_j column

$$n_i = \sum_{j=1}^K \hat{p}_{ij}$$

$\sum_{j=1}^K Y_{ij}$ n° in ith sample

$p_{i.} = P(A_i)$ prob that randomly selected obs falls in A_i row

		$j=1$	$j=2$...	Total
		$y_{11} (p_{11})$	$y_{12} (p_{12})$		$n_1 = \sum_j y_{1j} = p_{1.}$
		$y_{21} (p_{21})$	$y_{22} (p_{22})$		$n_2 = \sum_j y_{2j} = p_{2.}$
		:	:	:	:
Total		$y_{11} + y_{21}$ (\hat{p}_1)	$y_{12} + y_{22}$ (\hat{p}_2)	...	$n_1 + n_2 = n$ ($\hat{p}_{1.}$) ($\hat{p}_{2.}$)

- TEST FOR HOMOGENEITY : are the proportions of β in different groups of $A = \alpha$?

$$H_0: P_{A_1} = P_{A_2}^1; P_{A_2} = P_{A_2}^2 -$$

\neq

CHI-SQUARE TEST STATISTIC: $Q = \sum_{i=1}^n \sum_{j=1}^k \frac{(y_{ij} - n_i \hat{p}_{ij})^2}{n_i \hat{p}_{ij}} \sim \chi^2_{(n-1)(k-1)}$

\uparrow
 $n(k-1)-(k-1)$

Ex: Do residents apply ↑ unnecessary blood transfusions?

Residents $n = 71$; Attending physician $n = 49$.

(%) % of raw total P_{ij}/n_i

Physician	frequent	occasional	Rarely	Never	Total
Attending	2 (4.1%) $2/49 = 0.041$	3 (6.1%)	31 (63.3%)	13 (26.5%)	49
Resident	15 (21.1%)	28 (39.4%)	23 (32.4%)	5 (7.0%)	71
Total	17	31	54	18	120

$$H_0: p_{RF} = p_{AF} \text{ AND } p_{RO} = p_{AO} \dots \quad H_1: p_{RF} \neq p_{AF} \text{ OR } p_{RO} \neq p_{AO} \dots$$

Physician	frequent	occasional	Rarely	Never	Total
Attending	$\frac{17}{120} \cdot 49 = 6.94$	$\frac{31}{120} \cdot 49 = 12.65$	$\frac{54}{120} \cdot 49 = 22.29$	7.35	49
Resident	10.058	18.342	31.95	10.65	71
Total	17	31	54	18	120

$$Q = \frac{\text{observed} - \text{expected}}{\text{expected}}^2 = \frac{(2-6.94)^2}{6.94} + \dots + \frac{(5-10.65)^2}{10.65} = 31.88$$

$$\chi^2_{(4-1)(2-1)} = 3 \cdot \chi^2_3 = 7.815 \quad Q > \chi^2_3 \Rightarrow 31.88 > 7.815 \rightarrow \text{reject } H_0$$

Distribution of transfusions is different among physicians and residents

- TEST FOR INDEPENDENCE: are A and B independent from each other?

$$H_0: P(A_i \cap B_j) = P(A_i) \times P(B_j)$$

CHI-SQUARE TEST STATISTIC: $Q = \sum_{i=1}^n \sum_{j=1}^k \frac{(y_{ij} - \frac{y_{i \cdot} \cdot y_{\cdot j}}{n})^2}{\frac{y_{i \cdot} \cdot y_{\cdot j}}{n}} \sim \chi^2_{(n-1)(k-1)}$

\uparrow
 $n(k-1)-(k-1)$

Ex: Age and desire to ride bike independent?

Observed

		Age			
		18-24	25-34	35-49	
wants to ride bike	Yes	60	54		201
	No	40	44		194
	Total	100	98		395

expected

		Age			
		18-24	25-34	35-49	
wants to ride bike	Yes	$\frac{100 \cdot 201}{395} = 50.98$	$\frac{98 \cdot 201}{395} = 49.98$		201
	No				194
	Total	100	98		395

$$Q = \frac{(60-50.98)^2}{50.98} + \dots = 8.006 > \chi^2_3 = 7.815 \quad \text{Reject } H_0 \rightarrow \text{Riding bike depends on age!}$$

MCMENAMAR TEST: for 2×2 contingency table w/ dichotomous trait.
or paired data.

	B^+	B^-	
A^+	a	b	$a+b$
A^-	c	d	$c+d$
	$a+c$	$b+d$	n

It tests consistency in responses across 2 vbls (NOT independence)
Do subjects change from \oplus to \ominus randomly?
we don't care about the subject who don't change (a)

$$H_0: p_c = p_b \quad H_1: p_c \neq p_b$$

$$\text{TEST STATISTIC: } \chi^2 = \frac{(b-c)^2}{b+c}$$

MANN-WHITNEY U TEST = WILCOXON RANK SUM TEST.

H_0 : 2 populations are equal; H_1 : 2 populations are not equal.

MANN-WHITNEY TEST STATISTIC: $U_1 = n_1 n_2 + \frac{n_1(n_1+1)}{n_2} - R_1$

We always use the smaller of \leftarrow

Range of U : $0 - n_1 n_2$

↓ ↓

R_1 : sum of ranks in group 1

$U_2 = n_1 n_2 + \frac{n_2(n_2+1)}{n_1} - R_2$

R_2 : sum of ranks in group 2

complete separation unlikely
group 1 or 2 → separation.

Supports H_1 , Supports H_0

Reject H_0 if $U \leq$ critical value

Ex:

		Total Sample (Ordered Smallest to Largest)		Ranks	
Placebo	New Drug	Placebo	New Drug	Placebo	New Drug
7	3		1		1
5	6		2		2
6	4		3		3
4	2	4	4	4.5	4.5
12	1	5		6	
		6	6	7.5	7.5
		7		9	
		12		10	

37 | 18

$$U_1 = 5 \cdot 5 \cdot \frac{5 \cdot 6}{2} - 37 = 3$$

$$U_2 = 5 \cdot 5 \cdot \frac{5 \cdot 6}{2} - 18 = 22$$

for $\alpha = 0.05$ and $n_1 = 5, n_2 = 5$

$$U = 2 \quad U_1 > U \Rightarrow 3 > 2$$

Fail to reject $H_0 \Rightarrow$ pop's are =

WILCOXON SIGNED RANK TEST

H_0 : median difference = 0 H_1 : median difference $\neq 0$ or > 0 or < 0

WILCOXON TEST STATISTIC: $w_+ : \Sigma$ of $+$ ranks

use the smallest $\leftarrow w_- : \Sigma$ of $-$ ranks

If H_0 is true, we expect w_+ and w_- would be similar.

Reject H_0 if $w \leq$ critical value

Child	Before Treatment	After 1 Week of Treatment	Difference (Before-After)	Ordered Absolute Values of Difference Scores		
				Ranks	Signed Ranks	
1	85	75	10	-5	1	-1
2	70	50	20	10	3	3
3	40	50	-10	-10	3	-3
4	65	40	25	10	3	3
5	80	20	60	15	5	5
6	75	65	10	20	6	6
7	55	40	15	25	7	7
8	20	25	-5	60	8	8

$$w_+ = 34$$

Critical value of w for $n=8$ $\alpha=0.05 = 6$

$4 \leq 6 \Rightarrow$ reject $H_0 \Rightarrow$ There is a change

KRUSKAL-WALLIS TEST \approx 1-way ANOVA w/ rank ≥ 2 groups

H_0 : All the pop. medians are equal H_1 : the pop. medians \neq

TEST STATISTIC: $H = \frac{12}{N(N+1)} \sum_{j=1}^K \frac{R_j^2}{n_j} - 3(N+1)$

Reject H_0 if $H >$ critical value

Ex

			Total Sample (Ordered Smallest to Largest)			Ranks		
5% Protein	10% Protein	15% Protein	5% Protein	10% Protein	15% Protein	5% Protein	10% Protein	15% Protein
3.1	3.8	4.0	2.6			1		
2.6	4.1	5.5	2.9	2.9		2.5	2.5	
2.9	2.9	5.0	3.1			4		
	3.4	4.8		3.4			5	
	4.2			3.8			6	
					4.0			7
					4.1			8
					4.2			9
					4.8			10
					5.0			11
					5.5			12

SUM

| 7.5 | 30.5 | 40 |

$$H = \frac{12}{12 \cdot 13} \left[\frac{(7.5)^2}{3} + \frac{(30.5)^2}{5} + \frac{40^2}{4} \right] - 3 \cdot 13$$

$$= 7.52$$

$$H \text{ for } n_1 = 3, n_2 = 5, n_3 = 4$$

$$\alpha = 0.05$$

$$\text{is } 5.656.$$

$$7.52 > 5.656 \Rightarrow \text{groups } \neq$$

If there are 3 or more comparison groups and/or ≥ 5 obs. per group

$H \rightarrow \chi^2_{K-1}$ so the critical value can be derived from the χ^2 table

FRIEDMAN TEST \approx 1-way ANOVA w/ Rank and repeated measures

H_0 :

TEST STATISTIC $Q = \frac{12n}{K(K+1)} \sum_{j=1}^K \left(\bar{r}_{\cdot j} - \frac{K+1}{2} \right)^2$

where $\bar{r}_{\cdot j} = \frac{\sum_{i=1}^n r_{ij}}{n}$

n : rows / blocks

K : columns / treatments

DERIVE α , β , power

Ex: Random sample $X_i \sim N(\mu, \sigma^2 = 10)$ $n=25$. $H_0: \mu = 170$ $H_1: \mu > 170$.
If the engineer decides $\bar{X} \geq 172$ cutoff.
what is the prob of Type I error?

$$z = \frac{\bar{X} - \mu}{\sqrt{\frac{\sigma^2}{n}}} = \frac{172 - 170}{\sqrt{\frac{10}{25}}} = \frac{2}{\sqrt{2}} = \sqrt{2}$$

$$P(\bar{X} \geq 172; \mu = 170) = P(z \geq \sqrt{2}) = 0.1587$$

what if the true pop mean $\mu = 173$?

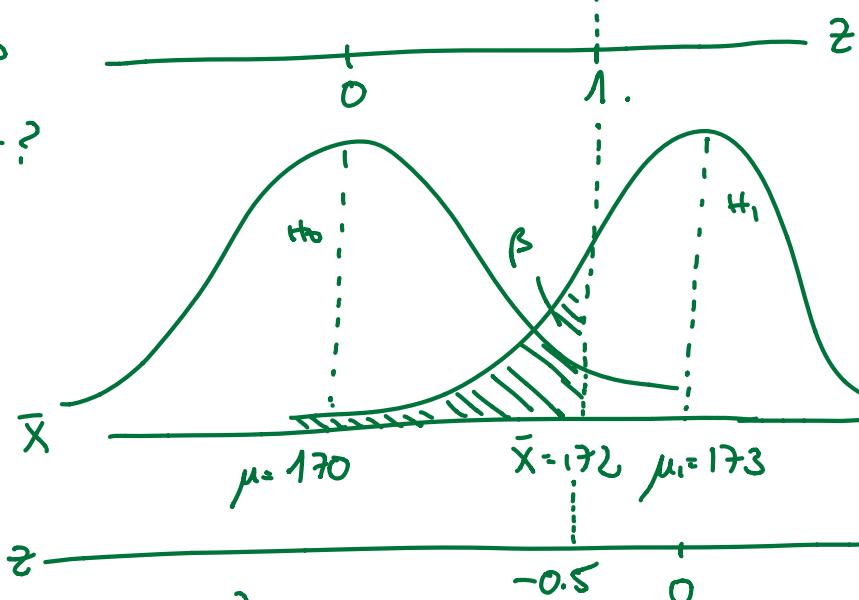
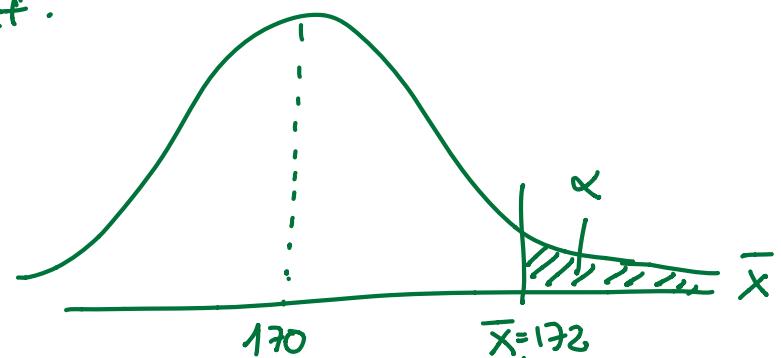
what is the prob of Type II error?
Fail to reject H_0 when H_1 is true

$$z = \frac{\bar{X} - \mu}{\sqrt{\frac{\sigma^2}{n}}} = \frac{172 - 173}{\sqrt{\frac{10}{25}}} = -\frac{1}{\sqrt{2}} = -0.5$$

$$P(\bar{X} \leq 172; \mu = 173) = P(z \leq -0.5) = 0.3085$$

what is the prob the right decision is made (H_1 is accepted when it is true)?

$$P(\bar{X} \geq 172; \mu = 173) = P(z \geq -0.5) = 0.6915 = 1 - 0.3085$$



	$D^+_{H_0}$	$D^-_{H_1}$	
$D^+_{H_0}$	a	b	$T^+ = a+b$
$D^-_{H_1}$	$D^+ \cap T^+$	$D^- \cap T^+$	
$D^-_{H_1}$	c	d	$T^- = c+d$
$D^+_{H_0}$	$D^+ \cap T^-$	$D^- \cap T^-$	
	$a+c$	$b+d$	

$$\text{SE} : \frac{D^+ \cap T^+}{D^+} = \frac{\alpha}{\alpha + c}$$

$$\text{SP} : \frac{D^- \cap T^-}{D^-} = \frac{\beta}{b+d}$$

$$1 - \text{SE} : \alpha : P(H_1 | H_0) = \frac{c}{\alpha + c} = \text{False Neg} \ominus$$

$$1 - \text{SP} : \beta = P(H_0 | H_1) : \frac{b}{b+d} = \text{False} \oplus$$

POWER FUNCTIONS

Power: the probability of making the correct decision if the H_0 is true \rightarrow reject the H_0 .

Ex: $X_i \sim N(\mu, \sigma^2 = 16)$ $n = 16$ students. $\alpha = 0.05$.

- ① $H_0: \mu = 100; H_1: \mu > 100$
- ② Z statistic = 1.645
- ③ Reject H_0 if $Z \geq 1.645$ or $\mu \geq 106.58$

$$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} \Rightarrow \bar{x} = \mu + Z \cdot \frac{\sigma}{\sqrt{n}} = 100 + 1.645 \cdot \left(\frac{16}{\sqrt{16}} \right) = 106.58$$

True pop. $\mu = 108$

$$\text{Power} = P(\bar{x} \geq 106.58 \mid \mu = 108) = P\left(Z \geq \frac{106.58 - 108}{16 / \sqrt{16}}\right) = P(Z \geq -0.36) =$$

$$= 1 - P(Z < -0.36) = 1 - \Phi(-0.36) = 0.6406$$

64% chance of rejecting H_0 when H_1 is true

True pop. $\mu = 116$

$$\text{Power} = P(\bar{x} \geq 106.58 \mid \mu = 116) = P(Z \geq -2.36) = 0.9909$$

Power \uparrow the further the actual μ moves away from the H_0 .

$$\text{Power} = 1 - \Phi(z) \text{ where } z = \frac{106.58 - \mu}{16 / \sqrt{16}}$$

POWER FUNCTION:

$$K(\mu) = 1 - \Phi\left(\frac{\text{CRITICAL VALUE} - \mu}{\sigma / \sqrt{n}}\right)$$

The power will depend on the value of μ

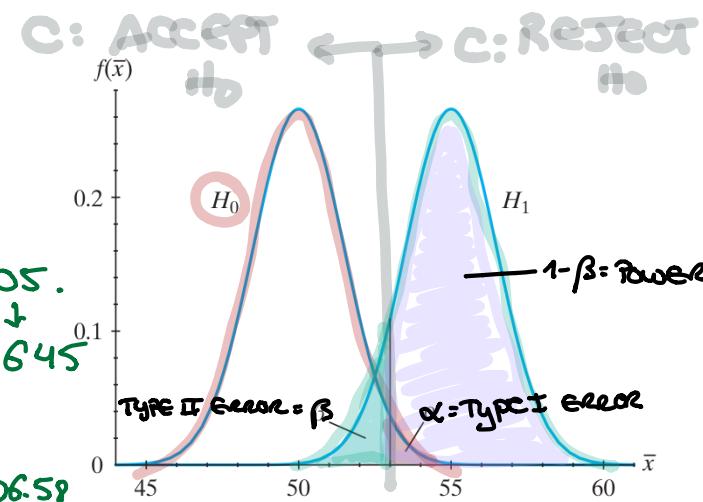
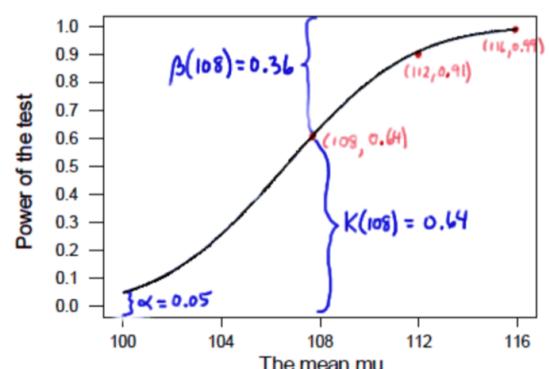
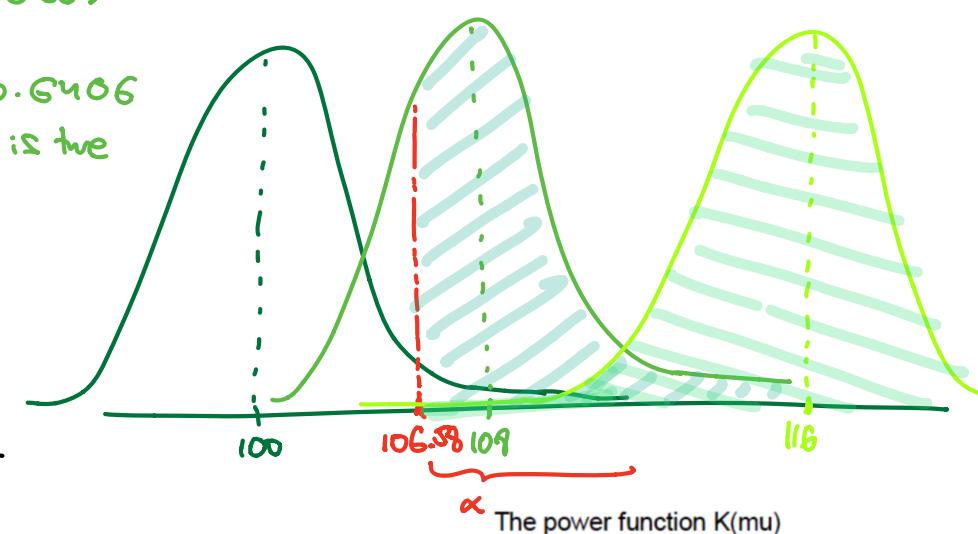


Figure 8.1-1 pdf of \bar{X} under H_0 and H_1

H_0 Distribution: centered around 0.

H_1 Distribution: not "

\hookrightarrow NON CENTRALITY DISTRIBUTION.



Add sample size calculation??

BEST CRITICAL REGIONS

Recall:

Likelihood function: $L(\theta)$ (for 1 parameter)

Random sample: $X_1, X_2 \dots X_n$. X_i is independent bc it's random sample

The pdf for each X_i is $f(x_i; \theta)$

Joint probability mass function of X_i independent vbles

$$L(\theta) = P(X_1 = x_1, X_2 = x_2 \dots X_n = x_n) = f(x_1; \theta) \cdot f(x_2; \theta) \dots f(x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$$

COMPOSITE HYPOTHESIS: A hypothesis that does not uniquely specify the distribution of the pop. from which the sample is taken.

$$\text{Ex: } X_i \sim \exp(\theta), H: \theta > 2 \quad f(x) = \frac{1}{\theta} e^{-\frac{x}{\theta}} \Rightarrow \frac{1}{\theta} e^{-\frac{x}{3}} \text{ or } \frac{1}{\theta} e^{-\frac{x}{22}} \dots \text{etc}$$

$$X_i \sim N(\mu, \sigma^2), H: \mu = 12 \quad f(x) = \frac{1}{\sigma\sqrt{2\pi}} \cdot \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \cdot \exp\left[-\frac{(x-12)^2}{2\sigma^2}\right] \begin{array}{l} \mu \text{ defined} \\ \sigma^2 \text{ not defined.} \end{array}$$

BEST CRITICAL REGION:

$$H_0: \theta = \theta_0 \quad H_1: \theta \neq \theta_0$$

$$\alpha = P(C; \theta_0)$$

C and D are critical regions of size $\alpha \rightarrow \alpha = P(D; \theta_0)$

C is the best critical region of size α if $P(C; \theta_0) \geq P(D; \theta_0)$

that is if the power of C is at least as great as the power of every other critical region D of size α . = the power of the test at $\theta = \theta_0$ is largest among all possible hypothesis tests.

• NEYMAN-PEARSON LEMMA H_0 : simple H ; H_1 : simple H

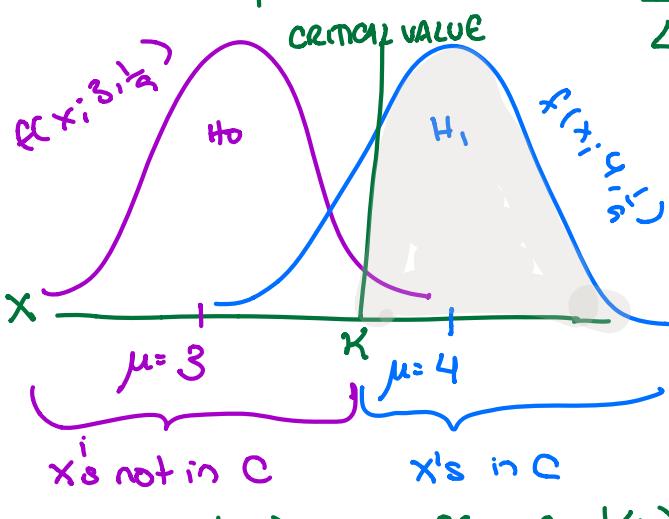
$X_i \sim \text{prob dist } (\Theta)$. X : random sample.

If C is a critical region of size α and K is a constant such that:

$$\left. \begin{array}{l} \frac{L(\theta_0)}{L(\theta_\alpha)} \leq K \text{ inside critical region } C \\ \frac{L(\theta_0)}{L(\theta_\alpha)} \geq K \text{ outside critical region } C \end{array} \right\} \begin{array}{l} \text{then } C \text{ is the best critical region} \\ \text{for testing the} \\ \text{simple } H_0 : \theta = \theta_0 \text{ against the} \\ \text{simple } H_1 : \theta \neq \theta_0 \end{array}$$

Ex: X is one data point from $X_i \sim N(\mu, \sigma^2 = \frac{1}{9})$ $H_0: \mu = 3$ vs $H_1: \mu = 4$

The most powerful test $\Rightarrow \frac{L(\mu_0)}{L(\mu_\alpha)} = \frac{L(3)}{L(4)} = \frac{f(x; 3, 1/9)}{f(x; 4, 1/9)}$



$$\frac{f(x; 3, 1/9)}{f(x; 4, 1/9)} \text{ large} \quad \frac{f(x; 3, 1/9)}{f(x; 4, 1/9)} \text{ small}$$

we reject $H_0: \mu = 3$
if our observed x is large
in favor of $H_1: \mu = 4$,
 \approx if x falls in critical region C

small for sample points x inside critical region $\leq K$
large for sample points x outside critical region $> K$

Ex: X is a single observation from $X_i \Rightarrow p.d.f. = f(x) \propto \theta x^{\theta-1} \quad 0 < x < 1$
Best critical region with $\alpha = 0.05$ to test $H_0: \theta = 3$ $H_1: \theta = 2 \rightarrow$ simple hypotheses

$$\frac{L(\theta_0)}{L(\theta_\alpha)} = \frac{3x^{3-1}}{2x^{2-1}} = \frac{3x^2}{2x} = \frac{3}{2}x \leq K \Rightarrow x < \frac{2}{3}K = K^* \Rightarrow$$

rejection region for the most powerful test

Find K^* ! we want $\alpha = P(\text{Type I Error}) = P(\text{reject } H_0 \text{ when } H_0 \text{ is true}) = 0.05$

$$\alpha = P(X < K^* \text{ when } \theta = 3) = \int_0^{K^*} 3x^2 dx = 0.05 \Rightarrow 3 \cdot \frac{1}{3}x^3 \Big|_0^{K^*} = (K^*)^3 = 0.05$$

$$K^* = (0.05)^{1/3} = 0.368$$

Neyman-Pearson lemma says the rejection region for the most powerful test to test $H_0: \theta = 3$ $H_1: \theta = 2$ is $x < 0.368$.

Ex: x_i : random sample $\sim N(\mu, \sigma^2 = 16)$ test critical region for $\{n = 16, \alpha = 0.05\}$
 $L(0)$ for $\sigma \sim N(\mu, \sigma^2) = (2\pi \sigma^2)^{-\frac{n}{2}} \cdot \exp\left[\frac{-1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu_0)^2\right]$

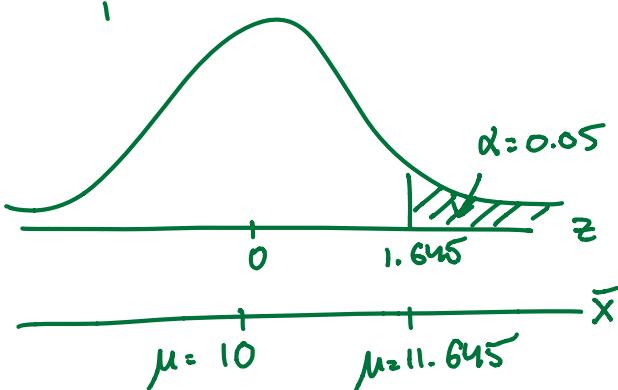
sample: $H_0: \mu = 10$;
 sample: $H_1: \mu = 15$

$$\begin{aligned} \frac{L(10)}{L(15)} &= \frac{(32\pi)^{-16/2} \cdot \exp\left[-\frac{1}{32} \sum_{i=1}^{16} (x_i - 10)^2\right]}{(32\pi)^{-16/2} \cdot \exp\left[-\frac{1}{32} \sum_{i=1}^{16} (x_i - 15)^2\right]} = e^{\left[-\frac{1}{32} \sum_{i=1}^{16} (x_i - 10)^2\right] - \left[-\frac{1}{32} \sum_{i=1}^{16} (x_i - 15)^2\right]} \\ &= e^{-\frac{1}{32} \left[\sum_{i=1}^{16} (x_i - 10)^2 - \sum_{i=1}^{16} (x_i - 15)^2 \right]} = e^{-\frac{1}{32} \left[\sum_{i=1}^{16} x_i^2 - 2 \cdot 10 \sum_{i=1}^{16} x_i + \sum_{i=1}^{16} 100 - (\sum_{i=1}^{16} x_i^2 - 2 \cdot 15 \cdot \sum_{i=1}^{16} x_i - 225^2) \right]} \\ &= e^{-\frac{1}{32} [10 \sum_{i=1}^{16} x_i + 1600 - 3600]} = e^{-\frac{1}{32} (10 \sum_{i=1}^{16} x_i - 2000)} \leq K \Rightarrow -10 \sum_{i=1}^{16} x_i + 2000 \leq \ln K \cdot 32 \\ &\Rightarrow -\frac{10}{n} \sum_{i=1}^{16} x_i \leq \frac{\ln K - 2000}{160} \Rightarrow \frac{\sum_{i=1}^{16} x_i}{16} \geq \frac{\ln K - 2000}{160} \Rightarrow \bar{x} \geq k^* \end{aligned}$$

Neyman-Pearson: Critical region $C = \{(\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n); \bar{x} \geq k^*\}$ where k^* is selected so that the size of the critical region $C = \alpha = 0.05$?
 Under the H_0 the sample mean has distribution $\bar{x} = 10$, $\sqrt{8\delta} = \sqrt{\frac{\sigma}{n}} = \frac{4}{4} = 1 \Rightarrow \bar{x} \sim N(10, 1)$

$$P[\bar{x} \geq k^*, \mu = 10] = P\left[\frac{\bar{x} - 10}{1} \geq \frac{k^* - 10}{1}; \mu = 10\right] = 1 - P\left(Z \leq \frac{k^* - 10}{1}\right) = 0.05$$

$$\frac{k^* - 10}{1} = 1.645 \Rightarrow k^* = 11.645$$



The rejection region for the most powerful test to test $H_0: \mu = 10$; $H_1: \mu = 15$ is $\bar{x} \geq 11.645$

The power of this test when $\mu = 15$

$$P(\bar{x} > 11.645; \mu = 15) = P\left(Z > \frac{11.645 - 15}{\sqrt{16}}\right) = P(Z > -3.36)$$

this is the max power we'd get out of any test

- UNIFORMLY MOST POWERFUL TESTS = UMP $H_0: \text{Simple H. } H_1: \text{composite H.}$
A test defined by the critical region C for a size α is the UMP test if it is the most powerful test against each simple alternative in the H_1 .
The critical region $C = \text{UMP critical region for size } \alpha.$

Ex: x_i : random sample $\sim N(\mu, \sigma^2 = 16)$ best critical region for $\left\{ n = 16, \alpha = 0.05 \right.$
 $L(\theta)$ for $\theta \sim N(\mu, \sigma^2) = (2\pi \sigma_0^2)^{-n/2} \cdot \exp \left[-\frac{1}{2\sigma_0^2} \sum_{i=1}^n (x_i - \mu_0)^2 \right]$ $\left. \begin{array}{l} \text{Simple: } H_0: \mu = 10 \\ \text{Composite: } H_1: \mu > 10 \end{array} \right\}$

$$\frac{L(10)}{L(\mu_0)} = \frac{(32\pi)^{-16/2} \cdot \exp \left[-\frac{1}{32} \sum_{i=1}^{16} (x_i - 10)^2 \right]}{(32\pi)^{-16/2} \cdot \exp \left[-\frac{1}{32} \sum_{i=1}^{16} (x_i - \mu_0)^2 \right]} = e^{-\frac{1}{32} \left[\sum_{i=1}^{16} (x_i - 10)^2 - \sum_{i=1}^{16} (x_i - \mu_0)^2 \right]} \leq K$$

take ln on both sides

$$C = e^{-\frac{1}{32} \left[\sum_{i=1}^{16} x_i^2 - 2 \cdot 10 \sum x_i + 16(10)^2 - \left(\sum x_i^2 - 2 \cdot \mu_0 \sum x_i - 16\mu_0^2 \right) \right]} = e^{-\frac{1}{32} (-2 \cdot \sum x_i (10 - \mu_0) + 16(\mu_0^2 - 10^2))} \leq K$$

more constants to write

$$-2 \cdot \sum x_i (10 - \mu_0) + 16(\mu_0^2 - 10^2) \leq 32 \ln(K) \Rightarrow \frac{\sum x_i}{n} \leq \frac{82 \cdot \ln(K) - 16(\mu_0^2 - 10^2)}{2 \cdot (10 - \mu_0)}$$

$$\bar{x} \geq -\frac{82 \cdot \ln(K) - 16(\mu_0^2 - 10^2)}{16 \cdot 2 \cdot (10 - \mu_0)} = K^* \Rightarrow \frac{L(10)}{L(\mu_0)} \leq K \Rightarrow \bar{x} \geq K^*$$

Neyman Pearson: The best critical region to test $H_0: \mu = 10$ against each $H_1: \mu \neq \mu_1$

where $\mu_1 > 10$ is $C = [(x_1, x_2, \dots, x_n); \bar{x} \geq K^*]$

where K^* is selected so that the size of the critical region $C = \alpha = P(\text{Type I Err})$

$$\alpha = P(\bar{x} > K^*) \text{ when } \mu = 10$$

LIKELIHOOD RATIOS

H_0 : simple/composite H H_1 : simple composite H

Ω : total possible parameter space of Θ in both the H_0 and H_1 ,

$H_0: \Theta \in \omega$ H_0 where ω is a subset of Ω

$H_1: \Theta \in \omega'$ H_1 where ω' is the complement of ω in space Ω .

$L(\hat{\omega})$: maximum likelihood function with respect to Θ where Θ is in entire space Ω

$L(\hat{\omega}')$: maximum likelihood function with respect to Θ where Θ is in null space ω to H_0 .

LIKELIHOOD RATIO: $\lambda = \frac{L(\hat{\omega}')}{L(\hat{\omega})}$

In general, a LR is simply the likelihood of one hypothesis relative to that of another, for an observed set of data. This is, calculated as the probability of the data given one hypothesis, often defined by a particular value of a parameter, divided by the probability of the same data given the second, competing hypothesis.

Critical region for THE LIKELIHOOD RATIO TEST to test $H_0: \Theta \in \omega$ against

is the set of sample points for which $\lambda = \frac{L(\hat{\omega}')}{L(\hat{\omega})} \leq K$
where $0 < K < 1$ and

K is selected so that the test has a desired significance α

- H_0 : simple H H_1 : composite H

Ex: Volume of honey is approx 10 oz. $x_i \sim N(\mu=10, \sigma^2=2)$
likelihood ratio test for testing $H_0: \mu=10$ vs. $H_1: \mu \neq 10$ at $\alpha=0.05$?

$\Omega = \{\mu : -\infty < \mu < \infty\}$; $H_0: \omega = \{\mu=10\}$

$L(\hat{\omega}) = L(10)$ when H_0 is true.

$L(\hat{\omega}')$ \Rightarrow what value of μ maximizes $L(\mu) \Rightarrow$ MLE of $\mu \Rightarrow \mu \in \hat{\mu} = \bar{x} \Rightarrow$

$$L(\hat{\omega}') = L(\bar{x})$$

$$L(\theta) \text{ for } \theta \sim N(\mu, \sigma^2) = (2\pi\sigma^2)^{-n/2} \cdot e^{-\frac{1}{2\sigma^2} \sum (x_i - \mu)^2}$$

$$\lambda = \frac{L(\hat{\omega}')}{L(\hat{\omega})} = \frac{L(10)}{L(\bar{x})} = \frac{\frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2\sigma^2} \sum (x_i - 10)^2}}{\frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2\sigma^2} \sum (x_i - \bar{x})^2}}$$

$$= \frac{\exp\left[-\frac{1}{4} \sum (x_i - \bar{x})^2\right]}{\exp\left[-\frac{1}{4} \sum (x_i - \bar{x})^2\right]} \Rightarrow \lambda = e^{-\frac{n}{4} (\bar{x} - 10)^2} \leq K$$

reject the H_0
when λ is
small (less
than some K)

$$-\frac{n}{4}(\bar{x}-10)^2 \leq \ln K \Rightarrow (\bar{x}-10)^2 \geq -\ln K \cdot \frac{4}{n} \Rightarrow \frac{|\bar{x}-10|}{\sqrt{\frac{n}{4}}} \geq \sqrt{\frac{-4 \cdot \ln K}{n}} = K^*$$

↓

$$z \sim N(0,1) \text{ for } H_0: \mu=10$$

2 ratio tells us to reject $H_0: \mu=10$ in favor of $H_1: \mu \neq 10$ for all sample means for which $\frac{|\bar{x}-10|}{\sqrt{\frac{n}{4}}} \geq 2_{0.025} = 1.96$ is true.
↳ because we set $\alpha = 0.05$

- H_0 : composite H H_1 : composite H

$E_n: X \sim N(\mu, \sigma^2)$
likelihood ratio test for testing $H_0: \mu=\mu_0$ vs. $H_1: \mu \neq \mu_0$ at $\alpha=0.05$?

$\Omega = \{(\mu, \sigma^2): -\infty < \mu < \infty, 0 < \sigma^2 < \infty\}$; $\omega = \{(\mu, \sigma^2): \mu = \mu_0, 0 < \sigma^2 < \infty\}$.

MLE for $\hat{\mu} = \bar{x}$ and $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$

$$\begin{aligned} L(\hat{\Omega}) &= (2\pi \hat{\sigma}^2)^{-\frac{n}{2}} \cdot \exp \left[-\frac{1}{2\hat{\sigma}^2} \sum (x_i - \hat{\mu})^2 \right] = \left[\frac{1}{2\pi \underbrace{\frac{1}{n} \sum (x_i - \bar{x})^2}_{\hat{\sigma}^2}} \right]^{\frac{n}{2}} \cdot \exp \left[-\frac{\sum (x_i - \bar{x})^2}{2 \cdot \frac{1}{n} \sum (x_i - \bar{x})^2} \right] \\ &= \left[\frac{n}{2\pi \sum (x_i - \bar{x})^2} \right]^{\frac{n}{2}} \cdot e^{-\frac{n}{2}} = \left[\frac{n \cdot e^{-1}}{2\pi \sum (x_i - \bar{x})^2} \right]^{\frac{n}{2}} \end{aligned}$$

under the H_0 parameter space ω MLE for $\hat{\mu} = \mu_0$ and $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_0)^2$

$$L(\hat{\omega}) = \left[\frac{n \cdot e^{-1}}{2\pi \sum (x_i - \mu_0)^2} \right]^{\frac{n}{2}}$$

$$\lambda = \frac{L(\hat{\Omega})}{L(\hat{\omega})} = \frac{\left[\frac{n \cdot e^{-1}}{2\pi \sum (x_i - \mu_0)^2} \right]^{\frac{n}{2}}}{\left[\frac{n \cdot e^{-1}}{2\pi \sum (x_i - \bar{x})^2} \right]^{\frac{n}{2}}} = \frac{\frac{2\pi \sum (x_i - \bar{x})^2}{2\pi \sum (x_i - \mu_0)^2}}{\left[\frac{\sum (x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2 + n(\bar{x} - \mu_0)^2} \right]^{\frac{n}{2}}} = \left[\frac{1}{1 + \frac{n(\bar{x} - \mu_0)^2}{\sum (x_i - \bar{x})^2}} \right]^{\frac{n}{2}}$$

$$\frac{1}{1 + \frac{n(\bar{x} - \mu_0)^2}{\sum (x_i - \bar{x})^2}} \leq K^{\frac{n}{2}} \Rightarrow \frac{n(\bar{x} - \mu_0)^2}{\sum (x_i - \bar{x})^2} \cdot \frac{1}{n-1} \geq (n-1)(K^{\frac{-2}{n}} - 1) \quad ?$$

when H_0 is true : $\frac{\sqrt{n}(\bar{X}-\mu_0)}{\sigma} \sim N(0,1)$ $\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sigma^2} \sim \chi^2_{(n-1)}$

$$H_0: T = \frac{\frac{\sqrt{n}(\bar{X}-\mu_0)}{\sigma}}{\sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sigma^2(n-1)}}} = \frac{\frac{n(\bar{X}-\mu_0)^2}{\sigma^2(n-1)}}{\sqrt{\frac{S^2}{n}}} = \frac{\bar{X}-\mu_0}{\frac{S}{\sqrt{n}}} \sim t_{\alpha/2, n-1} = K^*$$

$t_{\alpha/2, n-1}$ is the size α t-test that leaves the prob of committing type I error

① STATISTIC for the LIKELIHOOD RATIO : $-2 \log L_R$

under regularity conditions on the model $f(x|\theta)$.

$L_R: -2 \log R \rightarrow \chi^2_{df} \rightarrow$ difference in the number of parameters
 Number of free parameters specified by Θ in the H_0 space
 - Number of parameters in the entire parameter space.

② Reject H_0 if $-2 \log R > \chi^2_{df, \alpha}$

- $\Theta \in H_0$ parameter space
- $n \rightarrow \infty$

ANOVA: ANALYSIS OF VARIANCE

- ONE-WAY ANOVA: categorical variables

m independent random variables

m normal distributions $\left\{ \begin{array}{l} \text{means: } \mu_1, \mu_2, \dots, \mu_m \\ X_{ini} \sim N(\mu_i, \sigma^2) \end{array} \right.$

- $H_0: \mu_1 = \mu_2 = \dots = \mu_m = \mu$
- $H_1: \text{NOT } H_0 \text{ At least 1 mean differs}$

Table 9.3-1 One-factor random samples

Group	Data: j observations					Means
X _{1:}	X ₁₁	X ₁₂	...	X _{1n₁}		\bar{X}_1
X _{2:}	X ₂₁	X ₂₂	...	X _{2n₂}		\bar{X}_2
⋮	⋮	⋮	⋮	⋮	⋮	⋮
X _{m:}	X _{m1}	X _{m2}	...	X _{mn_m}		\bar{X}_m
					Grand Mean:	$\bar{X}_{..}$

Table 9.3-2 Analysis-of-variance table

Source	Sum of Squares (SS)	Degrees of Freedom	Mean Square (MS)	F Ratio
Treatment	SS(T) BETWEEN	m - 1	MS(T) = $\frac{SS(T)}{m - 1}$	$\frac{MS(T)}{MS(E)}$
Error	SS(E) WITHIN	n - m	MS(E) = $\frac{SS(E)}{n - m}$	
Total	SS(TO)	n - 1		

P-value

$$= P(F_{(m-1, n-m)} > F = \frac{MS(T)}{MS(E)})$$

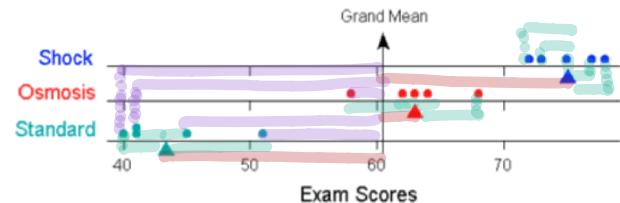
unbiased estimator of σ^2

m: no groups compared

x_{ij} : jth obs of i-th group. $j = 1 \dots n_i$
the no. of obs may vary per group

$$\bar{x}_{i.} = \frac{1}{n_i} \sum_{j=1}^{n_i} x_{ij} \text{ mean of group } i$$

$$\bar{x}_{..} = \frac{1}{n} \sum_{i=1}^m \sum_{j=1}^{n_i} x_{ij} \text{ grand mean of all data points}$$



ASSUMPTIONS

$$\underline{\text{SS}(T)} = \sum_{i=1}^m \sum_{j=1}^{n_i} (\bar{x}_{i.} - \bar{x}_{..})^2 = \sum_{i=1}^m n_i \bar{x}_{i.}^2 - n \bar{x}_{..}^2 \Rightarrow \frac{\text{SS}(T)}{\sigma^2} \sim \chi^2_{(m-1)} \leftarrow H_0 \text{ is true}$$

Difference group mean vs grand mean

$$\underline{\text{SS}(E)} = \sum_{i=1}^m \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{i.})^2 = \text{SS(TO)} - \text{SS}(T) \Rightarrow \frac{\text{SS}(E)}{\sigma^2} \sim \chi^2_{(n-m)} \left\{ \begin{array}{l} \cdot x_{ij} \sim N(\mu_i, \sigma^2) \\ \cdot S_i^2 = \frac{1}{n_i-1} \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{i.})^2 \end{array} \right. \text{sample variance } x_{ij}$$

Difference observations vs group mean

$$\underline{\text{SS(TO)}} = \sum_{i=1}^m \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{..})^2 = \sum_{i=1}^m \sum_{j=1}^{n_i} (x_{ij}^2 - n \bar{x}_{..}^2) = \overline{\text{SS}(E)} + \overline{\text{SS}(T)} \Rightarrow \frac{\text{SS(TO)}}{\sigma^2} \sim \chi^2_{(n-1)}$$

Difference observations vs. grand mean.

$$\text{If } x_{ij} \sim N(\mu_i, \sigma^2) \Rightarrow F = \frac{\frac{\text{SS}(T)}{\sigma^2 / (m-1)}}{\frac{\text{SS}(E)}{\sigma^2 / (n-m)}} = \frac{MST}{MSE} \sim F_{(m-1, n-m)}$$

$$\frac{\text{SS}(T)}{\sigma^2 / (m-1)}$$

$$\frac{\text{SS}(E)}{\sigma^2 / (n-m)}$$

composes the "average" variability between groups to "average" variability within groups

F is a reasonable statistic because

$$E(MSE) = \sigma^2$$

$$\text{under } H_0 (\mu_1 = \mu_2 = \dots = \mu_m) \quad E(MST) = \sigma^2 \Rightarrow \frac{MST}{MSE} \approx \frac{\sigma^2}{\sigma^2} \sim 1$$

$$\text{under } H_1 \quad E(MST) > \sigma^2 \Rightarrow \frac{MST}{MSE} > 1$$

② Reject the H_0 if $F > F_{\alpha(m-1, n-m)}$ or $P = P(F(m-1, n-m) \geq F) \leq \alpha$

ASSUMPTIONS of F-test : { independence
normality
equality of means}

Ex:

Table 9.3-3 Illustrative data			
	Observations	\bar{X}_i	
x_1 :	13 8 9 10		
x_2 :	15 11 13 13		
x_3 :	8 12 7 9		
x_4 :	11 15 10 12		
$\bar{x}_{..}$		11	

Let X_1, X_2, X_3, X_4 be independent random variables that have normal distributions $N(\mu_i, \sigma^2)$, $i = 1, 2, 3, 4$. We shall test

$$H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu$$

against all alternatives on the basis of a random sample of size $n_i = 3$ from each of the four distributions. A critical region of size $\alpha = 0.05$ is given by

$$F = \frac{SS(T)/(4-1)}{SS(E)/(12-4)} \geq 4.07 = F_{0.05}(3, 8).$$

Table 9.3-4 ANOVA table for illustrative data					
Source	Sum of Squares (SS)	Degrees of Freedom	Mean Square (MS)	F Ratio	p-value
Treatment	30	3	30/3	1.6	0.264
Error	50	8	50/8		
Total	80	11			

Note that since $SS(TO) = SS(E) + SS(T)$, only two of the three values need to be calculated directly from the data. Here the computed value of F is

$$\frac{30/3}{50/8} = 1.6 < 4.07,$$

and H_0 is not rejected. The p -value is the probability, under H_0 , of obtaining an F that is at least as large as this computed value of F . It is often given by computer programs.

MULTIPLE COMPARISONS : once you have rejected the H_0 , how do you know which means differ from each other?

- Exploratory analysis : sort means by magnitude, box-whisker plot.
- Testing : Tukey's honest significant difference test, etc --

• TWO-WAY ANOVA

- ONE VALUE PER CELL $\Rightarrow X_{ij} \sim N(\mu_{ij}, \sigma^2)$
- NO REPLICATES

$i = 1, 2, \dots, a \rightarrow \text{row}$
 $j = 1, 2, \dots, b \rightarrow \text{col.}$
 $a \times b = n \text{ comb.}$

Table 9.4-1 Two-way ANOVA table, one observation per cell

Source	Sum of Squares (SS)	Degrees of Freedom	Mean Square (MS)	F
Factor A (row)	SS(A)	$a - 1$	$MS(A) = \frac{SS(A)}{a - 1}$	$\frac{MS(A)}{MS(E)}$
Factor B (column)	SS(B)	$b - 1$	$MS(B) = \frac{SS(B)}{b - 1}$	$\frac{MS(B)}{MS(E)}$
Error	SS(E)	$(a - 1)(b - 1)$	$MS(E) = \frac{SS(E)}{(a - 1)(b - 1)}$	
Total	SS(TO)	$ab - 1$		

$$H_0: \alpha_1 = \alpha_2 = \dots = \alpha_a = 0 \\ H_A: \alpha_1 \neq \alpha_2 \dots \neq \alpha_a$$

\Rightarrow if true

$$\mu_{ij} = \mu + \alpha_i + \beta_j \quad \sum_{i=1}^a \alpha_i = 0 = \sum_{j=1}^b \beta_j = 0$$

ith row effect + jth col effect + overall effect.

$$: \frac{SS(A)}{\sigma^2} \sim \chi^2_{(a-1)} \quad : \frac{SS(B)}{\sigma^2} \sim \chi^2_{(b-1)} \quad : \frac{SS(E)}{\sigma^2} \sim \chi^2_{(a-1)(b-1)}$$

$$SS(TO) = \sum_{i=1}^a \sum_{j=1}^b (x_{ij} - \bar{x}_{..})^2 = SS(A) + SS(B) + SS(E) \sim \chi^2_{(ab-1)}$$

Test $H_A \Rightarrow F = \frac{MS(A)}{MS(E)}$ Reject H_A when F is too large

Test $H_B \Rightarrow F = \frac{MS(B)}{MS(E)}$ Reject H_B when F is too large

Ex: 3 cars each drive w/ 4 \neq brand of gasoline. Do we get the same mileage for each of the 4 brands?

Table 9.4-2 Gas mileage data

Car	Gasoline				\bar{X}_i
	1	2	3	4	
1	26	28	31	31	29
2	24	25	28	27	26
3	25	25	28	26	26
\bar{X}_j	25	26	29	28	27

① $H_B: \beta_1 = \beta_2 = \beta_3 = \beta_4 = 0 \quad H_1: \text{Not } H_B$

② Reject H_B if $F > F_{0.01}(3, 6)$

$$\frac{SS(B)/(4-1)}{SS(E)/[(3-1)(4-1)]} \geq 9.78 = F_{0.01}(3, 6).$$

We have

$$SS(B) = 3[(25-27)^2 + (26-27)^2 + (29-27)^2 + (28-27)^2] = 30; \\ SS(E) = (26-29-25+27)^2 + (24-26-25+27)^2 + \dots \\ + (26-26-28+27)^2 = 4.$$

Hence, the computed F is

$$\frac{30/3}{4/6} = 15 > 9.78 \Rightarrow \text{reject } H_B \rightarrow \text{gasolines } \neq$$

Table 9.4-3 ANOVA table for gas mileage data

Source	Sum of Squares (SS)	Degrees of Freedom	Mean Square (MS)	F	p-value
Row (A)	24	2	12	18	0.003
Column (B)	30	3	10	15	0.003
Error	4	6	2/3		
Total	58	11			

- ≥ 1 value per cell
REPLICATES

$X_{ijk} \sim N(\mu_{ij}, \sigma^2)$ $i = 1, 2, \dots$ $a \rightarrow$ row
independent roles. $j = 1, 2, \dots$ $b \rightarrow$ col
 $K = 1, 2, \dots$ $c \rightarrow$ replicates
 $a \times b \times c = n$ comb.

Table 9.4-4 Two-way ANOVA table, c observations per cell

Source	Sum of Squares (SS)	Degrees of Freedom	Mean Square (MS)	F
Factor A (row)	SS(A)	$a - 1$	$MS(A) = \frac{SS(A)}{a - 1}$	$\frac{MS(A)}{MS(E)}$
Factor B (column)	SS(B)	$b - 1$	$MS(B) = \frac{SS(B)}{b - 1}$	$\frac{MS(B)}{MS(E)}$
Factor AB (interaction)	SS(AB)	$(a - 1)(b - 1)$	$MS(AB) = \frac{SS(AB)}{(a - 1)(b - 1)}$	$\frac{MS(AB)}{MS(E)}$
Error	SS(E)	$ab(c - 1)$	$MS(E) = \frac{SS(E)}{ab(c - 1)}$	
Total	SS(TO)	$abc - 1$		

$$\mu_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij} \rightarrow$$

it's row effect + jth col effect + interaction effect

$$\sum_{i=1}^a \alpha_i = 0 = \sum_{j=1}^b \beta_j = 0 = \sum_{i=1}^a \gamma_{ij} = 0 = \sum_{j=1}^b \gamma_{ij} = 0$$

$$H_0 \begin{cases} H_{A0}: \alpha_1 = \alpha_2 = \dots = \alpha_a = 0 \\ H_{B0}: \beta_1 = \beta_2 = \dots = \beta_b = 0 \\ H_{C0}: \gamma_{ij} = 0 \text{ ? No interaction} \end{cases} \Rightarrow \text{if true}$$

$$\begin{aligned} & : \frac{SS(A)}{\sigma^2} \sim \chi^2_{(a-1)} & \frac{SS(E)}{\sigma^2} \sim \chi^2_{[ab(c-1)]} \\ & : \frac{SS(B)}{\sigma^2} \sim \chi^2_{(b-1)} \\ & : \frac{SS(AB)}{\sigma^2} = \frac{\sum_{k=1}^c (x_{ijk} - \bar{x}_{ij.})^2}{\sigma^2} \sim \chi^2_{(c-1)} \end{aligned}$$

$$SS(TO) = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^c (x_{ijk} - \bar{x}...)^2 = SS(A) + SS(B) + SS(AB) + SS(E) \sim \chi^2_{(abc-1)}$$

$$\text{Test } H_A \Rightarrow F = \frac{MS(A)}{MS(E)} \quad \text{Reject } H_A \text{ when } F > F_{\alpha}[(a-1)(b-1), (ab(c-1))]$$

$$\text{Test } H_B \Rightarrow F = \frac{MS(B)}{MS(E)} \quad \text{Reject } H_B \text{ when } F > F_{\alpha}[(a-1)(b-1), (ab(c-1))]$$

$$\text{Test } H_{AB} \Rightarrow F = \frac{MS(AB)}{MS(E)} \quad \text{Reject } H_{AB} \text{ when } F > F_{\alpha}[(a-1)(b-1), (ab(c-1))]$$

→ don't need to test

→ if there is interaction
↓
there is difference among the means

ASSUMPTIONS OF ANOVA:

- Observations are random samples from the populations
- are drawn from Normally dist. populations
- Variance of all populations is equal.
- Numerator and denominator of F test are independent.