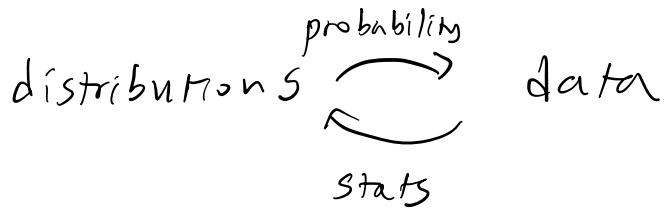


# 18.650 Spring 2024

LECTURE 1 2/5/24 (missed)



LLN:  $\tilde{X}_n \rightarrow \mu$  as  $n \rightarrow \infty$

CLT:  $\frac{\sqrt{n}}{\sigma} (\bar{X}_n - \mu) \rightsquigarrow N(0, 1)$

$X_i$ 's are  
i.i.d  
mean  $\mu$ , variance  $\sigma^2$

LECTURE 2 2/7/24 (missed)

- histograms
- summary statistics
  - mean, std dev, median, quantiles, IQR, outliers

100(1-d)  
percentile  $\rightarrow$  quantile of order 1-d  
 $= q_d = \text{value that } d \text{ of data exceeds}$

- boxplots, scatterplots.

$X_n \xrightarrow{P} X$  as  $n \rightarrow \infty$  if  $\forall \epsilon > 0$ ,

"convergence in probability"  $P(|X_n - X| > \epsilon) \rightarrow 0$  as  $n \rightarrow \infty$ .

ex if  $X_i \sim \text{Ber}(1/2)$   $\forall n$ . Then  $X_n \xrightarrow{P} \text{Ber}(1/2)$   
because  $P(|X_n - X| > \epsilon)$  for  $\epsilon \in (0, 1)$

$$= P(X_n = 0 \cap X = 1) + P(X_n = 1 \cap X = 0) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \rightarrow 0$$

but  $X_n \rightsquigarrow X$  as  $n \rightarrow \infty$  if  $\forall x$  for which  $x \mapsto P(X \leq x)$  cdf is continuous  
"convergence in distribution"  $P(X_n \leq x) \rightarrow P(X \leq x)$  as  $n \rightarrow \infty$ .

$$\boxed{X_n \xrightarrow{P} X \Rightarrow X_n \rightsquigarrow X}$$

Lemma If  $X_n \rightsquigarrow c$  for constant  $c$ , then  $X_n \xrightarrow{P} c$ .

PF  $P(|X_n - c| > \epsilon) = P(X_n \geq \epsilon + c) + P(X_n \leq c - \epsilon)$   
 $\rightarrow 0 + 0 = 0$  as  $n \rightarrow \infty$ .

$$\text{if } X_n \xrightarrow{P} X \Rightarrow X_n + Y_n \xrightarrow{P} X + Y$$

$$X_n Y_n \xrightarrow{P} XY$$

$$\text{if } X_n \rightsquigarrow X \Rightarrow X_n + Y_n \rightsquigarrow X + c$$

$$Y_n \xrightarrow{P} c \Rightarrow X_n Y_n \rightsquigarrow Xc$$

$$\text{if } X_n \xrightarrow{P} X \Rightarrow g(X_n) \xrightarrow{P} g(X)$$

Delta Method

$$\frac{\sqrt{n}}{\sigma} (Y_n - \mu) \rightsquigarrow Y \sim N(0, 1)$$

implies for any diff'ble  $g$

$$\frac{\sqrt{n}}{\sigma} (g(Y_n) - g(\mu)) \rightsquigarrow N(0, g'(\mu)^2)$$

$$\text{if } X_n \rightsquigarrow X \Rightarrow g(X_n) \rightsquigarrow g(X)$$

# LECTURE 4 2/12/24 (Missed)

Gaussian  $X \sim \mathcal{N}(\mu, \sigma^2)$

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

• 68-95-99.7 rule

$$X \mapsto Z = \frac{X-\mu}{\sigma} \quad \text{Z-score} \quad \text{"standardization"}$$

standard normal:  $\frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$

## LECTURE 5 2/14/24 1PM

random matrix  $X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^K$

$$P(X \in A) = \int_A f(x_1, x_2, \dots, x_K) dx_1 dx_2 \dots dx_K$$

$$\mathbb{E}[X] = \begin{pmatrix} \mathbb{E}[x_1] \\ \vdots \\ \mathbb{E}[x_K] \end{pmatrix} = \mu \quad \text{Cov}(x_i, x_j) = \mathbb{E}[x_i x_j] - \mathbb{E}[x_i] \mathbb{E}[x_j]$$

covariance matrix of  $X$

$$\mathbb{V}[X] = \Sigma = (\text{Cov}(x_i, x_j))_{i,j}^*$$

$$\Rightarrow \sum_{ii} = \text{Cov}(x_i, x_i) = \text{Var}(x_i)$$

note  $X^T X$  = inner product = scalar

$XX^T$  = matrix  $\in \mathbb{R}^{k \times k}$  = outer product

$$\Sigma := V[X] = E[XX^T] - E[X]E[X]^T = E[XX^T] - \mu\mu^T$$

$\Sigma^{-1}$  = precision matrix

### Linear Transformations

Let  $X \in \mathbb{R}^k$  be random vector,  $E[X] = \mu$ ,  $V[X] = \Sigma$ .

Let  $a \in \mathbb{R}^k$  be deterministic vector.

$$\begin{aligned} \Rightarrow E[a^T X] &= a^T \mu & V[a^T X] &= a^T \Sigma a \\ &\quad \uparrow & &\quad \uparrow \\ &\text{trivial} & E[a^T X X^T a] &= a^T \mu \mu^T a \\ &&&= a^T E[XX^T] a - a^T \mu \mu^T a \\ &&&= a^T \Sigma a. \end{aligned}$$

Let  $A = \mathbb{R}^{k \times e}$  matrix (deterministic),  $b \in \mathbb{R}^e$  deterministic vector.

Then  $A^T X + b \in \mathbb{R}^e$  s.t.  $E[A^T X + b] = A^T \mu + b$   
 $V[A^T X + b] = A^T \Sigma A$

### Multivariate Gaussian

$$X \sim N_k(\mu, \Sigma)$$

$$\Rightarrow f(x) = \frac{1}{\sqrt{(2\pi)^k \det \Sigma}} \exp \left( -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right)$$

note

$$\textcircled{1} \quad A^T X + b \sim N_d(A^T \mu + b, A^T \Sigma A)$$

$$\textcircled{2} \quad Z = \Sigma^{-1/2}(X - \mu) \sim N_k(0, I_k), \quad X = \Sigma^{1/2}Z + \mu$$

### Multivariate CLT

$X_1, X_2, \dots, X_n \in \mathbb{R}^k$ ,  $\mathbb{E}[X_i] = \mu$ ,  $\mathbb{V}[X_i] = \Sigma$   
random

Then  $\sqrt{n}(\bar{X}_n - \mu) \sim N(0, \Sigma)$

### Delta Method

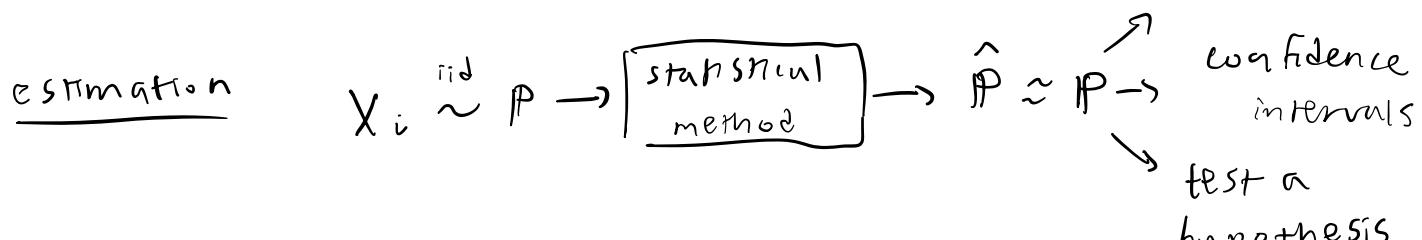
$X_1, \dots, X_n \in \mathbb{R}^k$  random,  $\mathbb{E}[X_i] = \mu$ ,  $\mathbb{V}[X_i] = \Sigma$ ,

$$g: \mathbb{R}^k \rightarrow \mathbb{R}$$

Then  $\sqrt{n}(g(\bar{X}_n) - g(\mu)) \rightsquigarrow N(0, \nabla g(\mu)^T \Sigma \nabla g(\mu))$

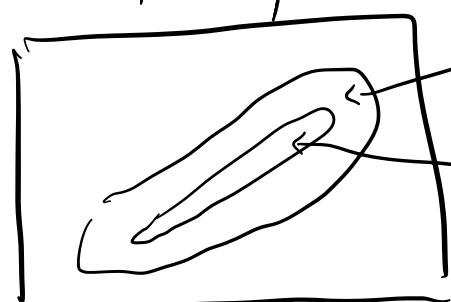
LECTURE 6 2/16/24 1PM

estimation



statistical model = subset of all probability distributions

ex) all pdfs/pdfs



$$\{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}$$

$$\{N(1, \sigma^2), \sigma^2 > 0\}$$

$$\{\text{Ber}(p) : p \in (0, 1)\}, \{\text{Exp}(\lambda), \lambda \in (0, 15)\}$$

we can all specify models based on

pmf/pdf

$$\left\{ \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \mu \in \mathbb{R}, \sigma^2 > 0 \right\}$$

or

cdf

$\{F(x) : F \text{ is a continuous cdf}\}$

generally

$$\left\{ P_\theta, \theta \in \Theta \right\}$$

↑  
parameter(s)      ↑  
                        parameter space

ex)

$$\theta = p$$

$$P_p = \text{Ber}(p)$$

dimension of statistical model  $\Theta$

$\text{Expl}(\lambda) \subset \Theta$  finite dimension  $\rightarrow$  parametric model

all polynomials  
as pdfs

$\Theta$  infinite dimension  $\rightarrow$  nonparametric model

Notation for Estimation

predicted  
 $\downarrow$   
 $\theta \sim \hat{\theta}$

$$P_\theta(X \geq 1), E_\theta[X], V_\theta[X]$$

$P(X \geq 1)$   
for  $X \sim P_\theta$

goal of : provide single best guess  
estimation  $\hat{\theta}$  for  $\theta$  based on data

An estimation  $\hat{\theta}$  is a function of the data

$$\hat{\theta} = g(X_1, X_2, \dots, X_n) \leftarrow \begin{matrix} \text{data} \\ \text{given} \end{matrix}$$

$$\textcircled{1} \quad \text{bias}(\hat{\theta}) = E_{\theta}[\hat{\theta}] - \theta$$

↑      ↓  
 "unbiased"      asymptotically unbiased  
 if      if      as  $n \rightarrow \infty$

ex  $\{N(\mu, 1), \mu \in \mathbb{R}\}$

$$\begin{array}{ll} \hat{\mu}_1 = \bar{x}_n & \text{bias}(\hat{\mu}_1) = 0 \\ \hat{\mu}_2 = x_1 & \text{bias}(\hat{\mu}_2) = 0 \\ \hat{\mu}_3 = 0 & \text{bias}(\hat{\mu}_3) = -\mu \end{array}$$

2 standard error

$$se(\hat{\theta}) = \sqrt{V_{\theta}[\hat{\theta}]} = \text{stdev}(\hat{\theta})$$

$$\begin{array}{l} se(\hat{\mu}_1) = \frac{1}{\sqrt{n}} \\ se(\hat{\mu}_2) = 1 \\ se(\hat{\mu}_3) = 0 \end{array}$$

3 mean squared error

$$MSE(\hat{\theta}) = \text{bias}^2(\hat{\theta}) + V_{\theta}[\hat{\theta}] = \text{bias}^2(\hat{\theta}) + se^2(\hat{\theta}) = E_{\theta}[(\hat{\theta} - \theta)^2]$$

$$MSE(\hat{\mu}_1) = \frac{1}{n} \quad MSE(\hat{\mu}_2) = 1 \quad MSE(\hat{\mu}_3) = \mu^2$$

consistency continuous mapping theorem often used  
Def  $\hat{\theta}_n$  is consistent estimator of  $\theta$  if  $\hat{\theta}_n \xrightarrow{P} \theta$

$$\text{ex)} \quad \hat{\mu}_1 = \bar{x}_n \xrightarrow{P} \mu \quad \checkmark \quad \hat{\mu}_2 = x_1 \xrightarrow{P} \mu \quad \times \quad \hat{\mu}_3 = 0 \xrightarrow{P} \mu \quad \times$$

Def

$\hat{\theta}_n$  asymptotically normal if  $\sqrt{n}(\hat{\theta}_n - \theta) \rightsquigarrow N(0, \sigma^2)$   
 for some  $\sigma > 0$

$\sigma^2$  is the asymptotic variance

$$\begin{array}{l} \text{ex)} \quad \bar{x}_n \Rightarrow \sqrt{n}(\bar{x}_n - \mu) \rightsquigarrow N(0, \sigma^2) \text{ by CLT} \\ g(\bar{x}_n) \Rightarrow \sqrt{n}(g(\bar{x}_n) - g(\mu)) \rightsquigarrow N(0, g'(\mu)^2 \sigma^2) \text{ by Delta Method} \end{array}$$

$$\rightarrow \hat{\theta}_n \approx N(\theta, \frac{\sigma^2}{n}) \approx \theta + \frac{\sigma}{\sqrt{n}} Z \leftarrow N(0, 1)$$

$$\Rightarrow se(\hat{\theta}_n) \approx \frac{\sigma}{\sqrt{n}}$$

another definition

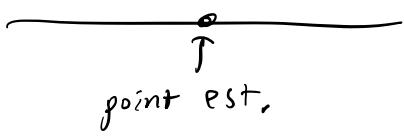
$\hat{\theta}_n$  is asymptotically normal if

$$\frac{\hat{\theta}_n - \theta}{\text{se}(\hat{\theta}_n)} \rightsquigarrow N(0, 1)$$

e.g.  $\hat{\theta}_n = \bar{x}_n$  is  
by CLT

LECTURE 7 2/20/24 1PM

## Confidence Intervals



Def A  $(1-\alpha)$  confidence interval is a random interval

$$C_n = (A_n, B_n) \text{ s.t. } P_\theta(\theta \in (A_n, B_n)) = P_\theta(A_n < \theta < B_n) \geq 1-\alpha \quad \forall \theta$$

↓ ↓

functions

$$A_n(x_1, \dots, x_n) \quad B_n(x_1, \dots, x_n)$$

↓

when computed,

$$\text{usually } \bar{x}_n - \frac{\sigma}{\sqrt{n}}$$

$$\bar{x}_n + \frac{\sigma}{\sqrt{n}}$$

either  $\theta \in (A_n, B_n)$  or NOT.

over many trials we get  $\sim 1-\alpha$  captured-

constructing a CI

Kiss Example  $X_1, \dots, X_n \sim \text{Ber}(p)$

$$\hat{p} = \bar{x}_n \approx N(p, \frac{p(1-p)}{n})$$

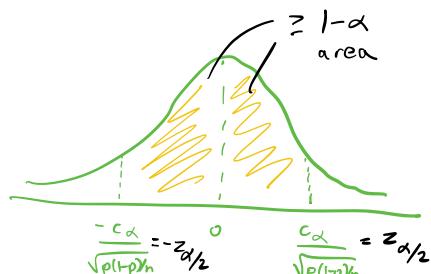
usually use estimator

want

$$1-\alpha = P_p(\bar{x}_n - c_\alpha < p < \bar{x}_n + c_\alpha) = P_p(p - c_\alpha \leq \bar{x}_n \leq p + c_\alpha)$$

$$= P_p\left(-\frac{c_\alpha}{\sqrt{p(1-p)/n}} \leq \frac{\bar{x}_n - p}{\sqrt{p(1-p)/n}} \leq \frac{c_\alpha}{\sqrt{p(1-p)/n}}\right)$$

$$\left(\bar{x}_n - \sqrt{\frac{\bar{x}_n(1-\bar{x}_n)}{n}} z_{\alpha/2}, \bar{x}_n + \sqrt{\frac{\bar{x}_n(1-\bar{x}_n)}{n}} z_{\alpha/2}\right)$$



parameter of interest vs. nuisance parameter

ex)  $N(\mu, \sigma^2)$

- $(\mu, \sigma^2)$

- $\mu$  but not  $\sigma^2$

- $\sigma^2$  but not  $\mu$

- $\frac{\mu}{\sigma}$  but not  $\mu$  or  $\sigma$  alone

might still need  
to be estimated

$$\mu \in \bar{x} \pm 1.96 \frac{\hat{\sigma}}{\sqrt{n}}$$

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

using  $\hat{\mu}$

given  $X_1, \dots, X_n \sim P_\theta$ , can we even hope to estimate  $\theta$ ?

identifiability -  $\theta$  identifiable if  $P_\theta = P_{\theta'}$ ,  $\rightarrow \theta = \theta'$

i.e.  $\theta \mapsto P_\theta$  is injective

ex)  $\{N(\mu, \sigma^2), \mu, \sigma \in \mathbb{R}\}$  is not because  $\frac{\partial}{\partial \sigma} \Rightarrow \sigma^2$  not injective

ex) dose-response

$$N(E_0 + \frac{D \times E_{max}}{D + ED_{50}}, \sigma^2)$$

$E_0$  = base effect (0 dose)  $E_{max}$  = max possible effect

$D$  = dose of drug (known)  $ED_{50}$  = dose producing half effect of  $E_{max}$

is not because varying  $E_{max}, ED_{50}$  can give same normal distribution

Plug-in Method

$E \rightsquigarrow \frac{1}{n} \sum_{i=1}^n$  substitution  
(only for expectations)

ex)  $N(\mu, \sigma^2)$

$$\mu = E[X] \rightsquigarrow \frac{1}{n} \sum_{i=1}^n x_i$$

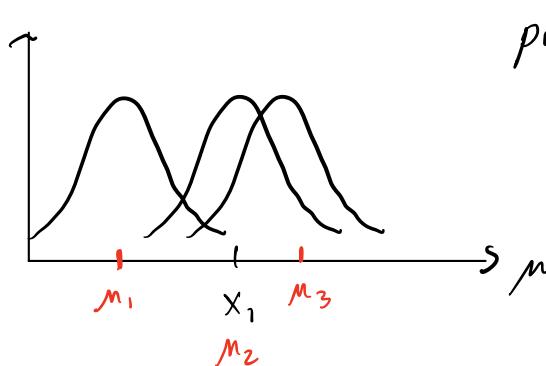
$$\sigma^2 = E[X^2] - E[X]^2 \rightsquigarrow \frac{1}{n} \sum_{i=1}^n x_i^2 - \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

asymptotically normal /  
consistent due to  
LLN and Continuous  
Mapping Theorem  
+ CLT / Delta

## Method of Maximum Likelihood

the parameter value most likely to have generated the data

ex)  $X_i \sim N(\mu^*, 1)$



pdf  $f_{\mu}(x_1)$  is highest  
when  $\mu = X_1$

$$\hat{\mu} = \arg \max_{\mu} f_{\mu}(x_1)$$

for  $X_1, X_2, \dots, X_n$  iid  $\sim f_{\theta^*}$

Def The likelihood is  $L_n(\theta) = f_{\theta}(X_1) f_{\theta}(X_2) \cdots f_{\theta}(X_n)$

maximum likelihood estimator (MLE)

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} L_n(\theta)$$

$$\rightarrow \hat{\theta}_n = \arg \max_{\theta \in \Theta} f_{\theta}(X_1) f_{\theta}(X_2) \cdots f_{\theta}(X_n)$$

Def log likelihood  $l_n(\theta) = \log L_n(\theta) = \sum_{i=1}^n \log f_{\theta}(X_i)$

increasing, so  $\hat{\theta}_n = \arg \max_{\theta \in \Theta} l_n(\theta)$  too.

also called

sample log likelihood

MLE

(1) Consistent:  $\hat{\theta}^{\text{MLE}} \xrightarrow{P} \theta^*$

(2) Asymptotically normal:  $\sqrt{n}(\hat{\theta}^{\text{MLE}} - \theta^*) \sim N(0, \sigma_{\text{MLE}}^2)$

(3) Asymptotic efficiency: for any other asymptotically normal  $\hat{\theta}$  with  $\sqrt{n}(\hat{\theta} - \theta^*) \sim N(0, \sigma^2)$ ,  $\sigma^2 \geq \sigma_{\text{MLE}}^2$

<sup>iF</sup>  
log likelihood  
is diff'ble

why is MLE good? notion of distance  $\text{dist}(P_{\theta^*}, P_{\theta})$

(1) computable from samples (approximation)

(2) minimized only at  $\theta^*$

Kullback-Leibler (KL) divergence  $\leftarrow$  not a distance!

$$D_{\text{KL}}(P_{\theta^*} \| P_{\theta}) = \int \log\left(\frac{f_{\theta^*}(x)}{f_{\theta}(x)}\right) f_{\theta^*}(x) dx \geq 0 \text{ always}$$

$$= \int \log(f_{\theta^*}(x)) f_{\theta^*}(x) dx$$

$$- \int \log(f_{\theta}(x)) f_{\theta^*}(x) dx$$

Jensen

$= 0$  if  $P_{\theta} = P_{\theta^*}$

$\downarrow$

$\theta = \theta^*$   
if identifiable

$$\cdot \quad \theta^* = \arg \min_{\theta} D_{\text{KL}}(P_{\theta^*} \| P_{\theta}) \leftarrow 0 \text{ only at } \theta = \theta^*$$

$$= \arg \max_{\theta} \int \log(f_{\theta}(x)) f_{\theta^*}(x) dx \leftarrow \text{since } \theta \text{ only appears here}$$

$$= \arg \max_{\theta} E_{\theta^*} [\log f_{\theta}(X)], \quad X \sim P_{\theta^*}$$

$$\approx \arg \max_{\theta} \frac{1}{n} \sum_{i=1}^n \log f_{\theta}(X_i) = \hat{\theta}^{\text{MLE}}$$

sample log likelihood  
 used to approximate

$\ell(\theta) = \mathbb{E}_{\theta^*} [\log f_{\theta}(X)]$  is the population log likelihood

$$\theta^* = \arg \max_{\theta} \ell(\theta) \Rightarrow \ell'(\theta^*) = 0.$$

note  $\frac{1}{n} \ln(\theta) \xrightarrow{P} \ell(\theta)$  via LLN.

so the maximizer  $\hat{\theta}^{\text{MLE}}$  of  $\frac{1}{n} \ln(\theta)$  converges to the maximizer  $\theta^*$  of  $\ell(\theta)$

### Asymptotic normality

$$\text{Thm } \sqrt{n} (\hat{\theta}^{\text{MLE}} - \theta^*) \rightsquigarrow N(0, I(\theta^*)^{-1})$$

$I(\theta^*)$  is Fisher Information

$$I(\theta^*) = \mathbb{E}_{\theta} \left[ -\frac{d^2}{d\theta^2} \log f_{\theta}(X) \right] = \mathbb{V}_{\theta} \left[ \frac{d}{d\theta} \log f_{\theta}(X) \right]$$

random variable  
 X evaluated @  $\theta$   
 Hessian for  $\nabla_{\theta}^2$   
 multidimensional

gradient  
 for  $\nabla$   
 multidimensional

ex  $P_{\theta} = N(\theta, 1)$

$$f_{\theta}(X) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(X-\theta)^2}{2}}$$

$$\log f_{\theta}(X) = \log \left( \frac{1}{\sqrt{2\pi}} \right) - \frac{(X-\theta)^2}{2}$$

$$\frac{d}{d\theta} \log f_{\theta}(X) = X - \theta \rightarrow \mathbb{V}_{\theta} [X - \theta] = 1$$

~  $N(\theta, 1)$

PF Let  $\hat{\theta} = \hat{\theta}^{\text{MLE}}$ . WTS  $\sqrt{n}(\hat{\theta} - \theta^*) \rightsquigarrow N_k(0, I(\theta^*)^{-1})$

ID: Taylor Expansion  $0 = l_n'(\hat{\theta}) \approx l_n'(\theta^*) + l_n''(\theta^*)(\hat{\theta} - \theta^*)$   
 by definition

$$\sqrt{n}(\hat{\theta} - \theta^*) \approx -\sqrt{n} \frac{l_n'(\theta^*)}{l_n''(\theta^*)} = -\sqrt{n} \frac{\sum_{i=1}^n \frac{d}{d\theta} \log f_\theta(X_i) \Big|_{\theta=\theta^*}}{\sum_{i=1}^n \frac{d^2}{d\theta^2} \log f_\theta(X_i) \Big|_{\theta=\theta^*}}$$

$$Y_i = \frac{d}{d\theta} \log f_\theta(X_i) \Big|_{\theta=\theta^*} = -\sqrt{n} \frac{\sum_{i=1}^n Y_i}{\sum_{i=1}^n W_i} = -\sqrt{n} \frac{\bar{Y}_n}{\bar{W}_n}$$

$$W_i = \frac{d^2}{d\theta^2} \log f_\theta(X_i) \Big|_{\theta=\theta^*}$$

$$\mathbb{E}[Y_i] = \mathbb{E}_{\theta^*} \left[ \frac{d}{d\theta} \log f_\theta(X_i) \Big|_{\theta=\theta^*} \right] = \frac{d}{d\theta} \mathbb{E}_{\theta^*} [\log f_\theta(X_i)] \Big|_{\theta=\theta^*} = l'(\theta^*) = 0$$

$$\mathbb{V}[Y_i] = \mathbb{V}_{\theta^*} \left[ \frac{d}{d\theta} \log f_\theta(X_i) \Big|_{\theta=\theta^*} \right] = I(\theta^*)$$

$$\bar{W}_n \xrightarrow{\mathbb{P}} \mathbb{E}[W_i] = \mathbb{E}_{\theta^*} \left[ \frac{d^2}{d\theta^2} \log f_\theta(X_i) \Big|_{\theta=\theta^*} \right] = -I(\theta^*) \quad \text{by LLN}$$

$$\sqrt{n}\bar{Y}_n = \sqrt{n}(\bar{Y}_n - 0) \rightsquigarrow N(0, I(\theta^*))$$

Slutsky's and CLT:  $\sqrt{n}(\hat{\theta} - \theta^*) \approx -\frac{\sqrt{n}\bar{Y}_n}{\bar{W}_n} \rightsquigarrow N(0, I^{-1}(\theta^*))$

LECTURE 10 3/1/24 1pm

recall  $\theta^*$  = true parameter  $\frac{1}{n} l_n \xrightarrow{\mathbb{P}} l$

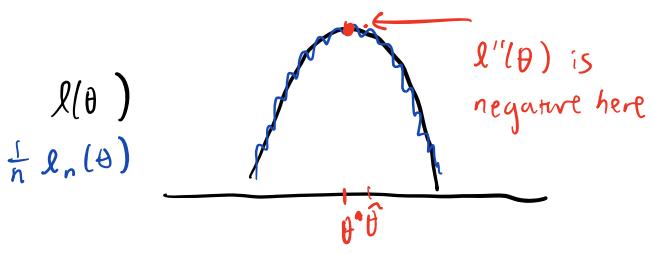
$l_n(\theta) = \sum_{i=1}^n \log f_\theta(X_i)$  = sll  $l(\theta) = \mathbb{E}_{\theta^*} [f_\theta(x)] = \rho l l$   
 max at  $\theta = \hat{\theta}^{\text{MLE}}$  max at  $\theta = \theta^*$

if model is regular.

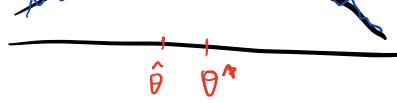
$$\sqrt{n}(\hat{\theta} - \theta^*) \rightsquigarrow N(0, I(\theta^*)^{-1})$$

$$\hat{\theta} \sim N(\theta^*, \frac{1}{n I(\theta^*)})$$

intuition for  $I(\theta^*) = -\ell''(\theta^*)$



$\text{Var}(\hat{\theta})$  is larger when  $\ell$  is flatter  $\Leftrightarrow -\ell''(\theta^*)$  is smaller



building CI's. 95% CI would be  $\theta^* \in \hat{\theta} \pm \frac{1.96}{\sqrt{n I(\hat{\theta})}}$

$$\text{ex) Ber}(p) \rightarrow \theta^* \in \hat{\theta} \pm \frac{1.96}{\sqrt{n}} \cdot \sqrt{\hat{\theta}(1-\hat{\theta})}$$

-if MLE is output of algorithm, we have no closed form.

But we can still reason about it.

### Algorithms to compute MLE

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \ell_n(\theta) \Rightarrow \text{optimization problem}$$

gradient ascent

$$\theta_{j+1} = \theta_j + \gamma \nabla \ell_n(\theta_j)$$

$\uparrow$   
unique parameter

$$\rightarrow \theta_{j+1} = \theta_j - \frac{f'(\theta_j)}{f''(\theta_j)}$$

$\hookrightarrow$  gradient ascent with  $\gamma = -\frac{1}{f''(\theta_j)}$

$$\theta_{j+1} = \theta_j - \nabla^2 \ell_n(\theta_j)^{-1} \nabla \ell_n(\theta_j)$$

Newton-Raphson

$\hookrightarrow$  much more computationally expensive than gradient ascent.

$$0 = f'(\theta) \approx f'(\theta_0) + f''(\theta_0)(\theta - \theta_0)$$

$\uparrow$  new     $\uparrow$  old

### Expectation Maximization (EM) Algorithm

usually to find MLE for a mixture:  $f(x) = (1-p)f_0(x) + p f_1(x)$

can write  $Z = \text{Ber}(p)$ ,  $X_0 \sim f_0$ ,  $X_1 \sim f_1$ ,  $\gamma = \begin{cases} X_0 & \text{if } Z=0 \\ X_1 & \text{if } Z=1 \end{cases}$

ex)  $f_0 = N(\mu_0, \sigma_0^2)$   $f_1 = N(\mu_1, \sigma_1^2)$

5 unknown parameters  $\theta = (p, \mu_0, \sigma_0^2, \mu_1, \sigma_1^2)$

$$\ell_n(\theta) = \sum_{i=1}^n \log \left( (1-p) \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(x-\mu_0)^2}{2\sigma_0^2}} + (1-p) \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} \right)$$

hopeless to find  $\log$  of sum of exponentials!

EM! assume  $p = \frac{1}{2}$ ,  $\sigma_0^2 = \sigma_1^2 = 1$   $\rightarrow f_0 = N(\mu_0, 1)$

$f_1 = N(\mu_1, 1)$

$$\begin{aligned} l_n^{\text{full}}(\theta) &= \sum_{i=1}^n \log(f_0(x_i)^{1-\hat{z}_i} f_1(x_i)^{\hat{z}_i}) \\ &= \underbrace{\sum_{i=1}^n (1-\hat{z}_i) \log(f_0(x_i))}_{\text{Mo param}} + \underbrace{\hat{z}_i \log(f_1(x_i))}_{\text{M1 param}} \end{aligned}$$

LECTURE 11 3/4/24 1 PM

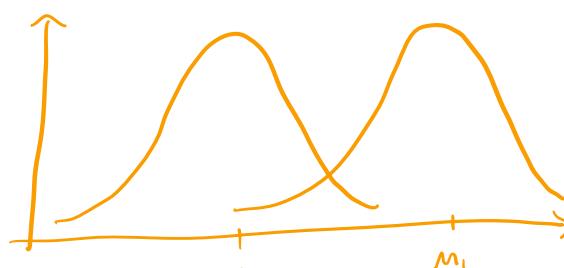
$$\hat{\mu}_0 = \frac{\sum_{i: Z_i=0} Y_i}{\#\{i: Z_i=0\}} = \frac{\sum_{i=1}^n (1-Z_i) Y_i}{\sum_{i=1}^n (1-Z_i)} \quad \hat{\mu}_1 = \frac{\sum_{i: Z_i=1} Y_i}{\#\{i: Z_i=1\}} = \frac{\sum_{i=1}^n Z_i Y_i}{\sum_{i=1}^n Z_i}$$

Main idea of EM: estimate  $Z_i$  with  $\hat{Z}_i \in [0, 1]$  which tells how strongly we believe  $Z_i=0$  or 1.

then, maximize  $l_n^{\text{full}}$  over  $\mu_0, \mu_1$ .

$$\begin{aligned} \hat{l}_n^{\text{full}}(\theta) &= \sum_i (1-\hat{z}_i) \left( \log \frac{1}{\sqrt{2\pi}} - \frac{(Y_i - \mu_0)^2}{2} \right) + \hat{z}_i \left( \log \frac{1}{\sqrt{2\pi}} - \frac{(Y_i - \mu_1)^2}{2} \right) \\ &= -\sum_i (1-\hat{z}_i) \frac{1}{2}(Y_i - \mu_0)^2 - \sum_i (\hat{z}_i) \frac{1}{2}(Y_i - \mu_1)^2 + \text{const.} \end{aligned}$$

how to get  $\hat{z}_i$



$$\hat{z}_i = \mathbb{E}[Z_i | Y_i]$$

$$= P[Z_i=1 | Y_i]$$

$$= \frac{f(Y_i | Z_i=1) P(Z_i=1) \underset{1/2}{\cancel{Y_i}}}{f(Y_i | Z_i=0) P(Z_i=0) + f(Y_i | Z_i=1) \underset{1/2}{\cancel{P(Z_i=1)}}} = \frac{f(Y_i | Z_i=1)}{f(Y_i | Z_i=0) + f(Y_i | Z_i=1)}$$

rename

$$= \frac{f_1(Y_i)}{f_0(Y_i) + f_1(Y_i)} = \frac{e^{-(Y_i - \mu_1)^2/2}}{e^{-(Y_i - \mu_0)^2/2} + e^{-(Y_i - \mu_1)^2/2}} \approx \frac{1}{0+1} = 1 \text{ when } Y_i \approx 1$$

$$\begin{array}{ccc}
 \hat{\mu}_0^{\text{old}}, \hat{\mu}_1^{\text{old}} & \longrightarrow & \hat{z}_1, \dots, \hat{z}_n \\
 & & \hat{z} = \mathbb{E}[z_i | Y_i] \\
 & & \text{E Step} \\
 & & \hat{\mu}_0^{\text{new}}, \hat{\mu}_1^{\text{new}} \\
 & & = \underset{\mu_0, \mu_1}{\operatorname{argmax}} \hat{l}_n^{\text{full}}(\mu_0, \mu_1) \\
 & & \text{M step}
 \end{array}$$

**EM Algorithm**  $P = \frac{1}{2}$ , two normals with  $\sigma_0^2 = \sigma_1^2 = 1$ .

initialize  $\hat{\mu}_0^{(0)}, \hat{\mu}_1^{(0)}$

for  $j=0, 1, 2, \dots$

$$\text{E step. } \hat{z}_i^{(j+1)} = \mathbb{E}_{\hat{\mu}_0^{(j)}, \hat{\mu}_1^{(j)}}[z_i | Y_i] = \frac{e^{-\frac{1}{2}(Y_i - \hat{\mu}_1^{(j)})^2}}{e^{-\frac{1}{2}(Y_i - \hat{\mu}_1^{(j)})^2} + e^{-\frac{1}{2}(Y_i - \hat{\mu}_0^{(j)})^2}}$$

$$\text{M step. } \hat{\mu}_0^{(j+1)}, \hat{\mu}_1^{(j+1)} = \underset{\mu_0, \mu_1}{\operatorname{argmax}} \hat{l}_n^{\text{full}}(\mu_0, \mu_1)$$

$$\hat{\mu}_0^{(j+1)} = \frac{\sum_{i=1}^n (1 - \hat{z}_i^{(j+1)}) Y_i}{\sum_{i=1}^n (1 - \hat{z}_i^{(j+1)})} \quad \hat{\mu}_1^{(j+1)} = \frac{\sum_{i=1}^n \hat{z}_i^{(j+1)} Y_i}{\sum_{i=1}^n \hat{z}_i^{(j+1)}}$$

need not be from same family

**Generalized EM** unknown  $p$ ,  $f(y) = (1-p)f_{\theta_0}(y) + p g_{\theta_1}(y)$

initialize  $\hat{\theta}_0^{(0)}, \hat{\theta}_1^{(0)}, p^{(0)}$

for  $j=0, 1, 2, \dots$

$$\text{E step. } \hat{z}_i^{(j+1)} = \frac{g_{\theta_1^{(j)}}(Y_i) p^{(j)}}{g_{\theta_1^{(j)}}(Y_i) p^{(j)} + f_{\theta_0^{(j)}}(Y_i)(1-p^{(j)})}$$

$$p^{(j+1)} = \frac{1}{n} \sum_{i=1}^n \hat{z}_i^{(j+1)}$$

if log 0,  
try these

else

$$\text{M step. } \theta_0^{(j+1)} = \underset{\theta_0}{\operatorname{argmax}} \prod_{i=1}^n f_{\theta_0}(Y_i)^{1 - \hat{z}_i^{(j+1)}} \quad \sum_{i=1}^n (1 - \hat{z}_i^{(j+1)}) \log f_{\theta_0}(Y_i)$$

$$\theta_1^{(j+1)} = \underset{\theta_1}{\operatorname{argmax}} \prod_{i=1}^n g_{\theta_1}(Y_i)^{\hat{z}_i^{(j+1)}} \quad \sum_{i=1}^n \hat{z}_i^{(j+1)} \log g_{\theta_1}(Y_i)$$

## Method of Moments

$X_1, X_2, \dots, X_n \sim P_\theta$ . Want to recover  $\theta$  from estimated moments

Def  $\alpha_k(\theta) = E_\theta(X^k)$  is the  $k^{\text{th}}$  moment

(1) Work out the moments as a function of  $\theta$   
 → find a function for  $\alpha_k(\theta)$  for  $k=1, 2, \dots$

(2) Estimate moments via plug-in rule

$$\rightarrow \hat{\alpha}_k = E[X^k] \rightsquigarrow \frac{1}{n} \sum_{i=1}^n \bar{x}_i^k = \hat{\alpha}_k$$

(3) Solve systems for  $\theta$

$$\rightarrow \begin{cases} \alpha_1(\theta) = \hat{\alpha}_1 \\ \vdots \\ \alpha_k(\theta) = \hat{\alpha}_k \end{cases}$$

ex)  $P_\theta(\mu, \sigma^2)$

$$(1) \quad \alpha_1(\theta) = \mu \quad \alpha_2(\theta) = \mu^2 + \sigma^2$$

$$(2) \quad \hat{\alpha}_1 = \bar{x}_n \quad \hat{\alpha}_2 = \frac{1}{n} \sum_{i=1}^n \bar{x}_i^2$$

$$(3) \quad \text{solve} \quad \begin{cases} \mu = \bar{x}_n \\ \mu^2 + \sigma^2 = \frac{1}{n} \sum_{i=1}^n \bar{x}_i^2 \end{cases}$$

Method of Moments also gives

$$\hat{\theta}_n \xrightarrow{P} \theta \quad \text{as} \quad n \rightarrow \infty \quad \sqrt{n}(\hat{\theta}_n - \theta) \rightsquigarrow N(0, \Sigma)$$

Bootstrap Brad Efron 1979

→ simulation-based method to approximate the variance of an estimator  $\hat{\theta}(x_1, \dots, x_n)$  and ultimately build confidence intervals

so far

$$\theta \in \hat{\theta} \pm 1.96 \hat{s.e.} \leftarrow \frac{\hat{\sigma}}{\sqrt{n}}$$

↑  
based on  
asympt. normality

- CLT (if  $\hat{\theta} = \bar{x}_n$ )
- Delta Method (if  $\hat{\theta} = g(\bar{x}_n)$ )
- Fisher info (if  $\hat{\theta} = \hat{\theta}^{\text{MLE}}$ )

but what if we lack asympt. normality ??

Ex  $P_\lambda = \text{Pois}(\lambda)$



goal: estimate  $\theta = \text{Median}(P_\lambda)$

i.e.  $\theta = \text{integer } k$  s.t.

$$\sum_{j=0}^{k-1} f_\lambda(j) < \frac{1}{2}, \quad \sum_{j=0}^k f_\lambda(j) \geq \frac{1}{2}$$

can try  $\hat{\theta} = g(\bar{x}_n)$ , but can't get  $\hat{s.e.}(\hat{\theta})$  because  $g$  not continuous

$\hat{\theta} = \text{Med}(x_1, \dots, x_n)$  isn't any one of the three forms above

we want: general purpose method of estimating  $s.e.(\hat{\theta})$  without relying on special structure of  $\hat{\theta}$ , using only  $x_i$ 's

Recall  $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$  say  $Y = \hat{\theta}(x_1, \dots, x_n)$ ,  $Y \sim D$

We want to get more samples of  $\hat{\theta}$ , but we've used up all the  $x_i$ 's for one  $\hat{\theta}$ .

## Thought Experiment (too expensive...)

if we could get as many  $X_i$  as we wanted,  
we can make multiple ( $B$ ) sets of  $\hat{\theta}$  estimators.

$$X_{1:n}^{(1)} = \{X_1^{(1)}, \dots, X_n^{(1)}\} \rightarrow \hat{\theta}^{(1)}$$

$$X_{1:n}^{(2)} = \{X_1^{(2)}, \dots, X_n^{(2)}\} \rightarrow \hat{\theta}^{(2)}$$

⋮

$$X_{1:n}^{(B)} = \{X_1^{(B)}, \dots, X_n^{(B)}\} \rightarrow \hat{\theta}^{(B)}$$

histogram of the  $\hat{\theta}^{(b)}$  for  $b=1, \dots, B$  could be revealing -

We can also estimate

$$\widehat{V[\hat{\theta}]} = \frac{1}{B} \sum_{b=1}^B \left( \hat{\theta}^{(b)} - \frac{1}{B} \sum_{c=1}^B \hat{\theta}^{(c)} \right)^2$$

## Better Alternative

sample from

$$\hat{\pi}_n = \text{Uniform}(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$$

"empirical distribution"

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i \leq x) \rightarrow \text{cdf},$$

LLN says  $\hat{F}_n(x) \xrightarrow{P} F(x)$  as  $n \rightarrow \infty$ .

## Bootstrap Sample

Given a sample  $X_1, \dots, X_n \stackrel{iid}{\sim} P$ , a bootstrap sample is

a collection of  $m$  random variables  $X_1^*, \dots, X_m^*$  s.t.

$$X_i^* \stackrel{iid}{\sim} \text{Unif}(X_1, \dots, X_n), \quad i=1, 2, \dots, m$$

obtain

$$\{X_1^{*(1)}, \dots, X_n^{*(1)}\} \rightarrow \hat{\theta}_1^*$$

$$\{X_1^{*(2)}, \dots, X_n^{*(2)}\} \rightarrow \hat{\theta}_2^*$$

⋮

$$\{X_1^{*(B)}, \dots, X_n^{*(B)}\} \rightarrow \hat{\theta}_B^*$$

as before but now uniform sample.

we can use

$$V_{boot} = \frac{1}{B} \sum_{b=1}^B \left( \hat{\theta}_b^* - \frac{1}{B} \sum_{c=1}^B \hat{\theta}_c^* \right)^2 \approx V_{\hat{P}_n}[\hat{\theta}] \approx V_P[\hat{\theta}]$$

$\uparrow$                        $\uparrow$   
OK if                      need large  
 $B$  large                       $n$ , harder

## LECTURE 13 3/8/24 1PM

Bootstrap samples:  $X_1^*, \dots, X_m^* \sim \text{Unif}(X_1, X_2, \dots, X_n)$

$$\{X_1^{*(1)}, \dots, X_n^{*(1)}\} \rightarrow \hat{\theta}_1^*$$

If  $\hat{P} = \text{Unif}(X_1, \dots, X_n)$

$$\{X_1^{*(2)}, \dots, X_n^{*(2)}\} \rightarrow \hat{\theta}_2^*$$

is close to  $P$  then

⋮

$$\{X_1^{*(B)}, \dots, X_n^{*(B)}\} \rightarrow \hat{\theta}_B^*$$

$\hat{\theta}^*$  is close to  $\hat{\theta}$

distribution

mean/variance of  $\hat{P}$

Suppose  $E[X_i] = \mu$ ,  $V[X_i] = \sigma^2$

$X_i^* \stackrel{\text{iid}}{\sim} \text{Unif}(X_1, \dots, X_n) \quad \underbrace{\text{conditionally}}_{\text{on } X_1, \dots, X_n} \text{ on } X_1, \dots, X_n$ , meaning

$P(X_i^* = x_j | X_1, \dots, X_n) = \frac{1}{n} \quad \text{for all } j$

unconditional distribution of  $X_i^*$  is  $P$  itself

conditional mean

$$\mathbb{E}[X_i^* | X_1, \dots, X_n] = \sum_{j=1}^n X_j \mathbb{P}(X_i^* = X_j | X_1, X_2, \dots, X_n) = \frac{1}{n} \sum_{j=1}^n X_j = \bar{X}_n$$

Sample mean

unconditional mean

$$\mathbb{E}[X_i^*] = \mathbb{E}[\mathbb{E}[X_i^* | X_1, \dots, X_n]] = \mathbb{E}[\bar{X}_n] = \mu$$

tower  
property

true  
mean

conditional variance

$$\begin{aligned} \mathbb{V}[X_i^* | X_1, \dots, X_n] &= \mathbb{E}[X_i^{*2} | X_1, \dots, X_n] - \mathbb{E}[X_i^* | X_1, \dots, X_n]^2 \\ &= \frac{1}{n} \sum_{j=1}^n X_j^2 - \left( \frac{1}{n} \sum_{j=1}^n X_j \right)^2 = \frac{1}{n} \sum_{j=1}^n (X_j - \bar{X}_n)^2 = \hat{\sigma}^2 \end{aligned}$$

Sample  
variance

unconditional variance

$$\begin{aligned} \mathbb{V}[X_i^* | X_1, \dots, X_n] &= \mathbb{E}[\mathbb{V}[X_i^* | X_1, \dots, X_n]] + \mathbb{V}[\mathbb{E}[X_i^* | X_1, \dots, X_n]] \\ &\stackrel{\text{law of total variance}}{=} \mathbb{E}[\hat{\sigma}^2] + \mathbb{V}[\bar{X}_n] = \frac{n-1}{n} \sigma^2 + \frac{\sigma^2}{n} = \sigma^2 \\ \mathbb{E}\left[\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2\right] &= \sigma^2 \quad \underbrace{\qquad\qquad\qquad}_{\text{true variance}} \end{aligned}$$

## Bootstrap CIs

(1) normal interval  $\hat{\theta} \pm 1.96 \sqrt{v_{\text{boot}}} \leftarrow \text{defeats purpose of bootstrap} \therefore$

(2) pivotal interval

→ let  $\theta = \text{ground truth (number)}$

→ let  $\hat{\theta} = \text{random variable (estimator of } \theta)$

→ let  $\hat{\theta}_0 = \text{single sample of } \hat{\theta} \text{ (number)}$

our CI will look like  $\theta \in (\hat{\theta} - r_1, \hat{\theta} + r_2)$

use  $\hat{\theta}_0$  to estimate  $\hat{\theta} \rightarrow \theta \in (\hat{\theta}_0 - r_1, \hat{\theta}_0 + r_2)$

① uses  $r_2 = r_1 = \hat{\sigma}_{\hat{\theta}} \text{ se}(\hat{\theta})$

we want

$$1 - \alpha = \mathbb{P}(\hat{\theta}_+ - r_1 \leq \theta \leq \hat{\theta}_+ + r_2) = \mathbb{P}(\underbrace{\theta - r_2 \leq \hat{\theta}}_{\approx \alpha/2 \text{ mass}} \leq \hat{\theta} \leq \underbrace{\theta + r_1}_{\approx \alpha/2 \text{ mass}})$$

recall  $\mathbb{P}(\hat{\theta} \leq q_{1-\alpha/2}) = \alpha/2$

$$\mathbb{P}(\hat{\theta} \geq q_{\alpha/2}) = \alpha/2$$

$q_x \rightarrow x \text{ of data is greater}$

take  $r_2 = \theta - q_{1-\alpha/2} \rightsquigarrow \hat{r}_2 = \hat{\theta}_0 - q^*_{1-\alpha/2} \leftarrow \text{quantiles of bootstrap samples}$

$$r_1 = q_{\alpha/2} - \theta \rightsquigarrow \hat{r}_1 = q^*_{\alpha/2} - \hat{\theta}_0$$

$$\hookrightarrow \theta \in (\hat{\theta}_0 - (q^*_{\alpha/2} - \hat{\theta}_0), \hat{\theta}_0 + (q^*_{\alpha/2} - \hat{\theta}_0))$$

$$= (2\hat{\theta}_0 - q^*_{\alpha/2}, 2\hat{\theta}_0 + q^*_{\alpha/2}) \rightarrow \boxed{\text{pivotal CI}}$$

### ③ percentile interval

$$CI = (q^{\alpha/2}, q^*_{1-\alpha/2})$$

$\uparrow \quad \uparrow$   
quantiles of  
bootstrap  
samples

LECTURE 14 3/11/24 1PM

Hypothesis Testing lets us answer binary questions

A test is a function  $\Psi : \text{data} \rightarrow \{0, 1\}$ . In particular,  $\Psi$  is an estimator.

rejection region  $\Rightarrow R = \{ \text{datasets for which } \Psi(\text{data}) = 1 \}$

$$\hookrightarrow \Psi(\text{data}) = 1 \text{ (data } \in R)$$

A test statistic is a function that summarizes the data and is sufficient to compute a test  $\Psi$ .

ex)  $\bar{X}_n, \bar{X}_n^3, \bar{X}_n - 30$  for a test  $\Psi(\text{data}) = 1(\bar{X}_n > 30)$

$\underbrace{\text{all sufficient for}}$

test, but  $\bar{X}_n$  most natural

A hypothesis test takes the form

$$H_0: \theta \in \Theta_0, \quad H_1: \theta \in \Theta_1$$

where  $\Theta_0 \cap \Theta_1 = \emptyset$ ,  $\Theta_0 \cup \Theta_1 = \Theta$  = full parameter space.

$\hookrightarrow H_0$  = innocent OR status quo

$\hookrightarrow H_1$  = guilty OR discovery

	conclude $H_0 (\Psi=0)$	conclude $H_1 (\Psi=1)$
$H_0$ true	✓	Type I Error ✓
$H_1$ true	Type II Error	✓

false positive  
false negative

The size of a test  $\Psi$  is

$$\text{size}(\Psi) = \max_{\theta \in \Theta_0} P_\theta(\Psi=1)$$

$\theta$  should be on boundary between  $\Theta_0, \Theta_1$

$\rightarrow$  max probability of type I error.

The test has level  $\alpha \in (0,1)$  if  $\text{size}(\Psi) \leq \alpha$ .

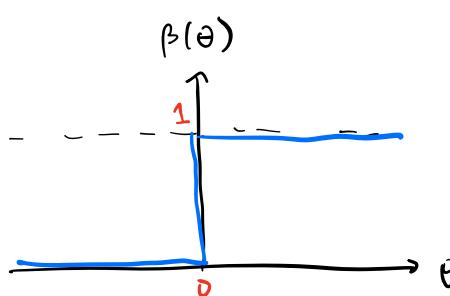
$\rightarrow$  typically use  $\alpha=0.01$  or  $\alpha=0.05$ .

we can get  $\text{size}(\Psi)=0$  if  $\Psi=0$  & datasets but this is not helpful.

The power function is defined as

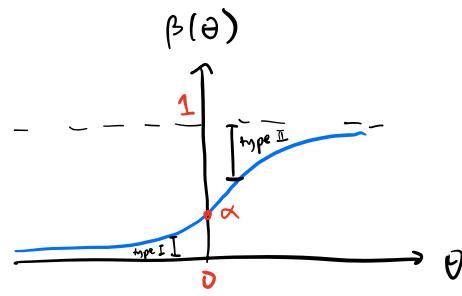
$$\beta(\theta) = P_\theta(\Psi=1)$$

defined for specific parameter set  $\theta$ .



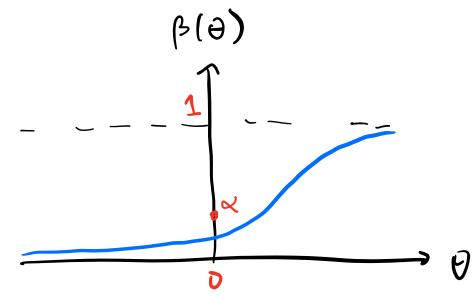
ideal

(but not easy with finite data)



size  $\alpha = \text{level } \alpha$

good



size  $\leq \text{level } \alpha$

bad

Today's construct a test  $\Psi$ . compute power function

ex)  $X_1, \dots, X_n$  iid,  $E[X_i] = \mu$

$$H_0: \mu \leq 30 \quad H_1: \mu > 30.$$

$$\Psi = \mathbb{1} \{ \bar{X}_n - 30 > c_\alpha \}$$

find  $c_\alpha$

$$\alpha = \text{size}(\Psi) = \max_{\mu \leq 30} P_\mu(\bar{X}_n - 30 > c_\alpha)$$

$$\text{but } \sqrt{n}(\bar{X}_n - 30) \rightsquigarrow N(0, \sigma^2) \text{ by CLT so } \bar{X}_n - 30 \approx N(0, \frac{\sigma^2}{n})$$

$$\Rightarrow \alpha = P \left( \frac{\sigma}{\sqrt{n}} Z \geq c_\alpha \right) \text{ so } c_\alpha = \frac{\sigma}{\sqrt{n}} z_\alpha \approx \boxed{\frac{\hat{\sigma}}{\sqrt{n}} z_\alpha}$$

$\uparrow$   
max at boundary       $\downarrow$        $\alpha = P(Z \geq \frac{\sqrt{n}}{\hat{\sigma}} c_\alpha = z_\alpha) \rightarrow$

$$\text{so our test is } \Psi = \mathbb{1} \{ \bar{X}_n - 30 > \frac{\hat{\sigma}}{\sqrt{n}} z_\alpha \}. \text{ for } \alpha = 0.05,$$

$$z_\alpha = 1.65$$

$$z_{\alpha/2} = 1.96$$

$$\rightarrow \text{if } n=164, \hat{\sigma}=12, \bar{X}_n = 33.4$$

$$\text{then } \bar{X}_n - 30 = 3.4 > \frac{12}{\sqrt{164}} \cdot 1.65 = 1.54 \quad \text{so } \boxed{\text{reject}}$$

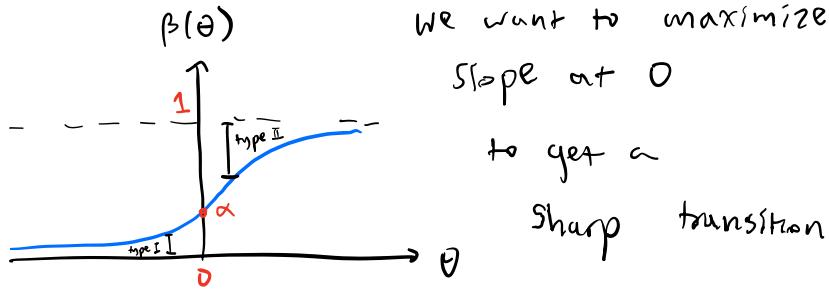
$$\text{Power Function} \quad \bar{X}_n \sim N(\mu, \frac{\hat{\sigma}^2}{n})$$

$$\beta(\mu) = P_\mu \left( \bar{X}_n - 30 > \frac{\hat{\sigma}}{\sqrt{n}} z_\alpha \right)$$

$$\approx P_\mu \left( \mu + \frac{\hat{\sigma}}{\sqrt{n}} Z - 30 > \frac{\hat{\sigma}}{\sqrt{n}} z_\alpha \right) = P \left( Z > z_\alpha + \frac{\sqrt{n}}{\hat{\sigma}} (30 - \mu) \right)$$

$$\begin{aligned} &\text{higher } \mu \rightarrow = 1 - \Phi \left( z_\alpha + \frac{\sqrt{n}}{\hat{\sigma}} (30 - \mu) \right) \\ &\Rightarrow \text{higher } \beta(\mu) \end{aligned}$$

Recall our power function from before.



We want to maximize  
slope at 0  
to get a  
sharp transition

in our example,

$$\beta'(\mu) \approx \frac{\sqrt{n}}{\hat{\sigma}} \phi(z_\alpha + \frac{\sqrt{n}}{\hat{\sigma}}(30-\mu))$$

at  $\mu=30$  (our boundary),

$$\beta'(\mu) = \frac{\sqrt{n}}{\hat{\sigma}} \phi(z_\alpha) \text{ is}$$

maximized when

$\rightarrow n \uparrow \leftarrow$  more data

$\rightarrow \hat{\sigma} \downarrow \leftarrow$  less variance

### Wald's Test

Suppose  $\frac{\hat{\theta} - \theta}{\hat{s}\hat{\theta}} \rightsquigarrow N(0, 1)$  as  $n \rightarrow \infty$

Case 1.  $\theta = \theta_0$  or  $\theta \neq \theta_0$

$$\Psi = \mathbb{1} \left( \left| \frac{\hat{\theta} - \theta}{\hat{s}\hat{\theta}} \right| > z_{\alpha/2} \right)$$

Case 2.  $\theta \leq \theta_0$  or  $\theta > \theta_0$

$$\Psi = \mathbb{1} \left( \frac{\hat{\theta} - \theta}{\hat{s}\hat{\theta}} > z_\alpha \right)$$

Case 3.  $\theta \geq \theta_0$  or  $\theta < \theta_0$

$$\Psi = \mathbb{1} \left( \frac{\hat{\theta} - \theta}{\hat{s}\hat{\theta}} < -z_\alpha \right)$$

ex  $\sqrt{n} (\hat{\theta}^{\text{MLE}} - \theta) \rightsquigarrow N(0, I^{-1}(\theta))$

so  $\text{se}(\hat{\theta}^{\text{MLE}}) \approx \frac{1}{\sqrt{n}I(\theta)}$  or can use  $\frac{1}{\sqrt{n}I^{-1}(\theta_0)}$

$\theta = \theta_0$   
or  
 $\theta \neq \theta_0$

$$\Psi = \mathbb{1} \left( \sqrt{n} I(\hat{\theta}^{\text{MLE}}) |\hat{\theta}^{\text{MLE}} - \theta| > z_{\alpha/2} \right)$$

$\uparrow$   
or  
 $I(\theta_0)$

p-values answer "How close was the call between  $H_0$  and  $H_1$ ?"

small p-value gives evidence against  $H_0$ .

$$\text{evidence against } H_0 = \begin{cases} \text{very strong} & p < 0.01 \\ \text{strong} & 0.01 < p < 0.05 \\ \text{weak} & 0.05 < p < 0.1 \\ \text{little or none} & 0.1 < p \end{cases}$$

The p-value is

- smallest level  $\alpha$  for which we can reject  $H_0$
- if  $T_n$  is test statistic and  $T_n^{\text{obs}}$  is observed value of  $T_n$ , rejection regions are

$$R = \{T_n > c_\alpha\} \Rightarrow p = \sup_{\theta \in \Theta_0} P_\theta(T_n > T_n^{\text{obs}})$$

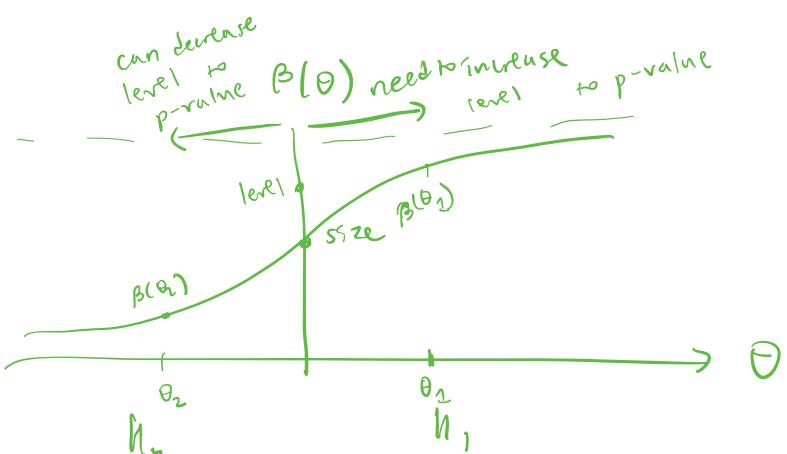
$$R = \{T_n < c_\alpha\} \Rightarrow p = \sup_{\theta \in \Theta_0} P_\theta(T_n < T_n^{\text{obs}})$$

$$R = \{|T_n| > c_\alpha\} \Rightarrow p = \sup_{\theta \in \Theta_0} P_\theta(|T_n| > T_n^{\text{obs}})$$

if  $T_n$  is standardized as  $N(0, 1)$  then (for first kind)

$$p = \sup_{\theta \in \Theta_0} P_\theta(N(0, 1) > T_n^{\text{obs}})$$

"largest probability of observing a more extreme value than  $T_n^{\text{obs}}$  over all  $\theta \in \Theta_0$ "



example  $H_0: p = \frac{1}{2}$   $H_1: p \neq \frac{1}{2}$   $n = 124$   $\bar{x}_n = 0.645$

Ber(p).

$$\frac{\bar{x}_n - p}{\sigma} \sim N(0, 1)$$

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

$$\rightarrow N(p, \frac{p(1-p)}{n})$$

A test of level  $\alpha$  is

$$\begin{aligned} \textcircled{(a)} \quad \Psi &= \mathbb{1} \left( \sqrt{\frac{n}{\bar{x}_n(1-\bar{x}_n)}} |\bar{x}_n - 0.5| > z_{\alpha/2} \right) \\ &= \mathbb{1} \left( 3.37 > z_{\alpha/2} \right) \end{aligned}$$

can have  $\alpha/2 \geq 0.0004$   
 $\alpha \geq 0.0008 \leftarrow p\text{-value.}$

$$\begin{aligned} &\text{under } \bar{x}_n \sim p \approx N\left(\frac{1}{2}, \frac{\bar{x}_n(1-\bar{x}_n)}{n}\right) \textcircled{(a)} \\ &\text{under } \theta \in \Theta_0 \approx N\left(\frac{1}{2}, \frac{1}{4n}\right) \textcircled{(b)} \end{aligned}$$

$$\begin{aligned} \textcircled{(b)} \quad \Psi &= \mathbb{1} \left( \sqrt{4n} |\bar{x}_n - 0.5| > z_{\alpha/2} \right) \\ &= \mathbb{1} \left( 3.23 > z_{\alpha/2} \right) \end{aligned}$$

can have  $\alpha/2 \geq 0.0006$   
 $\rightarrow p\text{-value} = 0.0012$

either works!

example  $H_0: \mu \leq 30$   $H_1: \mu > 30$   $n = 164$   $\bar{x}_n = 33.4$   $\hat{\sigma} = 12$

$$\bar{x}_n \sim N(30, \frac{12^2}{n}) \text{ under } H_0$$

$$\begin{aligned} \Psi &= \mathbb{1} \left( \frac{\sqrt{n}}{12} (\bar{x}_n - 30) > z_\alpha \right) \\ &= \mathbb{1} (3.63 > z_\alpha) \rightarrow p\text{-value} = 0.0002 \end{aligned}$$

$X \sim \chi^2_k$  is a chi-squared distribution with  $k$  deg. of freedom

$$X := Z_1^2 + Z_2^2 + \cdots + Z_k^2, \quad Z_i \sim N(0, 1).$$

$$\mathbb{E}[X] = \mathbb{E}[Z_1^2 + Z_2^2 + \cdots + Z_k^2] = k$$

$$\mathbb{V}[X] = \sum_{i=1}^k \mathbb{V}[Z_i^2] = \sum_{i=1}^k \mathbb{E}[Z_i^4] - \mathbb{E}[Z_i^2]^2 = \sum_{i=1}^k 2 = 2k$$

Let  $\chi_{k,\alpha}^2$  be the  $\alpha^{th}$  quantile of  $\chi^2_k$ .

$$\hookrightarrow P(X > \chi_{k,\alpha}^2) = \alpha, \quad X \sim \chi^2_k$$

Suppose we have a pmf  $f_0$  and want to test if a discrete r.v. has pmf  $f_0$ .

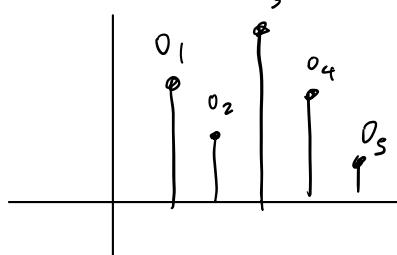
### $\chi^2$ Goodness-of-Fit Test

$X$  discrete, takes  $k$  values  $\{1, 2, \dots, k\}$

$$H_0 : f = f_0 \quad H_1 : f \neq f_0$$

Also read on  
 $\chi^2$  independence test!

for  $X_1, X_2, \dots, X_n \sim f$



$$O_j = \#\{i : X_i = j\} = \sum_{r=1}^n \mathbf{1}(X_r = j)$$

$$\text{expect } O_j \approx n f_0(j)$$

$$\hookrightarrow \text{under null, } n f_0(j) = E_j$$

$O_j$  = observed # of  $j$ 's

$E_j$  = expected # of  $j$ 's under null

$$T = \sum_{j=1}^k \frac{(O_j - E_j)^2}{E_j}$$

Reject if  $T > c_\alpha = \chi^2_{k-1, \alpha}$

$$\text{p-value} = P_{H_0}(T > T_{obs}) = P(\chi^2_{k-1} > T_{obs})$$

Theorem  $T \rightsquigarrow \chi^2_{k-1}$  as  $n \rightarrow \infty$

$\hookrightarrow k-1$  degrees of freedom since  $0_1 + 0_2 + \dots + 0_k = n$ .

Ex 100 surveys

	exp.	obs.
cliff	50	48
happy	30	40
sad	20	12

$$T = \frac{2^2}{50} + \frac{10^2}{30} + \frac{8^2}{20} = 6.61$$

$$\chi^2_{k-1, \alpha} = \chi^2_{2, 0.05} = 5.991 \rightarrow \text{reject}$$

Ex Bernoulli,  $H_0: p = \frac{1}{2}$ ,  $H_1: p \neq \frac{1}{2} \rightarrow$  a pmf between  $\{1, 2\}$ .  
 $n = 124$

	exp.	obs.
right	62	80
left	62	44

$$T = \frac{18^2 + 18^2}{62} = 10.46$$

$$\begin{aligned} \text{p-value} &= P(\chi^2_1 = Z^2 > 10.46) \\ &= P(|Z| > 3.23) \end{aligned}$$

same as  
Wald's Test?

LECTURE 18 4/1/24 (Did not attend)

### Nonparametric Tests

$X \sim N(0, 1)$

vs.

$X \sim$  any other distribution

} much more  
than just one  
"family" of  
distributions

$X \sim$  Normal

vs.

$X \sim$  not normal

$X \sim F_0$

vs.

$X \neq F_0$

### Kolmogorov-Smirnov

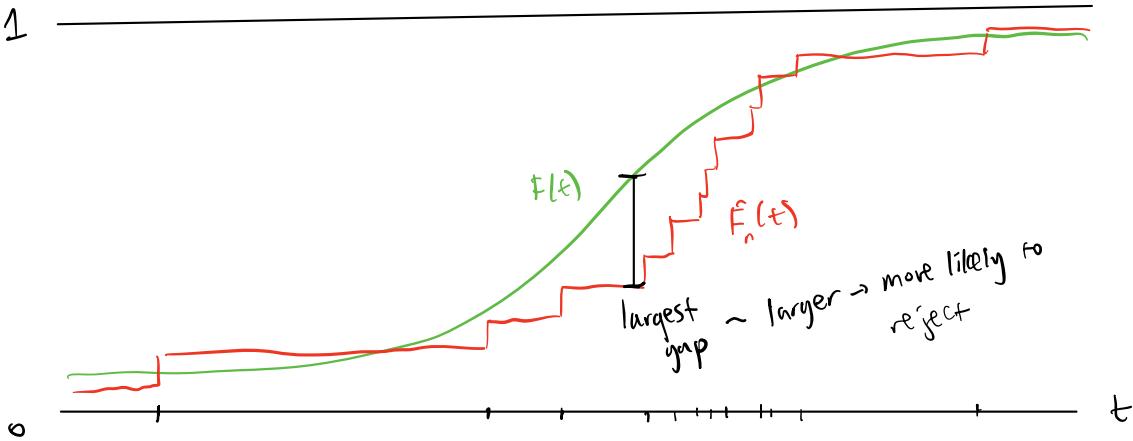
given  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} F$ , test  $H_0: F = F_0$  vs.  $H_1: F \neq F_0$

empirical cdf  $\hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i \leq t)$

•  $E[\hat{F}_n(t)] = F(t) \rightarrow$  unbiased estimator

•  $n\hat{F}_n(t) \sim \text{Bin}(n, F(t)) \rightarrow$  average of Bernoulli r.v.s

$\hookrightarrow \sqrt{n}(\hat{F}_n(t) - F(t)) \rightsquigarrow N(0, F(t)(1-F(t)))$



Kolmogorov - Smirnov Test

$$\Psi = \mathbb{1} \left( \sup_t |\hat{F}_n(t) - F(t)| > c_\alpha \right)$$

$\rightarrow c_\alpha$  doesn't depend on  $F_0$ , only on  $n!!$

### Kolmogorov - Lilliefors

for  $H_0: F$  is Gaussian vs.  $H_1: F$  is not Gaussian

using  $T = \sup_t |\hat{F}_n(t) - \Phi_{\hat{\mu}, \hat{\sigma}^2}(t)|$

requires us to use KL distribution, NOT KS.

$\downarrow$   
 $H_0 = \text{any Gaussian}$

$\downarrow$   
 $H_0 = \text{specific Gaussian}$

Note  $c_\alpha^{KL} < c_\alpha^{KS}$

### Permutation Test "2-sample test"

Given  $X_1, \dots, X_n \sim F_X$        $Y_1, \dots, Y_m \sim F_Y$

Test  $H_0: F_X = F_Y$  vs.  $H_1: F_X \neq F_Y$

can't use  $T = \sup_t |\hat{F}_X(t) - \hat{F}_Y(t)|$  and  $\Psi = \mathbb{1}(T > c_\alpha)$   
since  $F_X$  and  $c_\alpha$  are unknown

instead use  $T = |\bar{X}_n - \bar{Y}_n|$ ,  $Z_i = X_i$ ,  $Z_{i+n} = Y_i$  for  $1 \leq i \leq n$

let  $m=n$  for simplicity compute  $T^1, T^2, \dots, T^B$  by picking random permutation  $\epsilon_j$  of  $\{1, 2, \dots, 2n\}$ , then

$$Z_{\text{left}}^{\sigma_j} = \frac{1}{n} \sum_{i=1}^n Z_{\sigma_j(i)} \quad Z_{\text{right}}^{\sigma_j} = \frac{1}{n} \sum_{i=n+1}^{2n} Z_{\sigma_j(i)} \quad T^j = |Z_{\text{left}}^{\sigma_j} - Z_{\text{right}}^{\sigma_j}|$$

reject if  $T > t_\alpha$  where  $\frac{1}{B} \sum_{i=1}^B \mathbb{1}(T^i > t_\alpha) = \alpha$

"only  $\alpha$  of  
the permutations  
gave higher differences"

$$\text{p-value} = \frac{1}{B} \sum_{i=1}^B \mathbb{1}(T^i > T)$$

note that type II error can be large, especially if  $E[X] = E[Y]$

LECTURE 19 4/3/24 (Did not attend)

### Multiple Hypothesis Testing

↳ you will inevitably falsely reject some tests if you perform many.

① Bonferroni Correction → very conservative

$m$  tests, reject each at level  $\alpha/m$

$$\begin{aligned} \Pr_{H_0}(\exists i, H_{0,i} \text{ is rejected}) &= \Pr_{H_0}\left(\bigcup_{i=1}^m \{p_i \leq \alpha/m\}\right) \\ &\leq \sum_{i=1}^m \Pr_{H_0}\{p_i \leq \alpha/m\} = m \cdot \alpha/m = \alpha. \end{aligned}$$

② Benjamini-Hochberg (BH) Method

$$\text{FDP} = \frac{\#\text{false positives}}{\#\text{positives}}$$

instead of  $\exists$  false positive

$$\text{FDR} = E_{H_0}[\text{FDP}]$$

↑  
can approximate/bound

a) order your  $m$  p-values  $p_{(1)} \leq \dots \leq p_{(m)}$

b) Let  $i_{\max} = \max_{i \in [m]} : p_{(i)} \leq \frac{\alpha}{m} i$

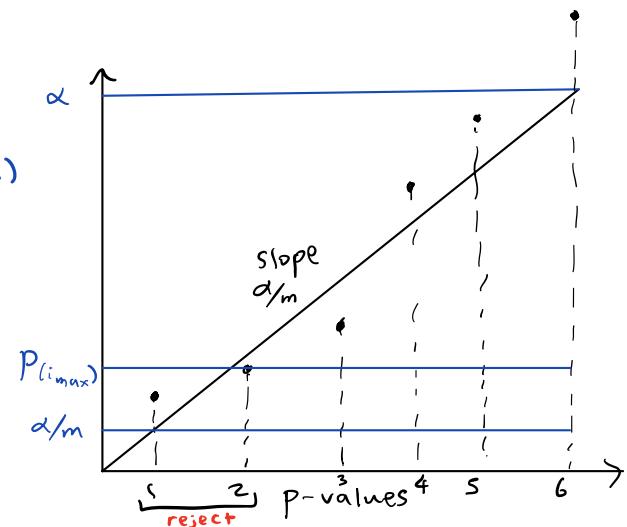
c) Reject all tests s.t.  $p_{(i)} \leq p_{(i_{\max})}$

### Summary

usually reject for  $\leq \alpha$  (5)

② BH  $\leq p_{(i_{\max})}$  (2)

① Bonferroni  $\leq \alpha/m$  (0)



## T-test

Wald Test had asymptotically normal data.

T-test  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$  so  $\sqrt{n}(\bar{X}_n - \mu)/\hat{\sigma}$  is exactly  $N(0, 1)$

↳ but NOT  $\sqrt{n}(\bar{X}_n - \mu)/\hat{\sigma}$ , especially for small n.

$$\sqrt{n}(\bar{X}_n - \mu)/\hat{\sigma} \sim \frac{\Gamma(\frac{v+1}{2})}{\sqrt{v\pi}\Gamma(\frac{v}{2})} \left(1 + \frac{x^2}{v}\right)^{-\frac{v+1}{2}}$$

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

Student t-distribution  
with  $v = n-1$  degrees of freedom

## Student's t-test

Given  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ,  $\Psi = \mathbb{1} \left( \sqrt{n} \frac{|\bar{X}_n - \mu_0|}{\hat{\sigma}} \geq t_{n-1, \alpha/2} \right)$

$H_0: \mu = \mu_0$      $H_1: \mu \neq \mu_0$     p-value =  $P(|t_{n-1}| \geq \sqrt{n} \frac{|\bar{X}_n - \mu_0|}{\hat{\sigma}})$

## LECTURE 20 4/5/24 1PM

Bayesian vs. Frequentist

alternative way to produce estimators, build CIs, do hypothesis testing, based on "degree of belief"

Frequentist

repeat experiment many times → based on frequency of things happening

Bayesian

weight likelihood  $L_n(\theta)$  via prior pdf  $f(\theta)$   
"belief before seeing data"  $X_1, \dots, X_n$ "

### Bayesian Method

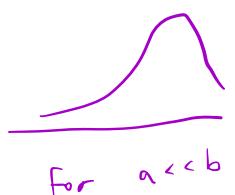
- start with prior belief about  $\theta^*$
- observe data  $X_1, \dots, X_n$
- update prior belief into a posterior belief

Prior Belief (probability distribution on unknown  $\theta^*$ )

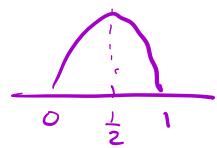
ex) assume  $X_1, \dots, X_n \stackrel{iid}{\sim} Ber(p^*)$

↳ convenient prior is  $Beta(a, b) = \frac{1}{K} p^{a-1} (1-p)^{b-1}$

↙ proportional constant



to determine if  $p^* = \frac{1}{2}$  we use  $\text{Beta}(2,2) = 6 p(1-p)$



### Prior to Posterior

The posterior is the probability of  $\theta$  given the data

$$\text{Recall } P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$

$$\hookrightarrow P(\theta | \text{data}) = \frac{P(\text{data} | \theta) P(\theta)}{P(\text{data})} \quad \begin{matrix} \text{likelihood} \\ \downarrow \\ P(\text{data} | \theta) \end{matrix} \quad \begin{matrix} \text{prior} \\ \swarrow \\ P(\theta) \end{matrix}$$

$P(\text{data}) \leftarrow \text{proportionalizing constant}$

$$\hookrightarrow f(\theta | x_1, \dots, x_n) = \frac{f(x_1, \dots, x_n | \theta) f(\theta)}{f(x_1, \dots, x_n)} = \frac{f(x_1, \dots, x_n | \theta) f(\theta)}{C_n}$$

$\uparrow \text{no dependence on } \theta \quad \uparrow \text{to make whole thing integrate to 1}$

just write  $f(\theta | x_1, \dots, x_n) \propto L_n(\theta) f(\theta)$

| Posterior      likelihood \* prior

back to example,

$$L_n(p) = p^{\sum x_i} (1-p)^{n-\sum x_i} \quad f(p) = 6 p(1-p)$$

$$\Rightarrow f(p | x_1, \dots, x_n) \propto p^{\sum x_i + 1} (1-p)^{n - \sum x_i + 1}$$

$$\hookrightarrow \text{Beta}(\sum x_i + 2, n - \sum x_i + 2)$$

if prior/posterior are from same family of distributions, we have  
a conjugate prior (e.g. Beta from above)

$$\text{Bayes estimator} := \mu_{\text{posterior}} = \mathbb{E}[f(\theta | x_1, \dots, x_n)]$$

$$\text{Maximum A Posteriori} = \text{MAP} := \text{mode}_{\text{posterior}} = \arg \max_{\theta} L_n(\theta) f(\theta)$$

Markov Chain Monte Carlo to estimate  $\hat{\theta}_{\text{posterior}} = \hat{\theta}^{\text{Bayes}}$

$$= \int \theta f(\theta | x_1, \dots, x_n) d\theta$$

① draw  $\theta_1, \dots, \theta_T$  iid from

$f(\theta | x_1, \dots, x_n)$  by simulating on Markov Chain

② compute  $\hat{\theta} = \frac{1}{T} \sum_{t=1}^T \theta_t$  (Monte Carlo)

$\hat{\theta}^{\text{MAP}} = \text{mode}_{\text{posterior}} \rightarrow \hat{\theta}^{\text{MAP}}$ .

Find usually using  $\arg \max_{\theta} (\log(f(\theta | x_1, \dots, x_n)))$

$$\begin{aligned} \log(f(\theta | x_1, \dots, x_n)) &= \log(c_n L_n(\theta) f(\theta)) \\ &= \log c_n + \log L_n(\theta) + \log f(\theta) \end{aligned}$$

set gradient  $D = \nabla \log L_n(\theta) + \nabla \log f(\theta)$ , solve for  $\theta$ .

↳ gradient ascent if needed:  $\theta^{(k+1)} = \theta^{(k)} + \gamma_k (\nabla \log L_n(\theta^{(k)}) + \nabla \log f(\theta^{(k)}))$

Def  $(1-\alpha)$  Posterior interval

Ans  $[a, b]$  s.t.  $\int_a^b f(\theta | x_1, \dots, x_n) d\theta = 1-\alpha$ .

### Choosing Priors

① True prior knowledge (e.g.  $\geq 0$ )

② convenient calculation (e.g. Beta conjugate for  $\text{Ber}(p)$ , normal conjugate for normal)

③ without prior knowledge, uniform prior



$\hat{\theta}^{\text{MAP}} = \hat{\theta}^{\text{MLE}}$  since  $f(\theta)$  uniform

if you want uniform over  $\mathbb{R}$ , can take  $f(\theta) = 1$

even if  $\int_{-\infty}^{\infty} 1 d\theta$  diverges  $\leadsto$  improper prior

$\hat{\theta}^{\text{MAP}}$  might **not** be equal to  $\hat{\theta}^{\text{Bayes}}$

ex) prior  $f(p) = \text{Unif}([0, 1])$

given  $x_1, \dots, x_n$

$$f(p|x_1, \dots, x_n) \propto p^{\sum x_n} (1-p)^{n-\sum x_n} \propto \text{Beta}(\underbrace{\sum x_n + 1}_a, \underbrace{n - \sum x_n + 1}_b)$$

$$\text{biased } \hat{\theta}^{\text{MAP}} = \frac{a}{a+b} = \frac{\sum x_n + 1}{n+2} \neq \hat{\theta}^{\text{MLE}} = \frac{\sum x_n}{n}$$

Normal prior is a conjugate prior for normal data

if

$$\theta \sim N(0, \sigma^2)$$

$$f(\theta) \propto e^{-\frac{x^2}{2\sigma^2}}$$

given data  $x_1, \dots, x_n \stackrel{iid}{\sim} N(\theta, 1)$ ,

$$\text{we have } f(\theta|x_1, \dots, x_n) \propto e^{-\frac{\theta^2}{2\sigma^2}} e^{-\sum_i \frac{(x_i - \theta)^2}{2}}$$

$$-\frac{\sum x_i^2}{2} = \exp\left(-\frac{1}{2} \left[ \sum_i x_i^2 - 2 \sum_i x_i \theta + n\theta^2 + \frac{\theta^2}{\sigma^2} \right]\right)$$

$$\underset{\text{is constant}}{\cancel{e}} \propto \exp\left(-\frac{1}{2} \left[ (n + \frac{1}{\sigma^2})\theta^2 - 2 \sum x_i \theta \right]\right)$$

$$\frac{1}{2\left(\frac{1}{n+\frac{1}{\sigma^2}}\right)} \left(\frac{\sum x_i}{n+\frac{1}{\sigma^2}}\right)^2 \underset{\text{is constant}}{\cancel{e}} \propto \exp\left(-\frac{1}{2\left(\frac{1}{n+\frac{1}{\sigma^2}}\right)} \left[\theta - \frac{\sum x_i}{n+\frac{1}{\sigma^2}}\right]^2\right)$$

$$\underset{\text{posterior}}{\cancel{e}} \text{posterior} \sim N\left(\frac{\sum x_i}{n+\frac{1}{\sigma^2}}, \frac{1}{n+\frac{1}{\sigma^2}}\right)$$

$$\underset{\text{MAP}}{\cancel{e}} \hat{\theta}^{\text{MAP}} = \hat{\theta}^{\text{Bayes}} = \frac{\sum x_i}{n+\frac{1}{\sigma^2}}$$

## Robust Statistics -

How to summarize data in presence of outliers?

Def A statistic/estimator is robust if it does not change the data even when you modify some of the data a lot.

Def Let  $m = \max$  amount of observations we can make arbitrarily large or small while leaving estimator bounded.

Then  $\frac{m}{n}$  is the breakdown point

ex)

$$\hat{\theta}(x_1, \dots, x_5) = \text{sample median}$$

$$\rightarrow |\hat{\theta}(x_1, x'_2, x_3, x'_4, x_5) - \hat{\theta}(x_1, x_2, x_3, x_4, x_5)| \leq 13 < \infty$$

$$\rightarrow |\hat{\theta}(x_1, x'_2, x'_3, x'_4, x_5) - \hat{\theta}(x_1, x_2, x_3, x_4, x_5)| \rightarrow \infty$$

} breakdown point  $3/5$ .

breakdown point for median for odd  $n$  is  $\frac{(n+1)/2}{n} = \frac{n+1}{2n} \rightarrow \frac{1}{2}$  as  $n \rightarrow \infty$ .

note something like  $\hat{\theta}(x_1, \dots, x_n) = 4$  has breakdown point 1 but isn't really useful. Best we can do is  $\frac{1}{2}$ .

Median as MLE

MLE is median in Laplace distribution  $f_{\mu, b}(x) = \frac{1}{2b} e^{-|x-\mu|/b}$

↳ simplify to  $b=1 \rightarrow f_{\mu, 1}(x) = \frac{1}{2} e^{-|x-\mu|}$

$$\text{then } l_n(\theta) = \sum_{i=1}^n \ln(\frac{1}{2}) - |x_i - \mu| \Rightarrow - \sum_{i=1}^n |x_i - \mu|$$

$$\hat{\mu}^{\text{MLE}} = \arg \min_{\mu} \sum_{i=1}^n |x_i - \mu| \rightarrow \text{median}(x_1, \dots, x_n)$$

## Huber's Contamination Model

For data "contaminated" with outliers.

$$Z \sim \text{Ber}(\epsilon) \quad X_{\text{true}} \sim P_{\theta^*} \quad X_{\text{out}} \sim Q \quad \epsilon \text{ small}$$

Let  $Y = (1-Z)X_{\text{true}} + Z X_{\text{out}}$

ex)  $P_{\theta^*} = N(\theta^*, 1)$  and  $Q$  has pdf  $g$ .

$$f(y) = (1-\epsilon) \frac{1}{\sqrt{2\pi}} e^{-\frac{(y-\theta^*)^2}{2}} + \epsilon g(y)$$

$$\ln(l_n(\theta)) = \sum_{i=1}^n \log \left( (1-\epsilon) \frac{1}{\sqrt{2\pi}} e^{-\frac{(Y_i-\theta^*)^2}{2}} + \epsilon g(Y_i) \right)$$

↳ but we need to know  $g$ :

to fix this, we can maximize over all  $g \in$  family  $Q$ .

quasi MLE  $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \max_{g \in Q} l_n(\theta)$

(uses iid assumption  
for both  $P_{\theta^*}, Q$ )

## Arbitrary Contamination

if we only know there are  $m$  outliers

↳ set  $C = \{i : \text{s.t. } X_i \text{ is outlier}\}$ ,  $|C|=m$ .

Let  $\lambda_n(\theta, C) = -\sum_{i \in C} (X_i - \theta)^2$  ↳ only define on non-outliers  
normal dist.

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \max_{|C|=m} \lambda_n(\theta, C)$$

"trimmed mean"

- takes mean of middle  $n-m$  values

Regression - use features  $X \in \mathbb{R}^k$  to predict response variable  $Y \in \mathbb{R}$

**goal** predict  $Y$  given  $X$  / understand how  $Y$  changes with  $X$

**challenge** for EACH fixed  $X=x$ , there is a whole distribution  $Y|X=x$   
 ↳ might not have exact  $f(x)$  where  $X=x \Rightarrow Y=f(x)$

**best prediction property** find  $f$  and try to minimize  $\mathbb{E}[(Y-f(X))^2]$

$$\mathbb{E}[(Y-f(X))^2] = \mathbb{E}[\mathbb{E}[(Y-f(X))^2 | X]]$$

↪  $\min_{f(x)} \mathbb{E}[(Y-f(X))^2 | X=x]$  for each  $x$

$f(x) = a$  where  $h(a) = \mathbb{E}[(Y-a)^2 | X=x]$  is minimized,  
 error

$$\text{make } h'(a) = 0 \rightarrow \mathbb{E}[-2Y + 2a | X=x] \Rightarrow 2\mathbb{E}[Y-a | X=x] = 0$$

$$a = f(x) = \mathbb{E}[Y | X=x]$$

works for  
multivariate  $Y$

regression function of  
 $Y$  onto  $X$ .

→ it's hard to perfectly compute  $f(x)$  for all  $x$  since we are given only a few  $(X_i, Y_i)$  to observe.

→  $\mathbb{E}[Y | X=x]$  is only the mean, which doesn't fully capture  $Y|X=x$ .  
Logistic regression uses confidence bands around  $\mathbb{E}[Y | X=x]$  instead

### Linear Regression

$$f(x) = \mathbb{E}[Y | X=x] = x^\top \beta$$

for some  $\beta = \beta^* \in \mathbb{R}^k$  fixed.

use MLE to find  $\beta$ , but assume the following:

$$\rightarrow Y|X=x \sim N(f(x), \sigma^2(x)) \quad [\text{gaussian}]$$

$$\rightarrow f(x) \text{ linear} \rightarrow f(x) = x^\top \beta^*$$

$$\rightarrow \sigma^2(x) = \sigma^2 \quad \forall x \quad (\text{constant})$$

then

$$\ell_n(\beta) = \sum_{i=1}^n \log \left( \frac{1}{\sigma \sqrt{2\pi}} \exp \left( -\frac{(Y_i - x_i^\top \beta)^2}{2\sigma^2} \right) \right) = -\sum_{i=1}^n \frac{(Y_i - x_i^\top \beta)^2}{2\sigma^2} + \text{const.}$$

minimize:  $\hat{\beta}^{\text{MLE}} = \underset{\beta \in \mathbb{R}^k}{\operatorname{argmin}} \sum_{i=1}^n (\gamma_i - x_i^\top \beta)^2$

for least squares

$$Y = \begin{pmatrix} \gamma_1 \\ \vdots \\ \gamma_n \end{pmatrix} \in \mathbb{R}^n \quad X = \begin{pmatrix} -x_1^\top - \\ -x_2^\top - \\ \vdots \\ -x_n^\top - \end{pmatrix} \in \mathbb{R}^{n \times k}$$

want  $\underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (\gamma_i - x_i^\top \beta)^2 = \underset{\beta}{\operatorname{argmin}} \|Y - X\beta\|^2$

gradient:  $\nabla_{\beta} \|Y - X\beta\|^2 = 2X(Y - X\beta) = 0$

$$\hookrightarrow X^\top Y = X^\top X\beta \rightarrow \overset{\text{also } \hat{\beta}^{\text{LS}}}{\hat{\beta}^{\text{LS}}} = (X^\top X)^{-1} X^\top Y$$

$$X\hat{\beta} = \underbrace{X(X^\top X)^{-1} X^\top Y}_P$$

"projection onto"  $P = \text{span}(X)$

LECTURE 24 4/17/24 1PM

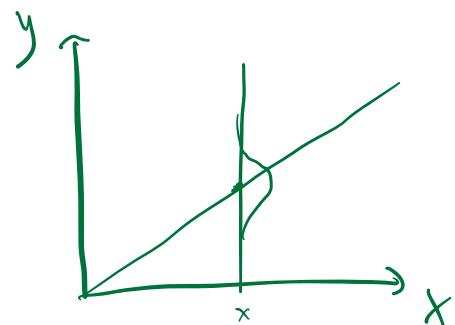
lin-reg.

$X \in \mathbb{R}^{n \times k}$ ,  $\gamma \in \mathbb{R}$  predict  $\gamma$  from  $X$

given  $(x_i, \gamma_i) \quad i \in [n]$

Assume

$E[\gamma | X=x] = x^\top \beta^*$  with Gaussian errors



$$\gamma_i = x_i^\top \beta^* + \varepsilon_i \quad \varepsilon_i \sim N(0, \sigma^2)$$

$$Y = \begin{pmatrix} \gamma_1 \\ \vdots \\ \gamma_n \end{pmatrix} \in \mathbb{R}^n \quad X = \begin{pmatrix} -x_1^\top - \\ -x_2^\top - \\ \vdots \\ -x_n^\top - \end{pmatrix} \in \mathbb{R}^{n \times k} \quad \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$Y = X\beta + \overset{\hat{\varepsilon}}{\varepsilon} \quad \hat{\varepsilon} \sim N(0, \sigma^2 I_n)$$

note  $\hat{\beta}^{\text{MLE}} = \hat{\beta}^{\text{LS}} = (X^\top X)^{-1} X^\top (X^\top \beta^* + \varepsilon) = (X^\top X)^{-1} (X^\top \beta^*) \beta^* + \underbrace{(X^\top X)^{-1} X^\top \varepsilon}_{\hat{A}\varepsilon}$

$$\varepsilon \sim N(0, \sigma^2 I_n) \Rightarrow A\varepsilon \sim N(0, A\sigma^2 I_n A^\top) \rightarrow N(0, \sigma^2 (X^\top X)^{-1} X^\top X (X^\top X)^{-1}) \rightarrow N(0, \sigma^2 (X^\top X)^{-1})$$

ex) blood pressure = Y

$$X = \begin{pmatrix} \text{age} \\ \text{weight} \\ \text{shoe size} \end{pmatrix} \quad b_{10} = \beta_1^* \text{age} + \beta_2^* \text{weight} + \beta_3^* \text{shoe size} = X^T \beta^*$$

if  $\beta_1^* > 0$  then  $b_{10} \uparrow$  as age  $\uparrow$  if  $\beta_1^* = 0$ , no.

$$\hat{\beta}_j^{LS} \sim N(\beta_j^*, \sigma^2(X^T X)^{-1}_{jj}) \quad \text{so CI: } \beta_j^* \in [\hat{\beta}_j^{LS} \pm \hat{\sigma} \sqrt{(X^T X)^{-1}_{jj}} z_{\alpha/2}]$$

$$X^T X = X_1 X_1^T + X_2 X_2^T + \dots + X_n X_n^T \quad \text{so it plays role of } \frac{1}{\sqrt{n}}$$

### Hypothesis Testing

$$H_0: \beta_j^* = 0 \quad \text{vs} \quad \beta_j^* \neq 0 \quad \text{under null, } \hat{\beta}_j^{LS} \sim N(0, \sigma^2(X^T X)^{-1}_{jj})$$

$$\text{Reject (Wald)} \quad \text{if} \quad \left| \frac{\hat{\beta}_j^{LS}}{\hat{\sigma} \sqrt{(X^T X)^{-1}_{jj}}} \right| > z_{\alpha/2} \quad p\text{-value} = P(|z| > \frac{\hat{\beta}_j^{LS}}{\hat{\sigma} \sqrt{(X^T X)^{-1}_{jj}}})$$

how to gen  $\hat{\sigma}$ ? use residuals  $\hat{\epsilon}_i = Y_i - X_i^T \hat{\beta}^{LS}$  column space of  $X$  has dimension  $k$

$$\text{so } \hat{\epsilon} = Y - X^T \hat{\beta} \text{ dimension } n-k$$

$$\text{we can use } \hat{\sigma}^2 = \frac{\|\hat{\epsilon}\|^2}{n-k} = \frac{1}{n-k} \sum_{i=1}^n \hat{\epsilon}_i^2.$$

LECTURE 25 4/19/24 (Did not attend)

What if the  $Y_i$ 's are binary?

$$\hookrightarrow Y | X=x \sim \text{Ber}(f(x)) \rightarrow \mathbb{E}[Y | X=x] = f(x)$$

Now  $X^T \beta^*$  is bad since this is unbounded.

need  $f(x) = \sigma(X^T \beta^*)$  where

- $\sigma$  increasing
- $\sigma(t) \rightarrow 0$  as  $t \rightarrow -\infty$ ,  $\sigma(t) \rightarrow 1$  as  $t \rightarrow \infty$
- $\sigma(0) = 1/2$

## Alternative 1 [Logistic Regression Model]

$Y|X=x \sim \text{Ber}(\sigma(x^T \beta^*))$  where  $\sigma(t) = \frac{e^t}{1+e^t}$  (sigmoid),  
(logistic).

## Alternative 2 [Probit Regression Model]

$Y|X=x \sim \text{Ber}(\phi(x^T \beta^*))$  where  $\phi(t) = \Phi(t)$

can also have  $Y = \mathbb{1}(x^T \beta^* + Z > 0)$  where  $Z \sim N(0, 1)$

↑ any other  $\tilde{\sigma}$  w/symmetric distribution  
and cdf  $F = \sigma$  can also  
be an alternative;

$Y|X=x \sim \text{Ber}(F(x^T \beta^*))$

$Y = \mathbb{1}(x^T \beta^* + \tilde{Z} > 0)$

for Sigmoid

$$\frac{\sigma}{1-\sigma} = \frac{\frac{e^t}{1+e^t}}{\frac{1}{1+e^t}} = e^t$$

$$P(Y_i|X_i) = \sigma(x_i^T \beta)^{Y_i} (1 - \sigma(x_i^T \beta))^{1-Y_i}$$

$$\ell_n(\beta) = \sum_{i=1}^n Y_i \log \sigma(x_i^T \beta) + (1-Y_i) \log (1 - \sigma(x_i^T \beta))$$

$$= \sum_{i=1}^n \left[ Y_i \log \frac{\sigma(x_i^T \beta)}{1 - \sigma(x_i^T \beta)} + \log (1 - \sigma(x_i^T \beta)) \right] \quad \curvearrowleft \frac{1}{1+e^{x_i^T \beta}}$$

$$\text{concave} \rightarrow = \sum_{i=1}^n [Y_i x_i^T \beta - \log (1 + e^{x_i^T \beta})]$$

gradient ascent

$$\beta^{(k)} \in \mathbb{R}^K \quad \beta^{(k+1)} = \beta^{(k)} + \gamma \nabla \ell_n(\beta)$$

## Multiclass Classification $\rightarrow Y \in \{0, 1, \dots, M\}$ instead of $\{0, 1\}$

want  $P(Y=l|X=x) = p_l(x) \quad \sum_{l=0}^M p_l(x) = 1 \quad \forall x,$

$M$  degrees of freedom so we shouldn't use

$$p_j(x) = \frac{e^{x^T \beta_j^*}}{\sum_{l=0}^M e^{x^T \beta_l^*}} \quad (\text{note } M+1 \neq M \text{ different } \beta_j^* \text{'s})$$

for  $M=1$  recall

"logit" function =  $e^{-x}$

$$f(x) = \frac{e^{x^T \beta}}{1+e^{x^T \beta}} \Rightarrow x^T \beta = \log \left( \frac{f(x)}{1-f(x)} \right) = \log \left( \frac{p_1(x)}{p_0(x)} \right)$$

analogously let  $\log \left( \frac{p_j(x)}{p_0(x)} \right) = x^T \beta_j$  for  $j=1, 2, \dots, M$

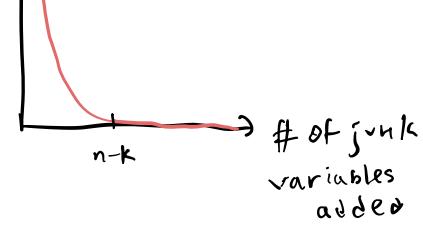
$$\hookrightarrow p_j(x) = \frac{e^{x^T \beta_j}}{1 + \sum_{\ell=1}^M e^{x^T \beta_\ell}}, \quad p_0(x) = \frac{1}{1 + \sum_{\ell=1}^M e^{x^T \beta_\ell}}$$

LECTURE 26 4/22/24 1PM

Avoiding Overfitting  $\rightarrow$  but how do we pick which features to include in  $\beta$ ? (variables)

$\|Y - X\beta\|^2$  gets better the more features we use

$\|Y - X\beta\|^2$  "more junk  $\rightarrow$  better predictor" BAD



Goal Given  $k$  features, pick subset  $S$

and get model

$$Y = \sum_{i \in S} \beta^{(i)} X^{(i)} + \epsilon$$

Hypothesis Testing Fails: e.g. Wald's for each  $\beta^{(j)}$

$$H_{0,j}: \beta^{(j)} = 0 \quad \text{vs.} \quad H_{1,j}: \beta^{(j)} \neq 0$$

$\hookrightarrow$  Bonferroni is too conservative.

$\hookrightarrow$  BH requires independent test statistics

Moreover,

$\hookrightarrow$  we want to test  $H_{0,j}: \beta^{(j)} \neq 0$  vs.  $H_{1,j}: \beta^{(j)} = 0$

$\hookrightarrow$  need  $n-k$ , not  $n$  to be large for asymptotics

## New Framework, $R^2$

given  $S \subset \{1, 2, \dots, k\}$

$$R^2(S) := 1 - \frac{\|\gamma - X\hat{\beta}(S)\|^2}{\|\gamma - \bar{Y}_n \mathbf{1}\|^2}$$

$\hat{\beta}$  using only  
subset  $S$   
of variables  
 ↑  
all 1's vector (constant fit)

(coefficient of determination)

- better linear fit

→  $R^2$  closer to 1

- if one feature is always 1, then  $0 \leq R^2(S)$

- larger  $|S|$

↳ larger  $R^2(S)$

## forward/greedy model selection

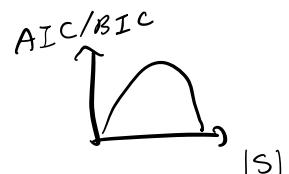
→ start with  $S = \emptyset$

→ add  $i \in \{1, 2, \dots, k\} \setminus S$  to  $S$  to maximize new  $R^2(S)$ ,

→ stop when  $R^2(S)$  plateaus (or exceeds desired  $R^2$  threshold)

## other metrics

$$\text{AIC (Akaike info criterion)} = \ell_n(\hat{\beta}(S)) - |S|$$



$$\text{BIC (Bayesian info criterion)} = \ell_n(\hat{\beta}(S)) - \frac{\log n}{2} |S|$$

## Issues / solutions

→  $\binom{k}{|S|}$  models for given  $|S|$ .

↳ use greedy to avoid cost

→ backward model selection (start with  $S$  and remove elements decreasing  $R^2$  the least)

→ stepwise model selection (pick multiple features at a time)

→ AIC/BIC can be maximized without worrying about threshold due to self-penalization

Survival Analysis - estimate distribution of time until event of interest.

ex) # months since subs $T_i$	cancelled $S_i$	censored $C_i$
5	1	$\geq 5$
12	1	$\geq 12$
7	1	$\geq 7$
35	1	$\geq 35$
48	censored	48
6	1	26
48	censored	48
48	censored	48

random variable  $T$ , observe  $T_1, \dots, T_n \stackrel{\text{iid}}{\sim} T$

Def Let  $F$  be the cdf of  $T$ . Survival function is  $S(t) = 1 - F(t)$ .

we want to estimate  $S(t)$

If we know all the  $T_i$ , we can just say  $\hat{S}(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(T_i > t)$

$\uparrow$  unbiased       $\uparrow$  consistent       $\uparrow$  asympt. normal

$$\sqrt{n}(\hat{S}(t) - S(t)) \rightsquigarrow N(0, S(t)(1-S(t)))$$

but some  $T_i$  may be censored due to dropout/ongoing.

so each  $T_i$  has a censoring time  $C_i \stackrel{\text{iid}}{\sim} C$  and we instead see

$$\tilde{T}_i = \min(T_i, C_i), \quad S_i = \mathbb{1}(C_i \geq T_i) = \mathbb{1}(\tilde{T}_i = T_i)$$

throw out censored data?

include it?

$$\hat{S}_{\text{naive}}(t) = \frac{\sum_{i=1}^n \mathbb{1}(T_i > t, C_i = \infty)}{\sum_{i=1}^n \mathbb{1}(C_i = \infty)}$$

$$\hat{S}_{\text{better}}(t) = \frac{\sum_{i=1}^n \mathbb{1}(\tilde{T}_i > t)}{\sum_{i=1}^n \mathbb{1}(C_i \geq t)} \approx P(T > t | C \geq t)$$

Kaplan-Meier Estimator

$$S(t) = P(T > t) = \frac{P(T > t)}{P(T > t-1)} \cdot \frac{P(T > t-1)}{P(T > t-2)} \cdot \dots \cdot \frac{P(T > 1)}{P(T > 0)} \leftarrow 1$$

Hazard Rate

$h(t) = 1 - q(t) = \text{probability that you die by } t \text{ given you survived to } t-1$ .

$$\hat{h}(s) = \frac{\sum_{i=1}^n \mathbb{1}(\tilde{T}_i \geq s, C_i \geq s)}{\sum_{i=1}^n \mathbb{1}(\tilde{T}_i \geq s)} \xrightarrow{\text{indep.}} \frac{P(\tilde{T} \geq s, C \geq s)}{P(\tilde{T} \geq s)} = \frac{P(\tilde{T} = s, C \geq s)}{P(\tilde{T} \geq s, C \geq s)} = \frac{P(\tilde{T} = s)}{P(\tilde{T} \geq s)} = h(s)$$

$$S(t) = \prod_{s=1}^t (1 - \hat{h}(s)) = \prod_{s=1}^t \left(1 - \frac{\#\{\tilde{T}_i = s, C_i \geq s\}}{\#\{\tilde{T}_i \geq s\}}\right)$$

variations - e.g. if it depends on features  $x \in \mathbb{R}^k$

So survival function can become  $S(t, x)$  depending on selected features  $x$ .

### Cox proportional hazard regression model

$$h(t, x) = h_0(t) \exp(\beta^T x) \rightarrow \text{decouples effect of } t \text{ vs. } x$$

also, continuous version of Kaplan-Meier.

$$h(t) = \frac{P(T=t)}{P(T > t-1)} \quad \text{vs.} \quad h(t) = \frac{P(t \leq T \leq t+dt)}{P(T \geq t)} = -\frac{S'(t)}{S(t)} = -\frac{d}{dt} \log(S(t))$$

LECTURE 28 4/26/24 (Did not attend)

cavation & correlation

Is (counterfactual) Model  $X \in \{0, 1\}$  no vs. yes treatment  
 $Y \in \mathbb{R}$  response

potential outcomes,  $Y = \begin{cases} C_0 & \text{if } X=0 \\ C_1 & \text{if } X=1 \end{cases}$

$X$	$C_0$	$C_1$
0	obs.	counterfactual
1	counterfactual	obs.

average treatment effect (ATE)  $\theta := \mathbb{E}[C_1] - \mathbb{E}[C_0]$

$$\alpha := \mathbb{E}[Y|X=1] - \mathbb{E}[Y|X=0]$$

association

not necessarily the same! e.g. if  $Z \sim \text{unif}([-1, 1])$ ,  $X = \mathbf{1}(Z > 0)$   
outcomes  $C_0 = C_1 = Z$ .

Then  $\theta = 0$ ,  $\alpha = 1$ .  $\rightarrow Z$  is confounding variable.

We want to find  $\theta$ , so we try to estimate it using  $\alpha$ .

If  $\alpha = \theta$ , just use plug-in estimator

$$\hat{\alpha} = \frac{\sum_{i=1}^n Y_i \mathbf{1}(X_i=1)}{\sum_{i=1}^n \mathbf{1}(X_i=1)} - \frac{\sum_{i=1}^n Y_i \mathbf{1}(X_i=0)}{\sum_{i=1}^n \mathbf{1}(X_i=0)}.$$

- To guarantee, we can use randomized control trials - flip coin with prob  $p$  to assign  $X_i=1$  and  $1-p$  to assign  $X_i=0$ .
- ↳  $\hat{\alpha} \approx \alpha = \mathbb{E}[Y|X=1] - \mathbb{E}[Y|X=0] = \mathbb{E}[c_1|X=1] - \mathbb{E}[c_0|X=0] = \mathbb{E}[c_1] - \mathbb{E}[c_0] = 0$ .
- ↳ we have  $(c_0, c_1) \perp\!\!\!\perp X$  → "strong unconfoundedness"
- ↑  
independent
- ↳ can be unethical to deny treatment by coin flip.

Sometimes, we can instead find a feature  $Z$  to condition on s.t.

$c_i \perp\!\!\!\perp X|Z$  for  $i=0, 1$ . e.g. age, salary, etc.

but if  $Z$  is high dimensional, we can't use  $\hat{\theta}_1 = \hat{\alpha}(z_1)$ ,  $\hat{\theta}_2 = \hat{\alpha}(z_2)$ , instead we use propensity score  $p(z) = P(X=1 | Z=z)$

↳ likelihood of assigning treatment based on  $Z=z$ .

Thm If  $c_i \perp\!\!\!\perp X|Z$ , then  $c_i \perp\!\!\!\perp X|p(z)$ .

so we can just divide  $X$  into bins based on  $p(z)$  and compute  $\hat{\theta}_j$  for each.

↳ calculate  $p(z)$  using logistic regression:  $\text{logit}(p(z)) = z^T \beta$ .

LECTURE 29 4/29/24 1PM

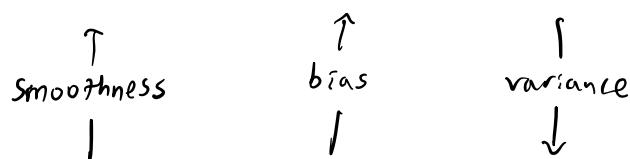
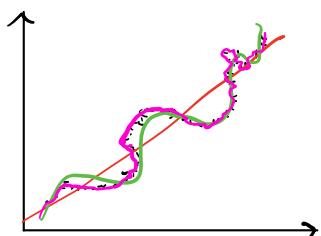
### Nonparametric Curve Estimation

1: density estimation — estimate density  $p(x)$  of iid  $X_1, \dots, X_n$

2: regression — estimate  $f(x) = \mathbb{E}[Y|X=x]$  given iid  $(X_1, Y_1), \dots, (X_n, Y_n)$

we don't want too smooth nor too choppy.

want a tradeoff after a smoothness assumption.



smoothness means close to Taylor approx of some order

Suppose we estimate  $g$  by  $\hat{g} \Rightarrow \hat{g}(x)$  for  $g(x)$  for each  $x$ .

$$\text{Bias: } b(x) = \mathbb{E}[\hat{g}(x)] - g(x)$$

$$\text{Variance: } v(x) = \mathbb{V}[\hat{g}(x)]$$

$$\text{MSE: } b(x)^2 + v(x)$$



MSE is also called risk

$$\int b(x)^2 dx + \int v(x) dx.$$

$$\text{MISE: } \int \text{MSE}(\hat{g}(x)) dx = \mathbb{E} \left[ \int (\hat{g}(x) - g(x))^2 dx \right]$$

Density Estimation      probability density function

$$\text{Suppose } X_1, \dots, X_n \stackrel{iid}{\sim} p(x)$$

① histogram estimator  $\hat{p}(x)$

$m$  bins  $B_1, \dots, B_m$  each width  $|B_j| = \frac{1}{m} = h$

$$\text{Let } n_j = \#\{i : X_i \in B_j\} = \sum_{i=1}^n \mathbb{1}(X_i \in B_j)$$

$$\text{Then let } \hat{p}_j = \frac{n_j}{n}. \text{ Define } \hat{p}(x) := \frac{\hat{p}_j}{h} = \sum_{j=1}^m \frac{\hat{p}_j}{h} \mathbb{1}(x \in B_j)$$

↑ normalization

note  $\hat{p}(x) \geq 0 \quad \forall x \geq 0$  and

$$\int \hat{p}(x) dx = \sum_{j=1}^m \int \frac{\hat{p}_j}{h} \mathbb{1}(x \in B_j) dx = \sum_{j=1}^m \int_{B_j} \frac{\hat{p}_j}{h} dx = \sum_{j=1}^m h \cdot \frac{\hat{p}_j}{h} = \sum_{j=1}^m \hat{p}_j = 1.$$

bias for given  $x \in B_j$

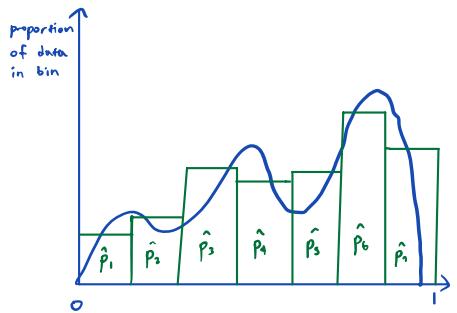
$$b(x) = \mathbb{E}[\hat{p}(x)] - p(x) = \mathbb{E}\left[\frac{\hat{p}_j}{h}\right] - p(x)$$

$$\text{but } \hat{p}_j = \frac{n_j}{n}, \quad n_j \sim \text{Bin}(n, \int_{B_j} p(y) dy)$$

$$\text{so } \mathbb{E}\left[\frac{\hat{p}_j}{h}\right] - p(x) = \frac{1}{h} \int_{B_j} p(y) dy - p(x) \rightarrow 0 \text{ as } h \rightarrow 0.$$

Also true that  $\int b(x)^2 dx \rightarrow 0$  as  $h \rightarrow 0$ .

$\Rightarrow$  smaller  $h$  means smaller bias



## Variance

$$v(x) = \mathbb{V} \left[ \frac{\hat{p}_j}{h} \right] = \frac{1}{h^2} \mathbb{V}[\hat{p}_j] = \frac{1}{h_n^2} \left[ \int_{B_j} p(y) dy \right] \left[ 1 - \int_{B_j} p(y) dy \right]$$

$$= \frac{1}{h_n} \cdot \frac{\int_{B_j} p(y) dy}{h} \cdot 1 - \int_{B_j} p(y) dy$$

⇒ smaller  $h$  means

bigger variance.

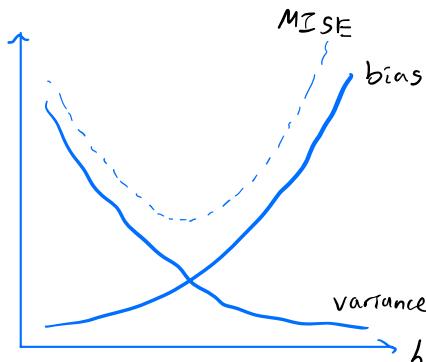
$$\rightarrow \frac{1}{h_n} \cdot p(x) \cdot 1 = \frac{p(x)}{h_n}$$

as  $h \rightarrow 0$ .

so we have a tradeoff!

→ to find optimal  $h$  without actually knowing  $p$ ,

we use Cross-validation



LECTURE 30 5/6/24 1PM

Cross-validation: note  $MSE \leftarrow \int (\hat{p}(x) - p(x))^2 dx$  ← unbiased estimator

$$= \underbrace{\int \hat{p}(x)^2 dx}_{\text{we know}} - 2 \int \hat{p}(x)p(x) dx + \underbrace{\int p(x)^2 dx}_{\text{constant, no need to optimize}}$$

note  $\int \hat{p}(x)p(x) dx = \mathbb{E}_{x \sim p} [\hat{p}(x)]$

$\hat{p}(x)$  depends on  $X_i$  so write as  $\hat{p}(x; X_1, X_2, \dots, X_n)$

$$\int \hat{p}(x)p(x) dx = \mathbb{E}_{X \sim p} [\hat{p}(X; X_1, \dots, X_n) | X_1, \dots, X_n]$$

Suppose we have some  $X'_1, \dots, X'_n$  → independent of  $X_1, \dots, X_n$

$$\hookrightarrow \int \hat{p}(x)p(x) dx = \frac{1}{n} \sum_{i=1}^n \hat{p}(X'_i; X_1, \dots, X_n) \quad \{X_1, \dots, X_n\} \setminus \{X'_i\}$$

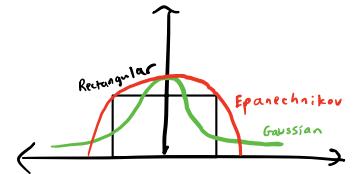
so minimize  $\widehat{MSE} = \int_0^1 \hat{p}(x)^2 dx - \frac{2}{n} \sum_{i=1}^n \hat{p}(X'_i; X_1, \dots, X_n) + \text{const.}$   
over bin width  $h$ .

(2) Kernel density estimator - smooth estimator of density.

↪ Kernel is any function  $K$  s.t.  $K(x) \geq 0$ ,  $\int K(x) dx = 1$ ,  $\int x K(x) dx = 0$ .

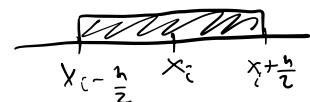
↑  
any even pdf works.

KDE is  $\hat{p}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)$   $h$  is bandwidth here.  
(how wide each kernel stretches)



ex) rectangular kernel  $K(x) = \mathbb{1}(|x| \leq \frac{1}{2})$  "sliding bin histogram"

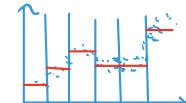
$$K\left(\frac{x-x_i}{h}\right) = \mathbb{1}\left(\left|\frac{x-x_i}{h}\right| \leq \frac{1}{2}\right) = \mathbb{1}\left(x \in [x_i - \frac{h}{2}, x_i + \frac{h}{2}]\right)$$



Non parametric Regression Given  $(X_i, Y_i)$ , predict  $E[Y | X=x]$  if nonlinear.

(1) Regressogram - break into various bins  $B_j$  along x-axis

$$E[Y | X=x] \approx \hat{f}(x) = \frac{\sum_{i=1}^n Y_i \mathbb{1}(X_i \in B_j)}{\sum_{i=1}^n \mathbb{1}(X_i \in B_j)} \rightarrow \text{estimate as average of each bin.}$$



(2) Nadaraya-Watson Estimator - sort of like KDEs

↪ estimate both  $p(x,y)$  and  $p(y)$  using KDEs

$$\text{then } E[Y | X=x] = \int y p(y|x) dx = \frac{\int y p(x,y) dx}{p(x)} \approx \frac{\int y \hat{p}(x,y) dx}{\hat{p}(x)}$$

$$\Rightarrow \hat{f}(x) = \frac{\frac{1}{nh^2} \int y \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right) K\left(\frac{y-y_i}{h}\right) dy}{\frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)} = \frac{\sum_{i=1}^n Y_i K\left(\frac{x_i-x}{h}\right)}{\sum_{i=1}^n K\left(\frac{x_i-x}{h}\right)}.$$

Curse of dimensionality:  $MISE \sim \left(\frac{1}{n}\right)^k / n$  with  $k$ -dimensional  $x_i$

↪ need LARGE  $n$  to make MISE small

Survey Sampling - pick a representative set of people for survey.

ex) MIT undergrads who prefer online > in-person.  $N = 4657$ .

$$\text{Let } y_i = \begin{cases} 1, & \text{if } i \text{ prefers online} \\ 0, & \text{if } i \text{ prefers in-person} \end{cases} \quad \text{We want } T = \sum_{i=1}^N y_i$$

we instead find a representative sample  $S \subset \{1, 2, \dots, N\}$ , use

$$\hat{T} = \sum_{i \in S} \lambda_i y_i$$

$\lambda_i \geq 1$  compensation weights

Now, how do we choose  $S$ ? The probability distribution of the random subset  $S$  is called the design.

common design choices

$\hookrightarrow$  without replacement  
one  $s$  from all  $S \subset \{1, \dots, N\}$  with  $|S|=n \rightarrow P(S=s) = \frac{1}{\binom{N}{n}}$

$\hookrightarrow$  random

each  $i \in \{1, \dots, N\}$  is included in  $S$  with  $p = \frac{n}{N}$ , size of  $S$  is  $|S| \sim \text{Bin}(N, \frac{n}{N})$ , expected size  $n$ .

$\hookrightarrow$  stratified

break up into groups (e.g. Major) then sample from each group:  
 $S$  union of samples

constructing an estimator

let  $p$  = pmf of random subset  $S$ .  $p(s) = P(S=s)$  where  $s \in \{1, \dots, N\}^n$

let  $\pi_i = P(i \in S) = \sum_{s \in \{1, \dots, N\}^n} p(s)$  be the inclusion probability.

let  $\pi_{i,j} = P(i, j \in S)$  be similarly defined.

## Moritz-Thompson Estimator

$$\hat{T}^{HT} = \sum_{i \in S} \frac{y_i}{\pi_i} = \sum_{i=1}^N \frac{y_i}{\pi_i} \mathbb{1}(i \in S)$$

$\lambda_i = \frac{1}{\pi_i}$  here.

for sampling w.o. replacement

$$\hat{T}^{HT} = \frac{N}{n} \sum_{i \in S} y_i$$

$$E[\hat{T}^{HT}] = \sum_{i=1}^N E\left[\frac{y_i}{\pi_i} \mathbb{1}(i \in S)\right] = \sum_{i=1}^N \frac{y_i}{\pi_i} P(i \in S) = \sum_{i=1}^N \frac{y_i}{\pi_i} \cdot \pi_i = T \rightarrow \text{unbiased.}$$

$$\begin{aligned} V[\hat{T}^{HT}] &= V\left[\sum_{i=1}^N \frac{y_i}{\pi_i} \mathbb{1}(i \in S)\right] = \sum_{i=1}^N \sum_{j=1}^N \frac{y_i}{\pi_i} \cdot \frac{y_j}{\pi_j} \text{Cov}(\mathbb{1}(i \in S), \mathbb{1}(j \in S)) \\ &= \sum_{i=1}^N \sum_{j=1}^N \frac{y_i}{\pi_i} \cdot \frac{y_j}{\pi_j} (\pi_{ij} - \pi_i \pi_j) \quad \leftarrow \text{in random sampling } \mathbb{1}(i \in S) \perp \mathbb{1}(j \in S) \\ &= \sum_{i=1}^N \frac{y_i^2}{\pi_i^2} (\pi_i - \pi_i^2) = \sum_{i=1}^N y_i^2 \left(\frac{1-\pi_i}{\pi_i}\right) \quad \leftarrow \text{but then this can be large if } \pi_i \text{ small} \end{aligned}$$

estimate variance using

$$\hat{V}(\hat{T}^{HT}) = \sum_{i \in S} \sum_{j \in S} \frac{y_i}{\pi_i} \cdot \frac{y_j}{\pi_j} \cdot \frac{\pi_{ij} - \pi_i \pi_j}{\pi_{ij}} \quad \begin{matrix} \text{then} \\ \text{(in some design choices)} \end{matrix} \quad \frac{\hat{T}^{HT} - T}{\hat{V}(\hat{T}^{HT})} \rightsquigarrow N(0, 1)$$

Methods for Sampling w.o. replacement

(1) one at a time

$N$  choices for 1<sup>st</sup>,  $N-1$  for 2<sup>nd</sup>, ...,  $N-n+1$  for  $n^{\text{th}}$   $\Rightarrow O(Nn)$ .

(2) give each person a uniform random variable  $\in [0, 1]$   
sort random variables and take  $n$  people with smallest variables.

$O(N \log N)$

(3) reservoir

$\rightarrow$  set  $S = \{1, 2, \dots, n\}$

$\rightarrow$  for  $k = n+1, \dots, N$ , replace a random element of  $S$  with  $k$   
with  $p = \frac{n}{k}$ . Else, do nothing.

$O(n^2 \log N)$  expected

classifier  $h: X \rightarrow Y \in \{0, 1\}$  (binary)  $\rightarrow$  defines a classification region

ex) linear classifier  $h(x) = \begin{cases} 1 & (w^T x + b \geq 0) \\ -1 & (w^T x + b < 0) \end{cases}$

$k+1$  parameters  $(w_1, \dots, w_k, b)$

## Bayes Classifier

minimize probability of making a mistake.

true error  $P(Y \neq h(X))$  vs. empirical error  $\frac{1}{n} \sum_{i=1}^n \mathbb{1}(Y_i \neq h(X_i))$

$$h^* := \underset{h}{\operatorname{argmin}} \mathbb{P}(Y \neq h(x)) = \mathbb{1}_{\{r(x) \geq 1/2\}}$$

↑    ↑  
 same since  
 both minimize →  $\mathbb{P}(Y=1 | X=x)$   
 error

↳ alternative way to arrive at this.

Let proportion of 0's, 1's in  $\gamma$  be  $\frac{1}{2}$  each.

then

$$r(x) = P(Y=1 \mid X=x) = \frac{f_1(x) \cdot \frac{1}{2}}{f_1(x) \cdot \frac{1}{2} + f_0(x) \cdot \frac{1}{2}} = \frac{f_1(x)}{f_1(x) + f_0(x)} \geq \frac{1}{2} \text{ iff } f_1(x) \geq f_0(x)$$

## 3 estimators of Bayes Classifier $h^*$

$$\textcircled{1} \quad \hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}(Y_i \neq h(X_i)) \quad \text{"empirical risk minimizer" (ERM)}$$

$\hookrightarrow$  brute force, perceptron

②  $\hat{h} = \mathbb{1}(r(x) \geq \frac{1}{2})$ , using  $\hat{r}(x) = \sigma(\hat{\beta}^T x)$  for example

③  $\hat{h} = \mathbb{1}(\hat{f}_i(x) \geq \hat{f}_o(x))$ , using density estimators (e.g. KDEs) of choice  
 → ex) Fisher's discriminant analysis

More on ③

can use

$$\hat{f}_o(x) = \frac{1}{n_0 h^k} \sum_{i: y_i=0} \text{dimension of } x_i k\left(\frac{\|x_i - x\|}{h}\right)$$

$$\hat{f}_1(x) = \frac{1}{n_1 h^k} \sum_{i: y_i=1} k\left(\frac{\|x_i - x\|}{h}\right)$$

but error  $|F_1(x) - \hat{f}_1(x)| \propto n^{-1/(2+k)}$  so big  $k$  is problematic  
without a lot of data

Solutions① Naive Bayes.  $\rightarrow$  assume all features independent.

$$\hat{f}_0(x) = \prod_{j=1}^k \hat{f}_0^{(j)}(x^{(j)}) \quad \hat{f}_1(x) = \prod_{j=1}^k \hat{f}_1^{(j)}(x^{(j)})$$

② Discriminant Analysis  $\rightarrow$  assume  $f_0, f_1$  normal.

then

$$\hat{f}_0(x) = \frac{1}{\sqrt{(2\pi)^k \det(\Sigma_0)}} \exp\left(-\frac{1}{2}(x - \mu_0)^\top \Sigma_0^{-1}(x - \mu_0)\right)$$

$$\hat{f}_1(x) = \frac{1}{\sqrt{(2\pi)^k \det(\Sigma_1)}} \exp\left(-\frac{1}{2}(x - \mu_1)^\top \Sigma_1^{-1}(x - \mu_1)\right)$$

if  $\Sigma_0 = \Sigma_1 = \Sigma$ ,

$$\begin{aligned} \hat{f}_1(x) \geq \hat{f}_0(x) &\Leftrightarrow (x - \mu_0)^\top \Sigma^{-1}(x - \mu_0) \geq (x - \mu_1)^\top \Sigma^{-1}(x - \mu_1) \\ &\Leftrightarrow x^\top (\Sigma^{-1}(\mu_1 - \mu_0)) \leq \frac{1}{2} (\mu_1^\top \Sigma^{-1} \mu_1 - \mu_0^\top \Sigma^{-1} \mu_0) \end{aligned}$$

we can estimate  $\mu_1, \mu_0, \Sigma$  to get a classifier  
using sample mean, covariance.

if  $\Sigma_0 \neq \Sigma_1$ , we get a quadratic boundary,  
 $\rightarrow$  "quadratic discriminant analysis"

### ③ k - nearest neighbor classifier

$\hat{h}(x) = 1 \text{ if } \geq \frac{k}{2} \text{ nearest neighbors are } Y_i=1, \text{ else } 0.$

### ④ logistic regression

$$\boxed{0,1} \quad \hat{p}(x) = \frac{e^{\hat{\beta}^{MLE T} x}}{1 + e^{\hat{\beta}^{MLE T} x}} \quad \hat{h}(x) = \mathbb{1}(\hat{p}(x) \geq \frac{1}{2})$$

**multiclass**  $\hat{p}_j(x) = \frac{e^{\hat{\beta}_j^{MLE T} x}}{1 + \sum_{k=1}^M e^{\hat{\beta}_k^{MLE T} x}} = P(Y=j | X=x) \quad \text{for } j=0, 1, \dots, M$

$$\hat{h}(x) = \operatorname{argmax}_{j=0, \dots, M} \hat{r}_j(x) \quad \log \left( \frac{p_j(x)}{p_0(x)} \right) = x^T \beta_j$$

### ⑤ neural networks

$$\log \left( \frac{p_j(x)}{p_0(x)} \right) = f_{\theta_j}(x) \quad \begin{matrix} \text{more} \\ \text{complex function} \\ \text{parametrized by } \theta \end{matrix}$$

Last layer does multiclass logistic regression