

statistics v.s. probability

probability : know about underlying distribution (population) → make predictions about sample

examples: •  $E$  of a random variable : tells you what the average value of your samples should be

- LLN : if you take lots of samples, then the sample average is likely to be close to the true underlying population average
- CLT : sample average follows a normal distribution

Statistics : start with data sample → try to learn sth about the underlying population that generated that sample.

examples: • Suppose you roll a die 6 times and you get a 1. 4/6 times

Q: do you conclude that the die is biased.

What if you got a 1  $\frac{4000}{6000}$  times? same Q

where do you draw the line?

• Suppose there are 2 drugs for a disease drug 1 cured  $\frac{7}{10}$  people, drug 2 cured  $\frac{9}{10}$

Q: Is drug 2 better?

what if drug 1 cured  $\frac{7000}{10,000}$ , drug 2 cured  $\frac{9000}{10,000}$ ?

In other words: How do you decide if sth you observed is real or just caused by normal random fluctuations?

The answers to these questions are necessarily a bit subjective, but the goal of this class is to be more quantitative about it.

• The main tool that we use for this is probability.

conditional probability      "probability of B given A"

A, B are two events.  $P(B|A)$  = the prob. of B happening if we assume that A happens

examples:  $P(\text{roll a } 2 \mid \text{roll an even \#}) = \frac{1}{3}$

$$P(\text{even \#} \mid 2) = 1$$

Sneak peek: how to answer the Qs from above?

calculate something like  $P(\text{getting a } 1 \mid \text{die is fair})$   
4000/6000 times

if this probability is tiny, then maybe you conclude that your assumption that the die is fair is incorrect

### Bayes Rule

$$\text{formula: } P(B \mid A) = \frac{P(A \mid B) \cdot P(B)}{P(A)}$$
 (gross)

#### example:

Suppose there's a disease and 1% of people are sick. Also, the test for the disease is 99% accurate. Suppose you get a + test what is the probability that you are actually sick?

Know:  $P(+ \mid \text{sick}) = 0.99$

Want to know:  $P(\text{sick} \mid +)$

	-	1,000,000 people
	+	
healthy		
990,000	10,000	

$$P(\text{sick} \mid +) = \frac{\# \text{ sick + people}}{\text{total \# + ppl}}$$

$$= \frac{9900}{9900 + 9900}$$

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

## Order statistics

$X$  = some probability distribution (continuous)

take a bunch of iid samples  $x_1, x_2, \dots, x_n \sim X$   
 ↳ independent and identically distributed

put the samples in order  $y_1 \leq y_2 \leq y_3 \dots \leq y_n$

text book # 6.3.3

exponential distribution formula

$$x = \text{Exp}(0=3) = \text{Exp}(\lambda=\frac{1}{3}) \rightarrow \begin{array}{l} \text{pdf: } f(t) = \frac{1}{3} e^{-\frac{1}{3}t} \\ \text{cdf: } F(t) = 1 - e^{-\frac{1}{3}t} = P(X \leq t) \end{array}$$

take 5 iid samples  $\sim$  ordered  $y_1 < y_2 < \dots < y_5$

(a) find pdf of  $y_3$

(b) find  $P(y_4 < 5)$

(c) find  $P(y_1 > 1)$

goal: translate the condition " $y_1 > 1$ " into a condition involving the original samples.

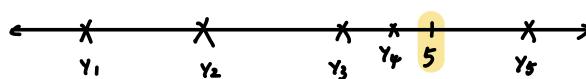
In this case: " $y_1 > 1$ " is equivalent to "all 5 samples are  $> 1$ "

$$P(y_1 > 1) = P(x_1 > 1 \text{ and } x_2 > 1 \text{ and } \dots \text{ and } x_5 > 1)$$

$$= P(x_1 > 1) P(x_2 > 1) \dots P(x_5 > 1)$$

$$= (1 - F(1))^5 = \exp(-\frac{1}{3})^5 \approx 0.19$$

(b)  $P(y_4 < 5) = P(\text{something about the samples})$



$P(\text{At least 4 of the 5 samples are } < 5)$

$$P(x_1, x_3, x_4, x_5 < 5)$$

$$= P(4 \text{ samples below} | \text{sample above}) + P(\text{all 5 samples below})$$

and  $x_2 > 5$ )

$$= \binom{5}{4} F(s)^4 (1 - F(s))^1 + F(s)^5$$

$$= F(s)^4 (1 - F(s))^1$$

$$= 5 (1 - \exp(-\frac{5}{3}))^4 (\exp(-\frac{5}{3}) + (1 - \exp(-\frac{5}{3}))^5 \approx 0.96$$

( $\binom{5}{4}$ ) ways to rearrange

(a) find pdf of  $Y_3$

find CDF first:  $F_{Y_3}(u) = P(Y_3 \leq u)$

$P(Y_3 \leq u) = P(\text{at least 3 samples are } \leq u)$

$$= P(3 \text{ below}, 2 \text{ above}) + P(4 \text{ below}, 1 \text{ above}) + P(5 \text{ below})$$

$$= \binom{5}{3} F(u)^3 (1-F(u))^2 + \binom{5}{4} F(u)^4 (1-F(u))^1 + F(u)^5$$

$$= 10(1 - \exp(-\frac{u}{3}))^3 \exp(-\frac{u}{3})^2 + 5(1 - \exp(-\frac{u}{3}))^4 \exp(-\frac{u}{3}) + (1 - \exp(-\frac{u}{3}))^5$$

→ pdf:  $\frac{d}{du}$  of this

### # 6.3.5

$X$ : Some distribution with  $P_{0.9} = 27.3$  (meaning  $P(X \leq 27.3) = 0.9$ )

8 iid samples from  $X \rightarrow$  order  $y_1 < y_2 < y_3 \dots < y_8$

find  $P(y_5 < 27.3 < y_8)$

↓ at most 7 samples are below  
at least 5 samples are below

$$\Rightarrow P(y_5 < 27.3 < y_8) = P(\# \text{ of samples that are } < 27.3 \text{ is } 5, 6 \text{ or } 7)$$

$$= P(5 \text{ below}, 3 \text{ above}) + P(6 \text{ below}, 2 \text{ above}) + P(7 \text{ below}, 1 \text{ above})$$

$$= \binom{8}{5} (0.9)^5 (0.3)^3 + \binom{8}{6} (0.9)^6 (0.3)^2 + \binom{8}{7} (0.9)^7 (0.3)^1$$

#### 4/14 Maximum Likelihood Estimators (MLE)

- Sampling coming from known family of distribution (Exp., Geom., etc.)
- goal estimate the parameter(s) based on the sample

What is a "good" estimate?

Idea of MLE:

- Likelihood Function  $L(\text{parameters}) = P(\text{getting the sample that we got})$
- MLE = choice of parameters that maximizes this likelihood

#### Simple Example

You have a coin with  $P(\text{heads}) = p$

flip 10 times ~ THHTTTTHHT

how to guess  $p$ ?

intuitive guess:  $p = \frac{4}{10}$

Let's calculate the MLE

$$\begin{aligned}L(p) &= P(\text{THHTTTTHHT}) \\&= (1-p) \cdot p \cdot p \cdot (1-p)^4 \cdot p \cdot (1-p) \cdot p \\&= p^4 (1-p)^6\end{aligned}$$

To find the MLE: do calculus to find the maximum

trick to make calculation easier:

$$\begin{aligned}\text{consider } \ell(p) &= \log L(p) = \log(p^4 (1-p)^6) \\&= 4 \cdot \log p + 6 \cdot \log(1-p)\end{aligned}$$

$$\text{set } 0 = \frac{d\ell}{dp} = \frac{4}{p} - \frac{6}{1-p}$$

$$\text{solve for } p \rightsquigarrow \boxed{p = \frac{4}{10}}$$

More complicated example  $\rightarrow$  pdf  $f(t) = \lambda \cdot \exp(-\lambda t)$

Sample from  $\text{Exp}(\lambda) \rightarrow x_1, x_2, \dots, x_n$

Find MLE for  $\lambda$

$$L(\lambda) = \text{IP} (\text{getting the sample } x_1, x_2, \dots, x_n) \quad (\text{exp is a continuous distribution, so IP of any specific sample is 0})$$

$$= f(x_1) \cdot f(x_2) \cdots f(x_n)$$

$$L(\lambda) = \lambda \exp(-\lambda x_1) \cdot \lambda \exp(-\lambda x_2) \cdots \lambda \exp(-\lambda x_n)$$

$$= \lambda^n \cdot \exp(-\lambda x_1 - \lambda x_2 - \dots - \lambda x_n)$$

$$= \lambda^n \exp(-\lambda (x_1 + x_2 + \dots + x_n))$$

$$\begin{aligned} \text{log likelihood } l(\lambda) &= \log L(\lambda) = \log (\lambda^n \cdot \exp(-\lambda (x_1 + \dots + x_n))) \\ &= \log (\lambda^n) - \lambda (x_1 + \dots + x_n) \\ &= n \cdot \log \lambda - \lambda (x_1 + \dots + x_n) \end{aligned}$$

$$\text{Set } 0 = \frac{d\ell}{d\lambda} = \frac{n}{\lambda} - (x_1 + \dots + x_n)$$

$$\text{Solve } \rightarrow \lambda = \frac{n}{x_1 + \dots + x_n} = \frac{1}{\bar{x}} \quad \text{Recall: } \mathbb{E}[\exp(\lambda)] = \frac{1}{\lambda}$$

Hard example

Sample  $x_1, x_2, \dots, x_n$  from  $\text{Unif}(0, T) \rightarrow$  pdf  $f(t) = \begin{cases} \frac{1}{T} & \text{for } 0 < t < T \\ 0 & \text{else} \end{cases}$

find MLE for  $T$

$$\text{X } L(T) = f(x_1) f(x_2) \cdots f(x_n)$$

$$= \frac{1}{T} \cdot \frac{1}{T} \cdot \dots \cdot \frac{1}{T}$$

$$= \frac{1}{T^n} \quad \text{X}$$

"maximum" is  $\infty$  at  $T=0$  (doesn't make sense  $\text{Unif}(0,0)$ )

$$L(T) = \left\{ \begin{array}{ll} \frac{1}{T} & \text{if } 0 < t < T \\ 0 & \text{else} \end{array} \right\} \cdots \left\{ \begin{array}{ll} \frac{1}{T} & \text{if } 0 < x_n < T \\ 0 & \text{else} \end{array} \right\} \rightarrow \text{只要有任何一个 } x_i \text{ 不在 } 0 < t < T \text{ 的 range, the whole thing is 0.}$$

$$= \begin{cases} \frac{1}{T^n} & \text{if all of the } x_j < T \\ 0 & \text{if any } x_j > T \end{cases}$$

(相乘)

$\theta$  can never be the maximum, so find the MLE we just need to find the maximum of  $L(T) = \frac{1}{T^n}$  subject to the constraint that  $T > \text{all of the } x_i$

try calculus:  $L'(T) = (-n)T^{-n-1} = 0$  no solutions because  $L'(T) < 0$  for all  $T$

$\frac{1}{T^n}$  is a decreasing function of  $T$ , so the maximum happens at the smallest possible value of  $T$ .

with constraint:  $T = \max(x_1, x_2, \dots, x_n)$

4/21

### Method of moments & Percentile matching

method of moments (MoM):

- if  $X$  is a RV, then the  $k^{\text{th}}$  moment is  $E(X^k)$
- if you have a sample  $x_1, \dots, x_n$ , you can calculate the  $k^{\text{th}}$  sample moment
- $m_k = \frac{1}{n} (x_1^k + x_2^k + \dots + x_n^k)$
- our distribution has some unknown parameter(s)  $\theta$
- MoM for estimating  $\theta$ :
  - calculate theoretical moments (function of  $\theta$ )
  - compare to sample moments
  - Solve for  $\theta$

### examples

# 6.4. 10

$$f(k) = (1-p)^{k-1} \cdot p \quad k=1, 2, 3, \dots$$

Sample  $x_1, \dots, x_n \sim \text{Geom}(p)$

estimate  $p$  using MoM

theoretical | sample

①  $E[\text{Geom}(p)] = \frac{1}{p}$   $\xleftarrow{\text{match}}$   $m_1 = \frac{1}{n} (x_1 + \dots + x_n) = \bar{x}$

so our estimator for  $p$  is  $\frac{1}{\bar{x}} = \bar{x} \rightarrow \boxed{p = \frac{1}{\bar{x}}}$

# 6.4.15

pdf  $f(x) = \frac{x^{\alpha-1} \exp(-\frac{x}{\theta})}{\Gamma(\alpha) \cdot \theta^\alpha}$

Sample  $x_1, \dots, x_n \sim \text{Gemma}(\alpha, \theta)$

estimate  $\alpha$  and  $\theta$  using MoM

given :  $E = \alpha\theta$  and  $\text{var} = \alpha\theta^2$

could calculate  $E(x)$  and  $E(x^2)$  from pdf

$\int_0^\infty x \cdot f(x) dx = \dots$

theoretical | sample

①  $E(x) = \alpha\theta$        $m_1 = \frac{1}{n} (x_1 + \dots + x_n)$

②  $E(x^2) = \text{Var}(x) + E(x)^2$        $m_2 = \frac{1}{n} (x_1^2 + \dots + x_n^2)$

$= \alpha\theta^2 + \alpha^2\theta^2$        $m_3 = \frac{1}{n} (x_1^3 + \dots + x_n^3)$

$\vdots$

$\alpha\theta = m_1$

$\alpha\theta^2 + \alpha^2\theta^2 = m_2$

solve for  $\alpha\theta$  in terms of  $m_1, m_2$

### Percentile matching :

- to estimate unknown parameters, match theoretical percentile values with sample percentile
- # of different percentile points you read to compare =

# of unknown parameters

you need to estimate

# 6.4.21

pdf  $f(x) = \frac{1}{\theta} \exp(-\frac{x}{\theta}), x > 0$

$\text{Exp}(\theta)$

estimate  $\theta$  by matching 0.5 percentile

We need to match the theoretical median to the sample median

- sample median let's call  $\tilde{\pi}_{0.5}$

- Theoretical median  $\pi_{0.5}$  is defined by  $P(\text{Exp}(\theta) < \pi_{0.5}) = \frac{1}{2}$

$$\frac{1}{2} = \int_0^{\pi_{0.5}} f(x) dx = \int_0^{\pi_{0.5}} \frac{1}{\theta} \exp(-\frac{x}{\theta}) dx = \left[ -\exp(-\frac{x}{\theta}) \right]_0^{\pi_{0.5}} = 1 - \exp(-\frac{\pi_{0.5}}{\theta})$$

$$\longrightarrow \pi_{0.5} = \theta \cdot \log(2)$$

$$\text{match } \pi_{0.5} = \tilde{\pi}_{0.5}$$

$$\theta \log 2 = \tilde{\pi}_{0.5}$$

$$\rightarrow \text{our estimate is } \theta = \frac{\tilde{\pi}_{0.5}}{\log(2)}$$

### Other example

distribution:  $\text{Unif}(a, b)$

let's use 25<sup>th</sup> and 75<sup>th</sup> percentile

sample percentiles  $\tilde{\pi}_{\frac{1}{4}}$  and  $\tilde{\pi}_{\frac{3}{4}}$  match

theoretical percentiles  $\pi_{\frac{1}{4}} = a + \frac{1}{4}(b-a)$

$$\pi_{\frac{3}{4}} = a + \frac{3}{4}(b-a)$$

solve for  $a, b$  in terms of  $\tilde{\pi}_{\frac{1}{4}}$  and  $\tilde{\pi}_{\frac{3}{4}}$

## 4/28 Interval estimation

- so far: all point estimators (guessing a single # for the unknown parameter)
- more useful to report a range of values along with some level of confidence that your range captures the true value.

### Simple example

Sample  $x_1, x_2, \dots, x_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$  unknown  
assume to  
be known

Find a 95% confidence interval for  $\mu$

- meaning: we will give a range, expressed as a function of the samples  $x_1, \dots, x_n$ , such that  $\mu$  is within the range with 95% probability.
- how do we do this?

idea: something to do with  $\bar{x}$

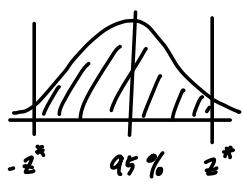
each  $x_i \sim N(\mu, \sigma^2)$

so  $x_1 + \dots + x_n \sim N(n\mu, n\sigma^2)$

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

one final transformation:  $\frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$

let's write down an interval that has 95% probability:



there is some number  $z^*$  so that

$$P(-z^* < \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} < z^*) = 0.95$$

↓ algebraic manipulation

$$P(\bar{x} - z^* \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + z^* \frac{\sigma}{\sqrt{n}})$$

conclusion: We can say that the range

$\boxed{\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}}}$  has a 95% chance of containing  $\mu$

↓  
margin of error

- How to find the number  $z^*$ ?

use the lookup table for 95%:  $z^* = 1.96$

Q : what if the  $x_i$  are not normal?

use central limit theorem : it doesn't matter if the  $x_i$  are normal or not,

$\bar{x}$  is always approximately normal if  $n$  is big

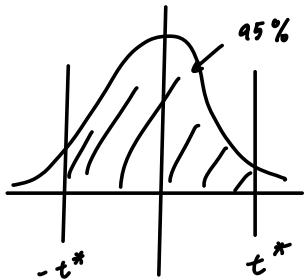
Q : what do we do if  $\sigma$  is also unknown?

We can try just replacing  $\sigma$  by  $s$  (the sample st. dev)  
everywhere it appears.

↓  
divided by  $n-1$  instead of  $n$

problem.  $\frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$  does not follow  $N(0,1)$

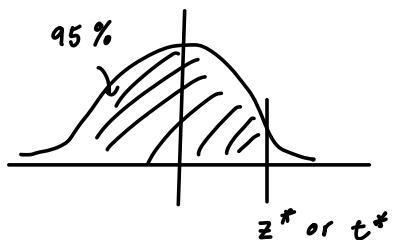
- It follows something called a T distribution with  $n-1$  degrees of freedom.



we just have to use a different look up table

Q : what if we want  $\mu$  in the range  $[\bar{x} - \_, \alpha]$ ?

instead of range  $\bar{x} \pm \_$



## 5/5 C.I. for difference of 2 means

simplest setting:  $X \sim N(\mu_x, \sigma_x^2 = \text{known})$  } independent  
 $Y \sim N(\mu_y, \sigma_y^2 = \text{known})$

goal: CI for  $\mu_x - \mu_y$

We know that  $\bar{X} \sim N(\mu_x, \frac{\sigma_x^2}{n_x})$  } still independent  
 $\bar{Y} \sim N(\mu_y, \frac{\sigma_y^2}{n_y})$

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

$$\Rightarrow \frac{\bar{X} - \bar{Y} - (\mu_x - \mu_y)}{\sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}} \sim N(0, 1)$$

→ CI for  $\mu_x - \mu_y$  is  $(\bar{X} - \bar{Y}) \pm z^* \cdot \sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}$   
 finding  $z^*$  is exactly the same as before

What happens if:

- $\sigma_x, \sigma_y$  are unknown?
- $X, Y$  are not normal?
- $X, Y$  are not independent?

Midterm Solution

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

estimate  $\mu, \sigma^2$  using percentile matching w/  $p = 0.3, 0.7$

how to find  $\pi_{0.3}, \pi_{0.7}$  of  $N(\mu, \sigma^2)$ ?

- transform from  $N(0,1)$ : if we know that  $30^{\text{th}}\%$  of  $N(0,1)$  is  $M$ ,

$$\text{then } \pi_{0.3} = \sigma \cdot M + \mu ?$$

use Z table: for  $N(0,1)$ ,  $30^{\text{th}}\% = -0.525$   $\rightsquigarrow \pi_{0.3} = -0.525 \sigma + \mu$

$$70^{\text{th}} = 0.525 \quad \pi_{0.7} = 0.525 \sigma + \mu$$

$$\rightarrow \hat{\mu} = 0 \quad \hat{\sigma} = 1.924$$

using MoM:

$$E(x) = \mu \quad m_1 = \frac{1}{n} \sum x_i$$

$$\begin{aligned} E(x^2) &= \text{var}(x) + E(x)^2 \\ &\quad m_2 = \frac{1}{n} \sum x_i^2 \\ &= \sigma^2 + \mu^2 \end{aligned}$$

# 4

$$x \sim P_{\theta_1}(\theta)$$

$$\text{Prior } h(\theta) = \frac{5}{\theta^6} \quad \theta > 1$$

$$\begin{aligned} \text{Bayes Rule: posterior } h(\theta|x_1, \dots, x_n) &= \frac{g(x_1, \dots, x_n; \theta) \cdot h(\theta)}{C} \\ &= \frac{\left( e^\theta \cdot \frac{\theta^{x_1}}{x_1!} \right) \cdots \left( e^\theta \cdot \frac{\theta^{x_n}}{x_n!} \right) \cdot \frac{5}{\theta^6}}{C} \end{aligned}$$

$$\text{Gam}(d, \beta): \quad = C \cdot e^{-n\theta} \cdot \theta^{x_1 + \dots + x_n - 6}$$

$$\text{pdf } f(x) = C \cdot x^{d-1} \cdot e^{-\theta} \quad = \text{Gam}(d = x_1 + \dots + x_n - 5, \beta = \frac{1}{n})$$

$$x > 0$$

May 12

### confidence intervals for percentiles

$X \sim$  some continuous distribution

goal: find CI for the median  $m = \pi_{0.5}$  (or some other percentile)

- start with the point estimate (sample percentile)

### example

$n=9$  samples, want a CI for the median  $m$

point estimate =  $Y_5$

possible intervals:  $(Y_4, Y_6)$

$$\left. \begin{array}{l} (Y_3, Y_7) \\ (Y_2, Y_8) \end{array} \right\} \begin{array}{l} \text{隨機選} \\ \text{都可用} \end{array}$$

What is the confidence level of these intervals?

$$\begin{aligned} (Y_3, Y_7) : \quad & \Pr(Y_3 < m < Y_7) = \Pr(\text{the # of samples that are } \leq m \text{ is } 3, 4, 5, 6) \\ & = \binom{9}{3} \cdot \Pr(X < m)^3 \cdot \Pr(X > m)^6 + \dots + \binom{9}{6} \Pr(X < m)^6 \Pr(X > m)^3 \\ & = \binom{9}{3} \left(\frac{1}{2}\right)^3 \left(\frac{1}{2}\right)^6 + \dots + \binom{9}{6} \left(\frac{1}{2}\right)^6 \left(\frac{1}{2}\right)^3 = \dots \end{aligned}$$

### example

$n=12$ , fin CI for  $\pi_{0.3}$

- point estimate:  $(12+1)(0.3) = 3.9 \rightarrow \tilde{\pi}_{0.3} = (0.1) Y_3 + (0.9) Y_4$
- let's just pick  $Y_4$  to serve as our center
- options for intervals:  $(Y_3, Y_5) \leftarrow$

$(Y_2, Y_6)$

:

$$\begin{aligned} \mathbb{P}(Y_3 < \pi_{0.3} < Y_5) &= \mathbb{P}(\# \text{ samples} < \pi_{0.3} \text{ is } 3 \text{ or } 4) \\ &= \binom{12}{3} \cdot p(x < \pi_{0.3})^3 \mathbb{P}(x > \pi_{0.3})^9 + \dots \\ &= \binom{12}{3} (0.3)^3 (0.7)^9 + \dots \end{aligned}$$

### example

n samples (n is big)

We want to find a 95% CI for the median m

- point estimate =  $\bar{Y}_{\frac{n}{2}}$
- only consider intervals of the form  $(\bar{Y}_{\frac{n}{2}-k}, \bar{Y}_{\frac{n}{2}+k})$

Q: find the best k so that the confidence level of this interval is  $\approx 95\%$

$$\begin{aligned} \mathbb{P}(\bar{Y}_{\frac{n}{2}-k} < m < \bar{Y}_{\frac{n}{2}+k}) &= \mathbb{P}(\underbrace{\# \text{ samples} < m}_{W \sim \text{Bin}(n, \frac{1}{2})} \text{ is in the interval } (\frac{n}{2}-k, \frac{n}{2}+k)) \\ &= \mathbb{P}(\frac{n}{2}-k < \text{Bin}(n, \frac{1}{2}) < \frac{n}{2}+k) \end{aligned}$$

normal approximation / CLT: when n is big,  $\text{Bin}(n, p) \approx N(np, np(1-p))$

in our case:  $W \approx N(\frac{n}{2}, \frac{n}{4})$

$$\begin{aligned} (\bar{Y}_{\frac{n}{2}-k}, \bar{Y}_{\frac{n}{2}+k}) \text{ has confidence } &\approx \mathbb{P}(\frac{n}{2}-k < N(\frac{n}{2}, \frac{n}{4}) < \frac{n}{2}+k) \\ &= \mathbb{P}\left(\frac{-k}{\sqrt{\frac{n}{4}}} < N(0,1) < \frac{k}{\sqrt{\frac{n}{4}}}\right) \end{aligned}$$

Set = 95%, solve for k:

$$\frac{k}{\sqrt{\frac{n}{4}}} = 1.96$$

$$k = \frac{1.96}{\sqrt{2}} \sqrt{n}$$

for example if n=100, then  $k = \frac{1.96}{\sqrt{2}} = 9.8 \approx 10$

so an approximate 95% CI for m is  $(\bar{Y}_{40}, \bar{Y}_{60})$

May 19

### hypothesis testing

Q: are UCLA students smarter than average?

Assume: average IQ = 110

IQ of UCLA students  $\sim N(\mu, \sigma^2 = 100)$

Q: is  $\mu > 110$ ?

Suppose we sample  $n=16$  UCLA students and get  $\bar{x} = 113.5$

- Is this enough evidence to conclude that  $\mu > 110$ ?

#### Set up:

null hypothesis  $H_0: \mu = 110$

alternative hypothesis  $H_1: \mu > 110$

We ask the question: if we assume  $H_0$  is true, then what is the probability of getting a sample at least as extreme as the one we actually got?

assuming  $H_0: X \sim N(\mu = 110, \sigma^2 = 100) \rightarrow \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$

calculate  $P(\bar{X} \geq 113.5) = P\left(\frac{\bar{X} - 110}{\frac{10}{\sqrt{16}}} \geq \frac{3.5}{\frac{10}{\sqrt{16}}}\right) = P(N(0, 1) \geq 1.4) \approx 0.08$  → called the p-value

- Is this p-value small enough to conclude that  $H_0$  is false?
- It depends on what significance level we are using for the test
- For example,  $\alpha = 10\%$ , then  $8\% < 10\% \rightarrow$  reject the  $H_0$  in favor of  $H_1$
- using  $\alpha = 5\%$ , then  $8\% > 5\% \rightarrow$  fail to reject the  $H_0$

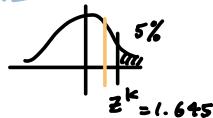
you will find the  $H_0$  value for  $\mu$  is outside the appropriate CI



the p-value of your sample is  $< \alpha$

another way to think about it:

- when we get our sample  $\bar{x}$ , calculate the "test statistics"  $Z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$
- compare  $Z$  to whatever critical threshold is
- in this case, the critical threshold is



### # 8.1.4

$$x \sim N(\mu, \sigma^2 = \text{unknown})$$

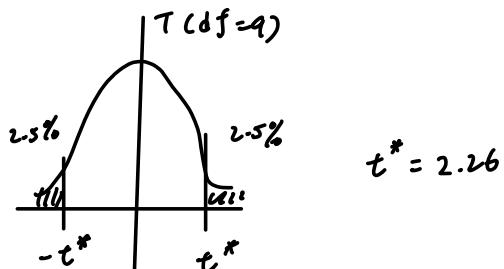
$$H_0: \mu = 9.5$$

$$H_1: \mu \neq 9.5$$

$$n = 10, \bar{x} = 9.55, s_x = 0.1027$$

do a hypothesis test at 5% significance

- calculate our test statistic  $T = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} = \frac{9.55 - 9.5}{0.1027/\sqrt{10}} \approx 1.539$
- what is our critical threshold for this test?



because our  $T$  is not extreme enough for the threshold of this test, we fail to reject  $H_0$ .

## Test for variances

general structure of hypothesis tests:

1) get a point estimator to base the test on (for means:  $\bar{x}$ )

2) transform it into a "test statistics" that

we know follows some special distribution (For means:  $Z = \frac{\bar{x} - \mu_0}{\frac{\sigma}{\sqrt{n}}}$  or  $T = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}$ )

3) Compare your test statistics to the appropriate critical values.

Test for variance:

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

base estimator  $s^2$

$$n=11 \text{ samples}$$

$$\text{test statistic is called } \chi^2 = \frac{(n-1)s^2}{\sigma^2} = \frac{10 s^2}{525^2}$$

$$H_0: \sigma^2 = 525^2$$

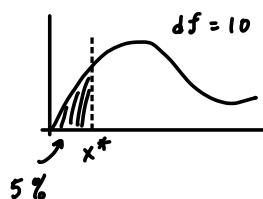
• this follows a distribution called the

$$H_1: \sigma^2 < 525^2$$

" $\chi^2$  distribution with  $n-1=10$  d.o.f"

$$\alpha = 5\%$$

Find the critical region:



use the  $\chi^2$  lookup table:  $x^* = 3.94$

critical region: Reject  $H_0$  if  $\chi^2 < 3.94$

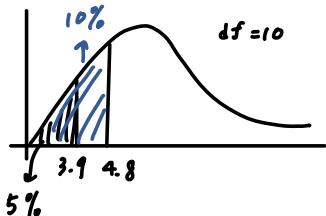
$$s^2 < \frac{(3.94)(525^2)}{10}$$

Let's say we take an actual sample and get  $s^2 = 113, 108.5$

$$\rightarrow \chi^2 = \frac{10(113,108.5)}{525^2} = 4.1 > 3.94, \text{ we fail to reject}$$

What is the p-value of this sample?

$$p\text{-val} = P(\chi^2(df=10) \leq 4.1) \in (5\%, 10\%)$$



→ 沒有辦法找到 exact value  
→ approx.

$$X \sim N(\mu_x, \sigma_x^2)$$

$$Y \sim N(\mu_y, \sigma_y^2)$$

$n_x = 11$  samples

$n_y = 9$  samples

$$H_0: \sigma_x^2 = \sigma_y^2$$

$$H_1: \sigma_x^2 > \sigma_y^2 \quad / \quad \sigma_x^2 < \sigma_y^2 \Rightarrow 5\% \text{ left tail}$$

start with  $s_x^2, s_y^2$

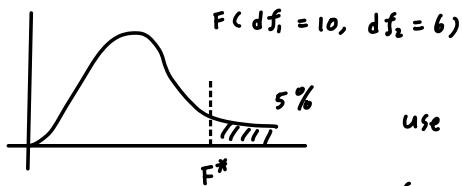
$$\text{flip } F = \frac{s_y^2}{s_x^2}$$

$$\text{test statistic} = F = \frac{s_x^2}{s_y^2}$$

this follows the F distribution with  $df_1 = n_x - 1$ ,

$$df_2 = n_y - 1$$

Find the critical region:



$$F(df_1 = 10, df_2 = 6)$$

use the F look up table:  $F^* = 4.06$

so our critical region for rejecting  $H_0$  is  $F = \frac{s_x^2}{s_y^2} > 4.06$

For  $\sigma_x^2 < \sigma_y^2$ : instead define  $F = \frac{s_y^2}{s_x^2} \sim F(df_1 = n_y - 1, df_2 = n_x - 1)$

and then do the same thing

6/2 ① Finals OH: Tues (June 7) 2-4 pm (in person MS 6139)

② Regrade grades  $\Rightarrow$  grade scope to grade book

### Hypothesis Tests for Median

"Wilcoxon signed rank test"

example (8.5.1)

(n = 13)

We want to test

sample: 41195 39485 41229 36840 38050 40890 38345

$H_0: m = 40,000$

34930 39245 31031 40780 38050 30906

$H_1: m < 40,000$

$\alpha = 5\%$

Step 1: look at the differences:

+5 -1 +6 -10 -8.5 +4 -7  
1195 -515 1229 -3160 -1950 -8.5 840 -1655  
-5070 -755 -8969 780 -1950 -8.5 -9040 -13

Step 2: assign "signed ranks"

Step 3: calculate

$$W = \text{sum of signed ranks} = -1 - 2 + 3 + 4 + 5 + 6 - 7 - 8.5 - 8.5 - 10 - 11 - 12 - 13$$

$$= -55$$

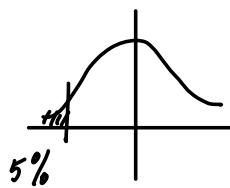
Step 4: apply the normalizing transformation:  $z = \sqrt{\frac{w}{\frac{(n+1)(2n+1)}{6}}}$

why is this true? it's complicated  
(textbook proofs)

$\sim N(0,1)$

$$z = \frac{-55}{\sqrt{\frac{(13)(14)(27)}{6}}} = -1.922$$

Step 5: compare to appropriate critical value



crit. region:  $z < -1.645$

our  $z$  was  $-1.92$ , so reject