

Review of Probability

Probability is common sense of uncertainty

Sample Space: the set of all possible outcomes (Ω, S)

Event: statement about the outcome. Subset of the Sample Space (A, B, C, \dots)

\emptyset : Null/Empty Event. Probability is defined on the event

Relations:

AND: $A \cap B$



OR: $A \cup B$



NOT: \bar{A}



Equally likely setting

- Randomly sample a single person in a population. All people are equally likely to be sampled.

$$P(A) = \frac{|A|}{|\Omega|} ; A = \text{sub-population} \quad |A|: \# \text{ of elements in } A$$

$\Omega = \text{entire population}$

Axioms:

$$(1) P(\Omega) = 1$$

$$(2) P(A) \geq 0$$

$$(3) \text{Additivity. If } A \cap B = \emptyset \text{ (disjoint) then } P(A \cup B) = P(A) + P(B)$$

Probability is a measure

σ algebra - collection of all meaningful statements

start from simple statements

Probability can be interpreted as long run frequency.

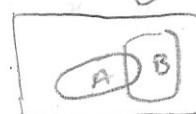
If we repeat the experiment a large # of times,

$P(A)$ can be interpreted as how often A happens.

Conditional Probability:

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

Geometry



$$P(A|B)$$

↳ falls into B (new Ω)

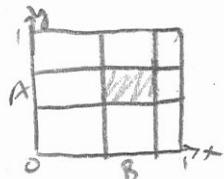
long-run frequency:

how often A happens when B happens

Independence

A and B are independent. $P(A|B) = P(A)$. $P(B|A) = P(B)$

Ⓐ Ⓛ disjoint \neq independent



$$|A \cap B| = |A| \cdot |B| \quad P(A \cap B) = P(A)P(B)$$

Random Variables:

Population $\Omega = \{w_1, w_2, \dots, w_m\}$

Randomly sample a person

X = height of this person. \Rightarrow height = $X(w)$

Randomly sample $w \in \Omega$

Experiment \rightarrow outcome $\xrightarrow{\text{label}}$ number

$A: X > 6 \text{ ft} \Rightarrow A = \{w : X(w) > 6\} \subset \Omega \Rightarrow P(A) = P(X > 6) = P(\{w : X(w) > 6\})$

Event: Statement about outcome

w/ numbers: inequalities, equation

Basic Events:

Discrete: $X = x$ $\{w : X(w) = x\}$

Continuous: $X \in (x, x + \Delta x)$ $\{w : X(w) \in (x, x + \Delta x)\}$

In general, a special event: $X \leq x$
 \hookrightarrow_{RV} specific value

$F(x) = P(X \leq x)$ CDF goes from 0 to 1.

Basic Events

Discrete: $X=x$ $P(X=x) = p(x)$ PMF

Continuous: $X \in (x, x+\Delta x)$ $P(X \in (x, x+\Delta x)) \approx f(x) \Delta x$ for small Δx

$f(x) = \lim_{\Delta x \rightarrow 0} \frac{P(X \in (x, x+\Delta x))}{\Delta x}$ PDF

PDF: density = $\frac{\# \text{ of people in } (x, x+\Delta x)}{\Delta x}$

Summarizing $p(x)$ or $f(x)$:

Expectation:

$$E(X) = \begin{cases} \sum_x x \cdot p(x) & \text{- Average} \\ \int x f(x) dx & \text{- long run average} \end{cases}$$

Variance/Volatility:

$$\text{Var}(X) = E[(X - E(X))^2] = E[X^2] - E[X]^2$$

$$aX+b \Rightarrow E[aX+b] = aE[X]+b$$

$$\text{var}(aX+b) = a^2 \text{Var}[X]$$

Two Random Variables

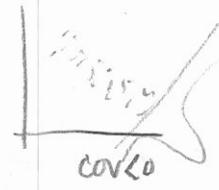
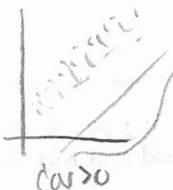
$$\text{COV}(X, Y) = E((X - \mu_X)(Y - \mu_Y))$$

$$X+Y: E[X+Y] = E[X] + E[Y]$$

$$\begin{aligned} \text{VAR}[X+Y] &= E[((X+Y) - (\mu_X + \mu_Y))^2] = E[((X - \mu_X) + (Y - \mu_Y))^2] \\ &= E[(X - \mu_X)^2 + (Y - \mu_Y)^2 + 2(X - \mu_X)(Y - \mu_Y)] \\ &= \text{Var}(X) + \text{Var}(Y) + 2 \text{COV}(X, Y) \end{aligned}$$

If X & Y are independent, $E(XY) = E[X]E[Y]$, $\text{COV}(X, Y) = 0$

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, E[\bar{X}] = \mu, \text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$



STATS 105

10/08

$$(X, Y) \sim f(x, y) \quad f(x, y) = \lim_{\substack{\Delta x \rightarrow 0 \\ \Delta y \rightarrow 0}} \frac{P(X \in (x, x+\Delta x) \text{ and } Y \in (y, y+\Delta y))}{\Delta x \Delta y}$$

$$\text{COV}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$$

$$\begin{array}{|c|c|} \hline & + \\ \hline + & + \\ \hline \end{array}$$

$\text{COV} > 0$

$$\begin{array}{|c|c|} \hline & - \\ \hline + & + \\ \hline \end{array}$$

$\text{COV} < 0$

$$\begin{array}{|c|c|} \hline & + \\ \hline + & - \\ \hline \end{array}$$

$\text{COV} = 0$

$$\text{CORR}(X, Y) = \frac{\text{COV}(X, Y)}{\sqrt{\text{Var}(X)} \sqrt{\text{Var}(Y)}} = \text{COV}\left(\frac{X - \mu_X}{\sigma_X}, \frac{Y - \mu_Y}{\sigma_Y}\right)$$

Vector Representation

$$\begin{array}{|c|c|c|} \hline n & X & Y \\ \hline 1 & X(1) & Y(1) \\ \hline 2 & X(2) & Y(2) \\ \hline \vdots & \vdots & \vdots \\ \hline w & X(w) & Y(w) \\ \hline \vdots & \vdots & \vdots \\ \hline m & X(m) & Y(m) \\ \hline \end{array} \rightarrow \begin{array}{|c|c|c|} \hline n & \tilde{X} & \tilde{Y} \\ \hline 1 & X(1) - \mu_X & Y(1) - \mu_Y \\ \hline 2 & X(2) - \mu_X & Y(2) - \mu_Y \\ \hline \vdots & \vdots & \vdots \\ \hline w & \cancel{X(w) - \mu_X} & Y(w) - \mu_Y \\ \hline \vdots & \vdots & \vdots \\ \hline m & X(m) - \mu_X & Y(m) - \mu_Y \\ \hline \end{array}$$

Recall:

$$\vec{a} = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix} \quad \vec{b} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix} \quad \langle \vec{a}, \vec{b} \rangle = a_1 b_1 + a_2 b_2 + \dots + a_n b_n = \sum_{i=1}^n a_i b_i$$

$$|\vec{a}|^2 = \langle \vec{a}, \vec{a} \rangle = \sum_{i=1}^n a_i^2$$

$$\cos \theta = \frac{\langle \vec{a}, \vec{b} \rangle}{|\vec{a}| \cdot |\vec{b}|} = \frac{\frac{1}{M} \langle \vec{a}, \vec{b} \rangle}{\sqrt{\frac{1}{M} |\vec{a}|^2} \sqrt{\frac{1}{M} |\vec{b}|^2}} = \frac{\text{COV}(X, Y)}{\sqrt{\text{Var}(X)} \sqrt{\text{Var}(Y)}} \in (-1, 1)$$

$$\text{Corr}(X, Y) = \cos \theta$$

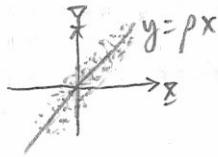
$$\frac{1}{M} \langle \vec{a}, \vec{b} \rangle = \frac{1}{M} \sum_{w=1}^M (X(w) - \mu_X)(Y(w) - \mu_Y) = E[(X - \mu_X)(Y - \mu_Y)] = \text{COV}(X, Y)$$

$$\frac{1}{M} |\vec{a}|^2 = \frac{1}{M} \sum_{w=1}^M (X(w) - \mu_X)^2 = E[(X - \mu_X)^2] = \text{Var}(X)$$

$$\text{Corr} = 1: \quad \vec{a} \quad \vec{b} = \lambda \vec{a}, \lambda > 0$$

For Simplicity Assume

$$E(\bar{X}) = E(\bar{Y}) = 0 \quad \text{and} \quad \text{Var}(\bar{X}) = \text{Var}(\bar{Y}) = 1$$



Find p , $\min E[(\bar{Y} - p\bar{X})^2]$ (least square estimation)

$$= \frac{1}{M} \sum_{m=1}^M (\bar{Y}(m) - p\bar{X}(m))^2 = \frac{1}{M} |\vec{b} - p\vec{a}|^2 \quad \frac{|p\vec{a}|}{|\vec{b}|} = |\vec{a}| = \cos \theta$$

$$1 - p^2 = \frac{|\vec{b} - p\vec{a}|^2}{|\vec{b}|^2} = \frac{E[(\bar{Y} - p\bar{X})^2]}{\text{Var}(\bar{Y})}$$

$$p = \cos \theta = \text{corr}(\bar{X}, \bar{Y})$$

$$R^2 = \cos^2 \theta = \frac{|p\vec{a}|^2}{|\vec{b}|^2} = \% \text{ variation of } \bar{Y} \text{ explained by linear relationship}$$

↳ strength of linear relationship

$$\text{Var}(\bar{Y}) = \text{Var}(p\bar{X}) + \text{Var}(\varepsilon)$$

$$\begin{array}{c} \vec{b} \\ \perp \theta \\ \vec{p}\vec{a} \end{array} \quad \varepsilon = \vec{b} - p\vec{a} \quad \bar{Y} = p\bar{X} + \varepsilon$$

Independently & Identically Distributed (iid) Random Variables

$$\bar{X}_1, \bar{X}_2, \dots, \bar{X}_n \stackrel{iid}{\sim} f(x)$$

If $\bar{X} \sim f(x)$

$$\begin{aligned} \mu = E(\bar{X}) \\ \sigma^2 = \text{Var}(\bar{X}) \end{aligned} \Rightarrow S = \sum_{i=1}^n \bar{X}_i, \bar{\bar{X}} = \frac{S}{n}$$

$$E(S) = E\left[\sum_{i=1}^n \bar{X}_i\right] = \sum_{i=1}^n E(\bar{X}_i) = n\mu$$

$$\text{Var}(S) = \text{Var}\left(\sum_{i=1}^n \bar{X}_i\right) = \sum \text{Var}(\bar{X}_i) + 2 \sum_{i < j} \text{Cov}(\bar{X}_i, \bar{X}_j) = n\sigma^2$$

$$E(\bar{\bar{X}}) = \frac{E(S)}{n} = \frac{n\mu}{n} = \mu$$

$$\text{Var}(\bar{\bar{X}}) = \text{Var}\left(\frac{S}{n}\right) = \frac{1}{n^2} \text{Var}(S) = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}$$

Law of Large Numbers

$$P(|\bar{X} - \mu| > \varepsilon) \xrightarrow[n \rightarrow \infty]{} 0$$

Central Limit Theorem

$$\begin{aligned} \sqrt{n}(\bar{X} - \mu) &\rightarrow E[\sqrt{n}(\bar{X} - \mu)] = 0 \\ \text{Var}[\sqrt{n}(\bar{X} - \mu)] &= \sigma^2 \end{aligned}$$

Follows Normal distribution

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

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

IF $Y \sim N(0, \sigma^2)$

$$P(\sqrt{n}(\bar{X} - \mu) \leq y) \rightarrow P(Y \leq y), \forall y$$

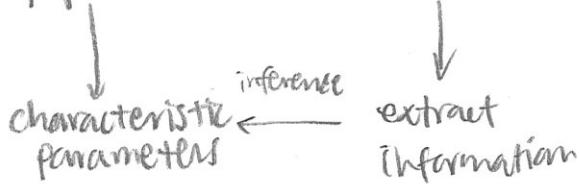
$$\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \rightarrow N(0, 1).$$

$$z \sim N(0, 1)$$

$$f(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) \cdot P\left(\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \leq z\right) \rightarrow P(z \leq z)$$

STATISTICS

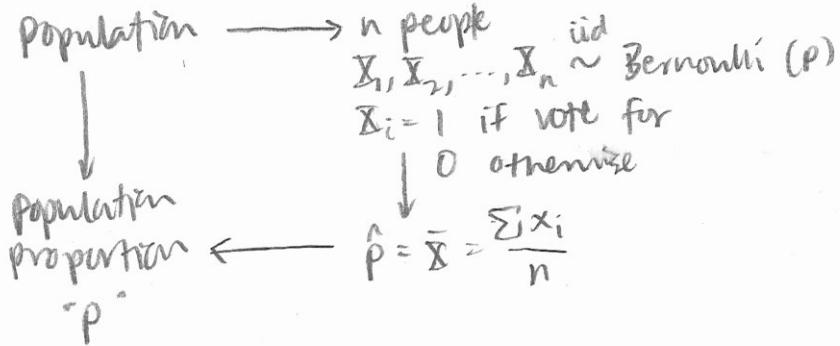
population \longrightarrow sample



Two Purposes

- understanding
- predictions

Estimate Population proportion



Bernoulli RV

$$P(X_i=1)=p \quad E[X_i] = 1 \cdot p + 0 \cdot (1-p) = p.$$

$$P(X_i=0)=1-p \quad \text{Var}[X_i] = (1-p^2) \cdot p + (0-p)^2 \cdot (1-p) = p(1-p)$$

\hat{p} is a random variable.

Hypothetic Repetition: Sample X_1, \dots, X_n again, get a different \hat{p} .

$$E[\hat{p}] = E[\bar{x}] = \mu = p \quad (\text{unbiased})$$

$$\text{Var}[\hat{p}] = \text{Var}[\bar{x}] = \frac{\sigma^2}{n} = \frac{p(1-p)}{n} \quad \text{s.d. } (\hat{p}) = \sqrt{\text{Var}(\hat{p})} = \sqrt{\frac{p(1-p)}{n}} \quad \text{if you don't know } p, \text{ plug in } \hat{p}.$$

$$\hat{p}_1 = \frac{s}{n} = \bar{x} \quad \hat{p}_2 = s + \frac{\sqrt{n}}{2} \quad \text{Bias: } E[\hat{p}_2] - p = \frac{\frac{1}{2} - p}{\sqrt{n+1}}$$

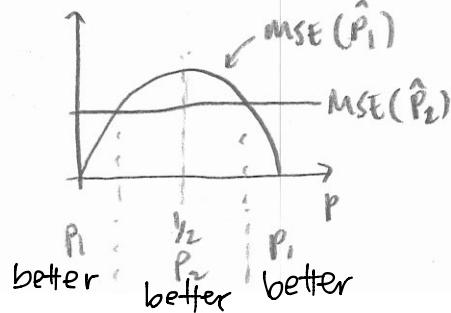
$$E[\hat{p}_2] = \frac{E[s + \frac{\sqrt{n}}{2}]}{n + \sqrt{n}} = \frac{E[s] + \frac{\sqrt{n}}{2}}{n + \sqrt{n}} = \frac{np + \frac{\sqrt{n}}{2}}{n + \sqrt{n}} = \frac{\sqrt{n}p + \frac{1}{2}}{\sqrt{n+1}}$$

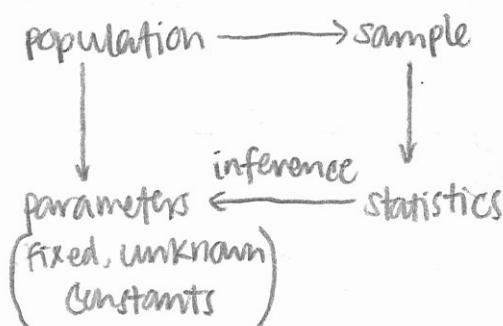
"unbiased" "biased"

$$\text{Var}(\hat{p}_2) = \frac{\text{Var}(s)}{(n+\sqrt{n})^2} = \frac{np(1-p)}{(n+\sqrt{n})^2} = \frac{p(1-p)}{(\sqrt{n}+1)^2} < \text{Var}(\hat{p}_1) \quad \text{trade off between bias and variance.}$$

$$\text{MSE}(\hat{p}_1) = \frac{p(1-p)}{n}$$

$$\text{MSE}(\hat{p}_2) = \left(\frac{1}{2} - p\right)^2 + \frac{p(1-p)}{(\sqrt{n}+1)^2} = \frac{\frac{1}{4}}{(\sqrt{n}+1)^2}$$



Point Estimation

$X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f(x, \theta)$
sample, data. model of population

$\hat{\theta} = h(X_1, X_2, \dots, X_n)$
Estimate is a function of the data
 $h(\cdot)$: estimator
Value of $\hat{\theta}$: estimate

X_1, X_2, \dots, X_n are random \Rightarrow different under hypothetical repetitions
and as a result, $\hat{\theta}$ is also random.

bias: $E[\hat{\theta}] - \theta$. unbiased if $E[\hat{\theta}] = \theta$. \checkmark

variance: $\text{Var}(\hat{\theta})$

Mean Square Error (MSE): $E[(\hat{\theta} - \theta)^2] = \text{bias}(\hat{\theta})^2 + \text{Var}(\hat{\theta})$ \checkmark

$$\begin{aligned} \text{Proof: } \mu = E[\hat{\theta}] \Rightarrow E[(\hat{\theta} - \theta)^2] &= E[((\hat{\theta} - \mu) - (\mu - \theta))^2] \\ &= E[(\hat{\theta} - \mu)^2 + (\mu - \theta)^2 - 2(\hat{\theta} - \mu)(\mu - \theta)] \\ &= E[\hat{\theta} - \mu]^2 + E[\mu - \theta]^2 - 2E[(\hat{\theta} - \mu)(\mu - \theta)] \\ &= \text{Var}(\hat{\theta}) + \underbrace{[E(\hat{\theta}) - \theta]^2}_{\text{bias}} - 2(\mu - \theta)E[\hat{\theta} - \mu] \\ &\quad \hookrightarrow = 0 \end{aligned}$$

Example:

$X_1, X_2, \dots, X_n \sim \text{Bernoulli}(\theta)$

$X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$

$X_1, X_2, \dots, X_n \sim \text{Exponential}(\lambda)$ $f(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$ (waiting time)

Method of Moments

estimating equation \rightarrow solution: $\hat{\theta}$
Find θ . $E_{\theta}(x) = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$. Another notation $\langle X \rangle_{\text{fo}} = \langle X \rangle_{\text{data}}$
(model)

$X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$

$$E_{\mu}[\bar{X}] = \mu$$

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

solution

$X_1, X_2, \dots, X_n \sim \text{Exp}(\lambda) . E[X] = \frac{1}{\lambda}$ Get this by trying to match model to data.

estimating equation: $\frac{1}{\lambda} = \bar{x}$

$$\text{solution: } \hat{\lambda} = \frac{1}{\bar{x}}$$

$X_1, X_2, \dots, X_n \sim \text{Bernoulli}(p) . E[X] = p$

estimating equation: $p = \bar{x}$

$$\text{solution: } \hat{p} = \bar{x}$$

Two Unknown Parameters

$X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$

\hookrightarrow unknown

$$E[X] = \mu$$

$$\text{① } \mu = \bar{x} = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{solution: } \hat{\mu} = \bar{x}$$

$$E[X^2] = E[\bar{X}]^2 + \text{Var}(\bar{X}) = \mu^2 + \sigma^2$$

$$\text{② } E[X^2] = \frac{1}{n} \sum_{i=1}^n X_i^2$$

$$\mu^2 + \sigma^2 = \frac{1}{n} \sum_{i=1}^n X_i^2$$

$$\text{solution: } \hat{\sigma}_2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{x})^2$$

In general: $X_1, X_2, \dots, X_n \sim f(x, \theta)$

$\theta = (\theta_1, \theta_2, \dots, \theta_m)$ needs "m" estimating equations

Estimating Equations: $E_\theta(\bar{X}^k) = \frac{1}{n} \sum_{i=1}^n \bar{X}_i^k$

Solutions: $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_m)$

More Generally: $E_\theta(h_k(x)) = \frac{1}{n} \sum_{i=1}^n h_k(x_i), k=1, \dots, m$

Best Estimation Equation: maximum likelihood estimator

$X_1, X_2, \dots, X_n \stackrel{iid.}{\sim} f(x, \theta_{true})$

Likelihood: $L(\theta) = \prod_{i=1}^n f(X_i, \theta)$ $\hat{\theta}_{MLE} = \arg \max_{\theta} L(\theta)$

\hookrightarrow independent (X_1, X_2, \dots, X_n)

$L(\theta)$ measures plausibility of θ in explaining X_1, X_2, \dots, X_n

MLE: most plausible explanation

log-likelihood

$$l(\theta) = \text{Log}(f(\theta)) = \sum_{i=1}^n \log f(x_i; \theta)$$

$$l'(\theta) = 0$$

$$\sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(x_i; \theta) = 0 \quad (\text{MLE equation})$$

Maximum Likelihood Estimator (MLE)

$X_1, X_2, \dots, X_n \stackrel{iid}{\sim} p(x; \theta)$ likelihood: $L(\theta) = \prod_{i=1}^n p(X_i; \theta)$

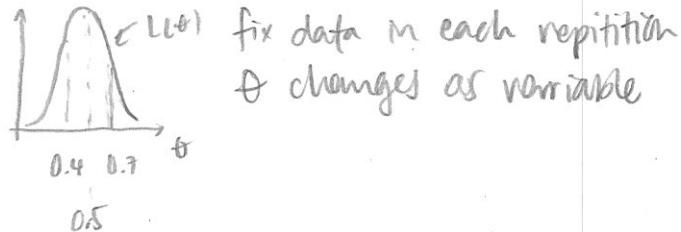
In general: $L(\theta) = \text{prob(data}; \theta)$

plausibility of θ for explaining the data.

$\text{Prob}(\text{data}; \theta_1) > \text{Prob}(\text{data}; \theta_2) \Rightarrow \theta_1$ is more plausible than θ_2 .

We want the most plausible value for θ .

$$\hat{\theta}_{\text{MLE}} = \arg \max L(\theta)$$



If data = (X_1, X_2, \dots, X_n) are independent, recall A & B independent

$$P(A \cap B) = P(A) P(B)$$

$$P(X_1, X_2, \dots, X_n; \theta) = \prod_{i=1}^n p(X_i; \theta)$$

$$\text{log-like: } l(\theta) = \text{Log } L(\theta) = \sum \log p(X_i; \theta)$$

$$\text{derivative: } l'(\theta) = \sum \frac{\partial}{\partial \theta} \log p(X_i; \theta)$$

$$\text{Estimating eqn: } l'(\theta) = 0$$

$$\text{Solve: } \hat{\theta}_{\text{MLE}} \quad l''(\hat{\theta}_{\text{MLE}}) < 0$$

Continuous

$$X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f(x; \theta) \text{ density}$$

$$\text{If } X \sim f(x; \theta), \text{ prob}(X \in (x, x + \Delta x)) = f(x; \theta) \Delta x$$

$$L(\theta) = \prod_{i=1}^n \left(f(x_i; \theta) \Delta x \right)$$

constant ↓
precision

Example 1:

$$X_1, X_2, \dots, X_n \sim \text{Bernoulli}(p)$$

$$X \sim \text{Ber}(p)$$

	0	1
Prob	$(1-p)$	p

$$P(X=x) = p^x (1-p)^{1-x}$$

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

$$\ell(p) = \sum_{i=1}^n (x_i \log p + (1-x_i) \log(1-p))$$

$$\ell'(p) = \sum_{i=1}^n \left(x_i \cdot \frac{1}{p} + (1-x_i) \cdot \frac{1}{1-p} \cdot (-1) \right)$$

Estimating equation: $\ell'(p) = 0$

Let $S = \sum X_i = \# \text{ of } 1's$. $\frac{S}{p} - \frac{n-S}{1-p} = 0 \Rightarrow S(1-p) = (n-S)p \Rightarrow \hat{p} = \frac{S}{n}$ freq of 1s.

$$\ell''(p) = -\frac{S}{p^2} - \frac{n-S}{(1-p)^2} < 0$$

Example 2

$$X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$$

$$f(x, \mu) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$L(\mu) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i-\mu)^2}{2\sigma^2}\right) = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2\right)$$

$$\min R(\mu) = \sum_{i=1}^n (x_i - \mu)^2 \Rightarrow R'(\mu) = \sum_{i=1}^n 2(x_i - \mu)(-1) = 0$$

$$\sum_i (x_i) - n\mu = 0 \Rightarrow \hat{\mu} = \frac{\sum x_i}{n}$$

Mixture Model

$$f(x, \theta) = \lambda \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left(-\frac{(x-\mu_1)^2}{2\sigma_1^2}\right) + (1-\lambda) \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{(x-\mu_2)^2}{2\sigma_2^2}\right)$$

$$\theta = (\lambda, \mu_1, \sigma_1^2, \mu_2, \sigma_2^2) \quad (\text{Method of moments could be used})$$

$$L(\theta) = \prod_{i=1}^n f(x_i; \theta) \quad (\text{Clustering})$$

Regression

Observations	Predictor	Response
1	X_1	Y_1
2	X_2	Y_2
:	:	:
n	X_n	Y_n

{ error $\sim N(0, \sigma^2)$

$$Y_i \sim \alpha + \beta X_i + \varepsilon_i$$

r fixed
 $\varepsilon_i \sim N(\alpha + \beta X_i, \sigma^2)$

$$L(\theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(Y_i - (\alpha + \beta X_i))^2}{2\sigma^2}\right)$$

$$L(\theta) = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (\bar{Y}_i - (\alpha + \beta \bar{X}_i))^2 \right)$$

$$\min R(\theta) = \sum_{i=1}^n (\bar{Y}_i - (\alpha + \beta \bar{X}_i))^2$$

$\hat{\alpha}, \hat{\beta}$ Least Square Estimators

Logistic Regression

$$y_i \in \{0, 1\} \quad y_i \sim \text{Ber}(p) \quad \log \frac{p_i}{1-p_i} = \alpha + \beta X_i$$

MLE

Example:

$$X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$$

$$L(\theta) = \prod_{i=1}^n f(x_i; \theta) \Rightarrow L(\theta) = \left(\frac{1}{2\pi\sigma^2} \right)^n \exp \left(-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2} \right), \quad \theta = (\mu, \sigma^2)$$

$$\ell(\theta) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$= -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \ln\sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$= \text{constant} - \frac{n}{2} \ln\sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\frac{\partial \ell}{\partial \mu} = -\frac{1}{2\sigma^2} \sum_{i=1}^n 2(x_i - \mu)(-1) \Rightarrow \frac{\partial \ell}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0 \quad (1)$$

$$\frac{\partial \ell}{\partial \sigma^2} = -\frac{n}{2} \cdot \frac{1}{\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_{i=1}^n (x_i - \mu)^2 = 0 \quad (2)$$

$$(1) \sum_{i=1}^n (x_i - \mu) = 0 \Rightarrow \sum_{i=1}^n x_i - n\mu = 0 \Rightarrow \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$$

$$(2) -n(\sigma^2) + \sum (x_i - \mu)^2 = 0 \Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\ell'(\theta) = \begin{pmatrix} \frac{\partial \ell}{\partial \mu} \\ \frac{\partial \ell}{\partial \sigma^2} \end{pmatrix}, \quad \ell''(\theta) = \begin{pmatrix} \frac{\partial^2 \ell}{\partial \mu^2} & \frac{\partial^2 \ell}{\partial \mu \partial \sigma^2} \\ \frac{\partial^2 \ell}{\partial \sigma^2 \partial \mu} & \frac{\partial^2 \ell}{\partial (\sigma^2)^2} \end{pmatrix}, \quad \ell''(\hat{\theta}) < 0.$$

negative definite

Example: Normal linear Regression

Observations	Predictor	Response	$\{Y_1, Y_2, \dots, Y_n \stackrel{iid}{\sim} N(\mu_i, \sigma^2)$
1	X_1	y_1	
2	X_2	y_2	
\vdots	fixed	random	$\mu_i = x_i \beta$
n	X_n	y_n	$y_i = x_i \beta + \varepsilon_i, \varepsilon_i \sim N(0, \sigma^2)$

$$f(y|x, \theta) = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) \exp \left(-\frac{(y - x_i \beta)^2}{2\sigma^2} \right) \quad E(\varepsilon_i) = 0$$

$$L(\theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(y_i - \mu_i)^2}{2\sigma^2} \right), \quad \mu_i = x_i \beta \quad \text{Var}(\varepsilon_i) = \sigma^2$$

$$= \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp \left(-\frac{\sum_{i=1}^n (y_i - \mu_i)^2}{2\sigma^2} \right) \quad E[\mu_i + \varepsilon_i] = \mu_i$$

$$\text{Var}[\mu_i + \varepsilon_i] = \sigma^2$$

$$l(\theta) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum (y_i - x_i \beta)^2$$

$$\max l(\theta) \Rightarrow \min \sum (y_i - x_i \beta)^2 = R(\beta)$$

\downarrow max likelihood \downarrow least squares

$$\frac{\partial l}{\partial \beta} = -\frac{1}{2\sigma^2} \sum_{i=1}^n 2(y_i - x_i \beta)(-x_i) \quad (1)$$

$$\frac{\partial l}{\partial \sigma^2} = -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^2} \sum (y_i - x_i \beta)^2 \quad (2)$$

$$(1) \sum (y_i - x_i \beta)x_i = 0 \Rightarrow \sum y_i x_i - \sum x_i^2 \beta = 0 \Rightarrow \hat{\beta} = (\sum x_i^2)^{-1} \sum x_i y_i$$

$$(2) \hat{\sigma}_e^2 = \frac{\sum (y_i - x_i \beta)^2}{n}$$

n -Dim

$$\|\vec{x} - \vec{\beta}\|^2 \text{ min } \quad \langle \vec{x}, \vec{y} - \vec{\beta} \rangle = 0 \quad (\text{inner product of orthogonal})$$

$$\sum x_i(y_i - x_i \beta) = 0$$

Multiple Linear Regression:

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i = y_i = \vec{x}_i^T \vec{\beta} + \epsilon_i$$

$$\hat{\beta} = \left(\sum \vec{x}_i \vec{x}_i^T \right)^{-1} \left(\sum \vec{x}_i y_i \right) \quad \vec{x}_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix} \quad \vec{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix}$$

Example:

$$y_i \sim \text{Ber}(p_i) \quad (x_i - \text{height}, y_i - \text{gender})$$

$$\log \left(\frac{p_i}{1-p_i} \right) = x_i \beta \Leftrightarrow p_i = \frac{e^{x_i \beta}}{1+e^{x_i \beta}} = \frac{1}{1+e^{-x_i \beta}}$$

$$L(\beta) = \prod_{i=1}^n p(y_i | x_i, \beta) = \prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i} = \prod_{i=1}^n \left(\frac{e^{x_i \beta}}{1+e^{x_i \beta}} \right)^{y_i} \left(\frac{1}{1+e^{x_i \beta}} \right)^{1-y_i}$$

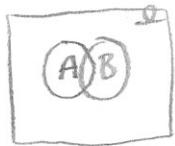
$$= \prod_{i=1}^n \left(\frac{e^{x_i \beta}}{1+e^{x_i \beta}} \right)^{y_i} \Rightarrow L(\beta) = \sum_{i=1}^n (y_i x_i \beta - \ln(1+e^{x_i \beta}))$$

$$L'(\beta) = 0 \Rightarrow \hat{\beta}_{MLE}; \quad L'(\beta) = \sum_{i=1}^n (y_i x_i - \frac{e^{x_i \beta}}{1+e^{x_i \beta}} \cdot x_i) \quad \text{want } p_i \text{ close to } y_i$$

Bayes Rule

Conditional probability: $P(A|B) = \frac{P(A \cap B)}{P(B)}$

Geometric:



$$P(A) = \frac{|A|}{|S|}$$

$P(A|B)$ (falls into B)

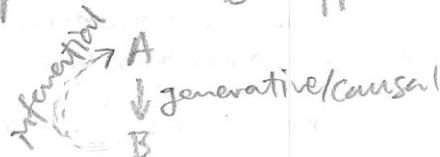
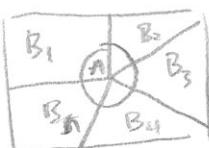
$$P(A|B) = \frac{|A \cap B|}{|B|} \text{ how much } A \text{ occupies } B.$$

slogan: conditional probability is just a regular probability except B is a regular probability in a conditional probability: $P(A) = P(A|S)$.

Frequency: $P(A|B)$: how often A happens when B happens.

Chain Rule: $P(A \cap B) = P(B) P(A|B)$

$$= P(A) P(B|A)$$

Rule of Total Probability

Given: $P(B_i)$; $i = 1, \dots, n$

$$P(A) = \sum_{i=1}^n P(A \cap B_i) \rightarrow \text{Additivity}$$

$$= \sum_{i=1}^n P(B_i) P(A|B_i)$$

$$\begin{aligned} P(B_i|A) &= \frac{P(B_i \cap A)}{P(A)} \\ &= \frac{P(B_i) P(A|B_i)}{\sum_j P(B_j) P(A|B_j)} \end{aligned}$$

Example:

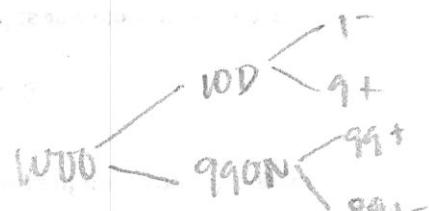
1% of population has disease.

$$\text{Medical Test: } P(+ | \text{Disease}) = 90\%$$

$$P(- | \text{Disease}) = 10\%$$

$$P(- | \text{Not Disease}) = 90\%$$

$$P(+ | \text{Not Disease}) = 10\%$$



$$\frac{1}{99+1} = \frac{1}{100}$$

$$P(D|+) = \frac{P(D \cap +)}{P(+)} = \frac{P(D) P(+|D)}{P(D) P(+|D) + P(N) P(+|N)} = \frac{P(D) P(+|D)}{P(D) P(+|D) + P(N) P(+|N)}$$

Experiment \rightarrow outcome \rightarrow #

↑
Statement what
↑
equation or
inequality about

$$(X, Y) \sim P(X, Y) = P(X=x, Y=y)$$

$$P(X) = P(X=x) = \sum_y P(x, y) \quad P(x|y) = P(X=x | Y=y) = \frac{P(x, y)}{P(y)}$$

$$P(Y=y) = \sum_x P(x, y) = \sum_x P(x) P(y|x)$$

$$p(x|y) = \frac{P(x) P(y|x)}{\sum_x P(x) P(y|x)}$$

p(x) prior probability
p(x|y) posterior probability

$$(X, Y) \sim f(x, y)$$

$$P(X \in (x, x+\Delta x) \cap Y \in (y, y+\Delta y)) = f(x, y) \Delta x \Delta y$$

$$f(x) = \int f(x, y) dy \quad f(y|x) = \frac{f(x, y)}{f(x)}$$

$$f(y) = \int f(x, y) dx \quad f(x|y) = \frac{f(x, y)}{f(y)}$$

$$f(y) = \int f(x, y) dx = \int f(x) f(y|x) dx$$

$$f(x|y) = \frac{f(x, y)}{\int f(x) f(y|x) dx}$$

Mixture Model:

$$\begin{aligned} X \in \{0, 1\} & \text{ gender} & f(y|0) = f_0(y) & X \sim \text{Ber}(\lambda) & P(1) = (\lambda) \\ Y \text{ continuous} & \text{ height} & f(y|1) = f_1(y) & & P(0) = (1-\lambda) \end{aligned}$$

$$f(y) = \sum_x p(x) f(y|x) = p(1) f(y|1) + p(0) f(y|0) = \lambda f_1(y) + (1-\lambda) f_0(y)$$

$$p(1|y) = \frac{f_1(y)}{f(y)} = \frac{p(1) f(y|1)}{p(1) f(y|1) + p(0) f(y|0)} = \frac{\lambda f_1(y)}{\lambda f_1(y) + (1-\lambda) f_0(y)}$$

Normal Inference

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

$$Y = X + \varepsilon, \quad \varepsilon \sim N(0, \tau^2)$$

$$(Y | X=x) \sim N(x, \tau^2)$$

$$f(x|y) = \frac{f(x) f(y|x)}{\int f(x) f(y|x) dx}$$

Probability Calculations

Events:

chain Rule: $P(A \cap B) = P(B) P(A|B)$

Total:  $P(A) = \sum_{i=1}^n P(A \cap B_i) = \sum_{i=1}^n P(B_i) P(A|B_i)$

Bayes Rule: $P(B_i|A) = \frac{P(A \cap B_i)}{P(A)} = \frac{P(B_i) P(A|B_i)}{\sum_{j=1}^n P(B_j) P(A|B_j)}$

Random Variable:

Discrete: chain Rule: $P(x,y) = P(x) P(y|x)$

marginality: $P(y) = \sum_x P(x,y) = \sum_x P(x) P(y|x)$

Bayes Rule: $P(x|y) = \frac{P(x,y)}{P(y)} = \frac{P(x) P(y|x)}{\sum_x P(x) P(y|x)}$

Continuous:

$$f(x,y) = f(x) f(y|x)$$

$$f(y) = \int f(x,y) dx = \int f(x) f(y|x) dx$$

$$f(x|y) = \frac{f(x,y)}{f(y)} = \frac{f(x) f(y|x)}{\int f(x) f(y|x) dx}$$

Mixture Model Example:

$$X \sim \text{Ber}(\lambda) \quad \begin{array}{c|cc} x & 0 & 1 \\ \hline p(x) & 1-\lambda & \lambda \end{array} \quad P(1) = \lambda \quad P(0) = 1-\lambda$$

$$[Y|X=1] \sim f_1(y) = f(y|1)$$

$$f(y,x) = f(y|x) p(x)$$

$$[Y|X=0] \sim f_0(y) = f(y|0)$$

$$f(y,1) = f(y|1) p(1) = \lambda f_1(y)$$

$$f(y,0) = f(y|0) p(0) = (1-\lambda) f_0(y)$$

$$f(y) = \sum_x f(y,x) = \sum_x f(y|x) p(x) = \lambda f_1(y) + (1-\lambda) f_0(y)$$

$$P(x|y) = \frac{f(y,x)}{f(y)} = \frac{f(y|x) p(x)}{f(y)}$$

$$P(1|y) = \frac{f_1(y) \lambda}{\lambda f_1(y) + (1-\lambda) f_0(y)}$$

$$P(0|y) = \frac{(1-\lambda) f_0(y)}{\lambda f_1(y) + (1-\lambda) f_0(y)}$$

$$f(y,x) = (\lambda f_1(y))^x (1-\lambda) f_0(y)^{1-x}$$

Current Population:

proportion

$$M = 300 \text{ mil}, \text{ male} = \lambda, \text{ female} = (1-\lambda)$$

$$\# \text{ of males: } M \cdot \lambda \rightarrow \# \text{ of males in } [y, y + \Delta y] = M \lambda f_1(y) \Delta y$$

$$\text{females: } M \cdot (1-\lambda) \rightarrow \# \text{ of females in } [y, y + \Delta y] = M (1-\lambda) f_0(y) \Delta y$$

$$\# \text{ of ppl in } [y, y + \Delta y] = M (\underbrace{\lambda f_1(y) + (1-\lambda) f_0(y)}_{f(y)}) \Delta y$$

Among all ppl in $[y, y + \Delta y]$:

f(y)

proportion male: $\frac{M \lambda f_1(y) \Delta y}{M (\lambda f_1(y) + (1-\lambda) f_0(y)) \Delta y}$

$$= p(1|y) = p(X=1 | Y \in [y, y + \Delta y])$$

Example:

$$(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n) \quad X_i \sim \text{Ber}(\lambda)$$

$$[Y_i | X_i=1] = f_1(y_i) \sim N(\mu_1, \sigma_1^2) \quad \Theta = (\lambda, \mu_1, \sigma_1^2, \mu_0, \sigma_0^2)$$

$$[Y_i | X_i=0] = f_0(y_i) \sim N(\mu_0, \sigma_0^2)$$

$$L(\theta) = \prod_{i=1}^n f(y_i, x_i) = \prod_{i=1}^n [\lambda f_1(y_i)]^{x_i} [(1-\lambda) f_0(y_i)]^{1-x_i}$$

$$l'(\theta) = 0 \Rightarrow \text{MLE} \quad \hat{\lambda} = \frac{\sum x_i}{n}, \quad \hat{\mu}_1 = \frac{\sum y_i x_i}{\sum x_i}, \quad \hat{\mu}_0 = \frac{\sum y_i (1-x_i)}{\sum (1-x_i)}$$

What if we do not observe X_i , $i=1, \dots, n$?

$$L(\theta) = \prod_{i=1}^n f(y_i) = \prod_{i=1}^n (\lambda f_1(y_i) + (1-\lambda) f_0(y_i))$$

Expectation Maximization (EM) Algorithm

E step: Given current parameter values

$$\lambda^{(t)}, \mu_0^{(t)}, \sigma_0^{(t)}, \mu_1^{(t)}, \sigma_1^{(t)}$$

Guess \hat{x}_i

$$= \frac{\lambda^{(t)} f_1(y_i)}{\lambda^{(t)} f_1(y_i) + (1-\lambda^{(t)}) f_0(y_i)}$$

$$M \text{ step: } \lambda^{(t+1)} = \frac{\sum \hat{x}_i}{n}$$

STATS 105 Homework Due Thurs after Midterm 10/31
 All homework problems: All calculations in notes. Don't skip any steps

Bayes Rule

$$\bar{X} \sim N(\mu, \sigma^2) \quad X \sim f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$\bar{Y} | \bar{X} = x \sim N(\alpha, \tau^2) \quad (\bar{Y} | \bar{X}) \sim f(y|x) = \frac{1}{\sqrt{2\pi\tau^2}} \exp\left(-\frac{(y-x)^2}{2\tau^2}\right)$$

$$[\bar{X}, \bar{Y}] \sim f(x, y) = f(y|x)f(x)$$

$$[\bar{X} | \bar{Y}] \sim f(x|y) = \frac{f(x,y)}{f(y)} \text{ density of } x \text{ w/ } y \text{ fixed.}$$

$$f(x|y) \propto f(x,y) \Rightarrow f(x|y) \propto \exp\left(-\frac{1}{2}\left(\frac{x^2}{\sigma^2} - \frac{2\mu}{\sigma^2}x + \frac{x^2}{\tau^2} - \frac{2y}{\tau^2}x\right)\right)$$

$$f(x|y) \propto \exp\left(-\frac{1}{2}\left(\frac{1}{\sigma^2} + \frac{1}{\tau^2}\right)x^2 - 2\left(\frac{\mu}{\sigma^2} + \frac{y}{\tau^2}\right)x\right) \\ \times \exp\left(-\frac{1}{2}\left(\frac{1}{\sigma^2} + \frac{1}{\tau^2}\right)\left(x - \frac{\mu + y}{\sigma^2 + \tau^2}\right)^2\right)$$

$$[\bar{X} | \bar{Y}] \sim f(x|y) \sim N\left(\frac{\mu + \frac{y}{\tau^2}}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}}, \frac{1}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}}\right)$$

fix $\sigma^2, \tau^2 \rightarrow 0 \Rightarrow$ posterior mean = y
 fix $\tau^2, \sigma^2 \rightarrow 0 \Rightarrow$ posterior mean = μ

↳ compromise between prior & evidence
 weighted by $1/\text{var} = \text{precision}$

Inference:

$$\bar{X}_1, \bar{X}_2, \dots, \bar{X}_n \sim N(\mu, \sigma^2) \xrightarrow{\text{unknown given}} \hat{\mu}_{MLE} = \bar{\bar{X}} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

Bayesian: $\xrightarrow{\text{given}}$

$$\mu \sim N(\alpha, \tau^2) \xrightarrow{\text{given}} [\mu | \bar{X}] \sim N\left(\frac{\alpha}{\tau^2} + \frac{\bar{X}n}{\sigma^2}, \frac{1}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}}\right)$$

$$n \rightarrow \infty \quad \xrightarrow{\downarrow} \bar{X}$$

Bayesian vs. Frequentist:

μ : speed of light

Frequentist: μ unknown constant

Bayesian: since unknown, random variable

$$L(\mu) = f(x_1, x_2, \dots, x_n | \mu) = \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{n}{2}} \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right)$$

$$\text{MLE: } \min \sum_{i=1}^n (x_i - \mu)^2 \rightarrow \hat{\mu} = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}$$

regularization by penalty (penalize large μ)

$$\min \left[\sum_{i=1}^n (x_i - \mu)^2 + \lambda \mu^2 \right]$$

→ prefer small μ

As if prior $\mu \sim f(\mu) \propto \exp(-\lambda \mu^2) \sim N(0, \frac{1}{2\lambda})$

$$\text{Posterior } [p(\mu|x)] \propto f(x_1, x_2, \dots, x_n | \mu) f(\mu) \propto \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 - \frac{\lambda \mu^2}{2\sigma^2} \right)$$

Recall Regression

obs	predictor	response
1	x_1	y_1
2	x_2	y_2
3	:	:
n	x_n	y_n

$$y_i = \beta x_i \Rightarrow \min_{\beta} \sum_{i=1}^n (y_i - \beta x_i)^2$$

if x_i is p dim & so is β

$$\min_{\beta} \left[\sum_{i=1}^n (y_i - \vec{x}_i^T \vec{\beta})^2 + \lambda |\vec{\beta}|^2 \right]$$

→ curve smooth.

Election:

State wise probability/proportion biased polls

P_i : probability in state i.

$$x_i | n_i, P_i \sim \text{Bin}(n_i, P_i)$$

$$P_i \sim N(\alpha, \gamma^2)$$

Borrow strength from other similar states.

Batting Averages $\frac{\text{hits}}{\text{attempts}}$

$$\text{player } i \quad x_i | n_i, P_i \sim \text{Bin}(n_i, P_i)$$

$$P_i \sim N(\alpha, \gamma^2)$$

$$\text{MLE: } \hat{P}_i = \frac{x_i}{n_i}$$

Bayesian: pull $\frac{x_i}{n_i}$ towards α

Problem 3:

$$(1) f(x, y) = f_x(x)f(y|x) \propto \exp\left(-\frac{(y-\mu)^2}{2(\sigma^2 + \tau^2)}\right) \Rightarrow N(\mu, \sigma^2 + \tau^2)$$

$$(2) f(x|y) = \frac{f(x,y)}{f(y)} \quad f(y) = \int f(x,y) dx$$

Problem 4:

$$(1) \text{MLE of } \mu. L(\mu) \Rightarrow l(\mu) \Rightarrow \hat{\mu}_{\text{MLE}}$$

$$(2) p(\mu), p(x_i|\mu) \text{ known, } p(\mu|x)?$$

$$p(\mu|x) \propto p(\mu) \underbrace{p(x|\mu)}_{\prod_{i=1}^n p(x_i|\mu)} \sim \text{Normal}$$

$$\text{Var}(\sum X_i) = \sum \text{Var}(X_i) + \sum_{i \neq j} \text{Cov}(X_i, X_j) \quad \frac{n(n-1)}{2} \text{ terms}$$

Midterm Review

$$E[aX+b] = aE[X]+b$$

$$\text{Var}[aX+b] = a^2 \text{Var}[X]$$

MLE for θ

$$L(\theta) \Rightarrow l(\theta) \Rightarrow \hat{\theta}_{\text{MLE}}$$

$$\text{Var}[X] = E[X^2] - E[X]^2$$

$$\text{Cov}[X, Y] = E[XY] - E[X]E[Y]$$

$$\text{Corr}[X, Y] = \frac{\text{Cov}[X, Y]}{\sqrt{\text{Var}[X]\text{Var}[Y]}}$$

$$\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}[X, Y]$$

$$(1) \text{ Bias}(\hat{\theta}) = E[\hat{\theta}] - \theta$$

$$\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2] = \text{Var}(\hat{\theta}) + \text{Bias}^2(\hat{\theta})$$

$$\text{Ber}(p) : p(X=1)=p, p(X=0)=1-p$$

$$\text{Bin}(n, p) : p(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$$

Normal:

$$\text{Exponential}(\lambda) \quad f(x) = \lambda e^{-\lambda x}$$

Confidence Interval

Probability Background

$$X_1, X_2, \dots, X_n \sim f(x)$$

for $X \sim f(x)$ $\mu = E[X]$, $\sigma^2 = \text{Var}[X]$

$$\bar{X} = \frac{\sum x_i}{n} \Rightarrow E[\bar{X}] = \mu, \text{Var}[\bar{X}] = \frac{\sigma^2}{n}$$

Law of Large

$$\bar{X} \xrightarrow{n \rightarrow \infty} \mu \quad P(|\bar{X} - \mu| > \epsilon) \xrightarrow{n \rightarrow \infty} 0$$

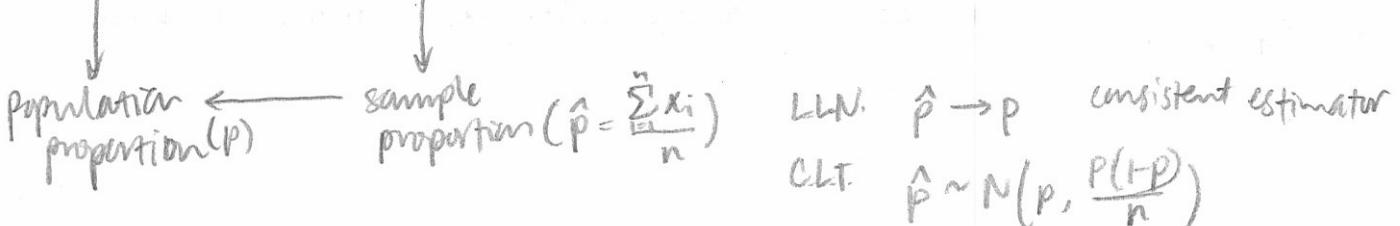
Central Limit Theorem

$$\sqrt{n}(\bar{X} - \mu) \rightarrow N(0, \sigma^2)$$

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

Inference of population proportion

population \longrightarrow sample $(\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n \sim \text{Ber}(p))$



$$s^2 = \frac{\hat{p}(1-\hat{p})}{n}, s(\hat{p}): \text{standard error}$$

$$\text{Approximately } \hat{p} \sim N(p, \sigma^2 = \frac{p(1-p)}{n}), s^2 = \hat{p}^2$$

Confidence interval of 95%

$$\hat{p} \pm 1.96 s(\hat{p}) \quad [\hat{p} - 1.96s(\hat{p}), \hat{p} + 1.96s(\hat{p})]$$

Margin of error = 1.96 s(\hat{p})

Example: Call WOO ppl. $\hat{p} = 0.53$

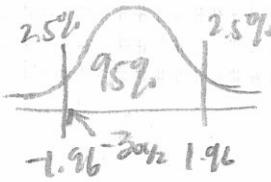
$$\text{CI} = 0.53 \pm 1.96 \left(\sqrt{\frac{0.53(1-0.53)}{WOO}} \right)$$

$$P([\hat{p} - 1.96 \sqrt{\frac{p(1-p)}{n}}, \hat{p} + 1.96 \sqrt{\frac{p(1-p)}{n}}] \ni p) = 95\%$$

\uparrow \downarrow \uparrow \downarrow

random fixed

$$Z = \frac{\hat{p} - p}{\sqrt{p(1-p)/n}} \sim N(0,1)$$



$$P(Z \in (-1.96, 1.96)) = 95\%$$

$$\hat{p} - 1.96 \sqrt{\frac{p(1-p)}{n}} \leq p \leq \hat{p} + 1.96 \sqrt{\frac{p(1-p)}{n}} \Leftrightarrow -1.96 \leq \frac{\hat{p} - p}{\sqrt{p(1-p)/n}} \leq 1.96$$

$1-\alpha$ Confidence Interval

$$\hat{p} \pm 3\% \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \quad 99\% \quad \hat{p} \pm 3.5\% \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

Normal: $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_n \sim N(\mu, \sigma^2)$

$$\hat{\mu} = \bar{X} \sim N(\mu, \frac{\sigma^2}{n}) \quad 95\% = \hat{\mu} \pm 1.96 s(\hat{p})$$

$$\hat{\mu} \pm 1.96 \frac{s}{\sqrt{n}}$$

2.2 BaySilm:

$$\mu \sim N(\alpha, \gamma^2) \quad [\mu | \bar{X}] \sim N\left(\frac{\frac{1}{\gamma^2} + \bar{X} \frac{n}{\sigma^2}}{\frac{1}{\gamma^2} + \frac{n}{\sigma^2}}, \frac{1}{\frac{1}{\gamma^2} + \frac{n}{\sigma^2}}\right)$$

$E[\bar{X}]$ (noninformative prior)
as $\gamma^2 \rightarrow \infty$ $N(\bar{X}, \frac{\sigma^2}{n})$

$$P(\mu \in [\bar{X} - 1.96 \frac{s}{\sqrt{n}}, \bar{X} + 1.96 \frac{s}{\sqrt{n}}]) = 95\%$$

\downarrow random \downarrow fixed interval

$$\sigma^2 \text{ unknown} \quad S^2 = \hat{\sigma}^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \quad \begin{array}{l} \text{already used one degree of freedom for} \\ \text{estimating } \mu \text{ at } \bar{x}. \end{array} \quad E[S^2] = \sigma^2$$

$$\sum_{i=1}^n (x_i - \bar{x}) = \sum x_i - n\bar{x} = \sum x_i - n \frac{\sum x_i}{n} = 0$$

For large n :

$$95\% \text{ CI.} = \bar{X} \pm 1.96 \frac{s}{\sqrt{n}} \quad , \text{ replace } \sigma^2 \text{ by } S^2$$

For small n :

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

$$T = \frac{\bar{X} - \mu}{s/\sqrt{n}} \quad \begin{array}{l} \text{(distribution spreads out)} \\ \text{added variability} \end{array}$$

$$1-\alpha \text{ CI.} : \bar{X} \pm t_{n-1, \alpha/2} \cdot \frac{s}{\sqrt{n}}$$

Hypothesis Testing

$$X_1, X_2, \dots, X_n \sim \text{Ber}(p)$$

$$H_0: p = p_0 \quad \text{simple " = "}$$

$$H_1: p > p_0 \quad \text{complex/composite " ≠ " " < " " > "}$$

Attitude: H_0 dependent

Reject H_0 only when there is "strong" evidence.

Test statistic (calculated from the data)

$$\hat{p} = (\sum_{i=1}^n X_i) / n \stackrel{H_0}{\sim} N(p_0, \frac{p_0(1-p_0)}{n}) \leftarrow \text{Null distribution}$$

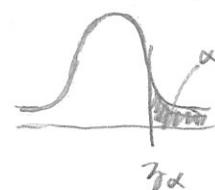
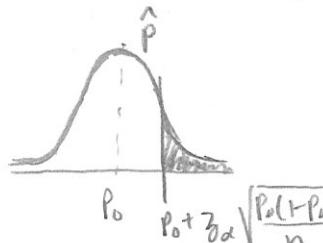
$$z(H_0) \leftarrow \text{normalization under null hypothesis} \Rightarrow z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \stackrel{H_0}{\sim} N(0, 1)$$

Compare \hat{p}_{obs} or z_{obs} to their null distribution
whether \hat{p}_{obs} or z_{obs} is too extreme to be explained by the

Decision Rule:

$$\text{Reject } H_0 \text{ if } \hat{p} > p_0 + 3\alpha \sqrt{\frac{p_0(1-p_0)}{n}}$$

$$z > z_{\alpha}$$



$$\text{Prob(Type I error)} = \alpha$$

Decision Rule \hookrightarrow Type I error.

p-value (percentile)

p-value = $P(\text{observing something more extreme than the observed value } | H_0)$

How extreme the observed value is.

$$\hat{p}_{\text{obs}} = 0.55, z = 1 \quad \text{p-value} = \int_{-\infty}^{z_{\text{obs}}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) dz$$

Decision Rule with Type I = α

Reject H_0 if p-value $< \alpha$

Example:

Admission

1. Admit if SAT score > 2300

2. Admit if SAT percentile $> 98\%$
p-value $< 2\%$

Two Sided

$$X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$$

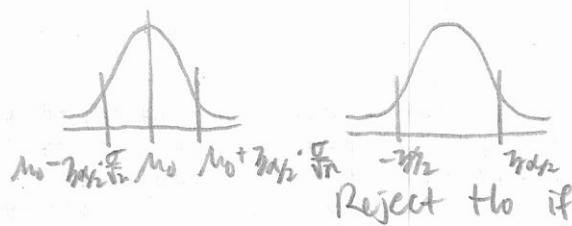
known

$$H_0: \mu = \mu_0$$

$$H_1: \mu \neq \mu_0$$

$$\bar{X} = \frac{\sum X_i}{n} \stackrel{H_0}{\sim} N(\mu_0, \frac{\sigma^2}{n})$$

$$Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} \stackrel{H_0}{\sim} N(0, 1)$$



Reject H₀ if

$$|\bar{X} - \mu_0| > 3\alpha/2 \cdot \frac{\sigma}{\sqrt{n}}$$

$$|Z| > 3\alpha/2$$

$$P(|Z| \geq |z_{\text{obs}}|)$$

Example:

$$H_0: X_1, X_2, \dots, X_{300} \sim \text{Unif}[0, 1]$$

$$H_1: \text{not from Unif}[0, 1]$$

$$\bar{X}_{\text{obs}} = 0.52$$

$$\mu = E(X) = 0.5$$

$$\sigma^2 = \text{Var}(X) = E(X^2) - E(X)^2$$

$$\frac{1}{3} - \frac{1}{4} = \frac{1}{12}$$

$$E(X^2) = \int_0^1 x^2 dx = \frac{x^3}{3} \Big|_0^1 = \frac{1}{3}$$

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \stackrel{H_0}{\sim} N(0, 1) \Rightarrow z_{\text{obs}} = \frac{\bar{X}_{\text{obs}} - 0.5}{\frac{1}{\sqrt{12}}/\sqrt{300}} = \frac{0.52 - 0.5}{\sqrt{1/12}} = 1.2 < 1.96 \quad \checkmark 5\%$$

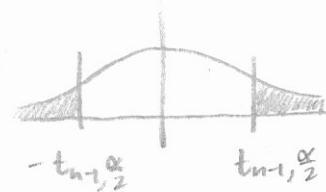
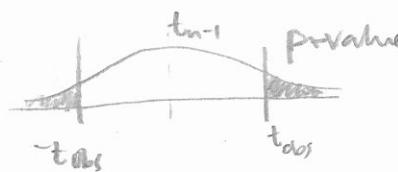
Two Sided

$$X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$$

Unknown

$$\text{original way: } Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \stackrel{H_0}{\sim} N(0, 1)$$

$$\text{new: } T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \stackrel{H_0}{\sim} t_{n-1}$$



Two Sample Test:

$$X_1, X_2, \dots, X_n \sim N(\mu_x, \sigma^2)$$

$$Y_1, Y_2, \dots, Y_m \sim N(\mu_y, \sigma^2)$$

$$H_0: \mu_x = \mu_y$$

$$H_1: \mu_x \neq \mu_y$$

$$T = \frac{\bar{X} - \bar{Y}}{S\sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{n+m-2}$$

$\bar{X} - \bar{Y} \sim N(\mu_x - \mu_y, \sigma^2(\frac{1}{n} + \frac{1}{m}))$

$$Z = \frac{\bar{X} - \bar{Y} - (\mu_x - \mu_y)}{\sigma(\sqrt{\frac{1}{n} + \frac{1}{m}})}$$

$$S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{i=1}^m (Y_i - \bar{Y})^2}{n+m-2}$$

Regression

obs	predictor X	response Y	$Y_i = \beta X_i + \varepsilon_i$ $\varepsilon_i \sim N(0, \sigma^2)$
1	X_1	Y_1	$Y_i \sim N(\beta X_i, \sigma^2)$
2	X_2	Y_2	
3	\vdots	\vdots	
i	X_i	Y_i	
n	\vdots	\vdots	

$H_0: \beta = 0$
 $H_1: \beta > 0$

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \varepsilon_i$$

$$H_0: \beta_2 = 0 \rightarrow \text{simple model } Y_i = \beta_0 + \beta_1 X_{i1} + \varepsilon_i$$

$$H_1: \beta_2 \neq 0 \rightarrow \text{complex model } Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \varepsilon_i$$

Hypothesis testing helps avoid overfitting

Likelihood Ratio Test H_0 : simple H_1 : simple $X_1, X_2, \dots, X_n \sim \text{Ber}(p)$ $H_0: p = \frac{1}{2} = 0.5$ $H_1: p = 0.7$ Reject H_0 if $\frac{L(H_1)}{L(H_0)} > c$ (cutoff)If $\frac{L(H_1)}{L(H_0)} > c$ if H_0 value < H_1 value

$$\hat{p} = \frac{\sum X_i}{n} > C_3$$

if H_0 value > H_1 value

$$\hat{p} = \frac{\sum X_i}{n} < C_4$$

Recall likelihood & maximum likelihood

$$L(p) = \prod_{i=1}^n p^{X_i} (1-p)^{1-X_i} = p^{\sum X_i} (1-p)^{n-\sum X_i}$$

MLE = max log(L(p)) over $p \in [0, 1]$

$$\Rightarrow \frac{0.7^{\sum X_i} 0.3^{n-\sum X_i}}{0.5^{\sum X_i} 0.5^{n-\sum X_i}}$$

$$\Rightarrow \left(\frac{0.7}{0.5}\right)^{\sum X_i} \left(\frac{0.5}{0.3}\right)^{n-\sum X_i} \left(\frac{0.3}{0.5}\right)^n$$

$$\Rightarrow \left(\frac{0.7}{0.5}\right)^{\sum X_i} \left(\frac{0.5}{0.3}\right)^n > c$$

After log $\rightarrow \sum X_i > C_1$ or $\frac{\sum X_i}{n} > C_2$ Form of Decision Rule $H_0: p = p_0$ $H_1: p = p_1$ ($p_1 > p_0$)Reject H_0 if $\hat{p} > c$ Type I error = α

$$\hat{p} \stackrel{H_0}{\sim} N(p_0, \frac{p_0(1-p_0)}{n})$$

$$c = p_0 + 3\sqrt{\frac{p_0(1-p_0)}{n}}$$

In General:

 $H_0: \bar{X} \sim p_0(x)$ $H_1: \bar{X} \sim p_1(x)$

Rejection Rule:

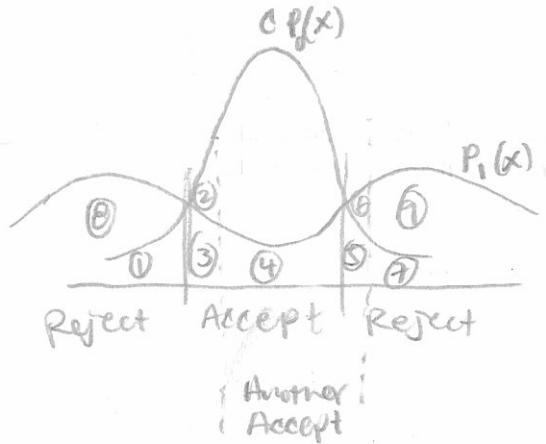
$$\frac{L(H_1)}{L(H_0)} = \frac{p_1(x)}{p_0(x)}$$

 H_1 is more plausible than H_0

Two Decision Rules:

keep type I error fixed

Compare type II error



Type I error

$$(① + ⑤ + ⑦)/c$$

Type II error:

$$③ + ④$$

Another Test:

$$I: (③ + ① + ② + ⑦)/c$$

$$II: ④ + ⑥ + ⑤$$

Same Type I

$$① + ⑤ + ⑦ = ③ + ① + ② + ⑦ \Rightarrow ⑤ = ③ + ②$$

Compare Type II:

$$\begin{aligned} ③ + ④ &\text{ vs. } ④ + ⑥ + ⑤ \\ &= ④ + ⑥ + ③ + ② \\ &\text{bigger} \end{aligned}$$

Decision: $\frac{P_1(x)}{P_0(x)} > 0 \Rightarrow \text{reject } H_0$

Bayes Test:

$$\begin{aligned} \text{prior: } p(H_1) &= \lambda \\ p(H_0) &= 1 - \lambda \end{aligned}$$

$$\text{posterior: } p(H_1|x) = \frac{\lambda P_1(x)}{\lambda P_1(x) + (1-\lambda) P_0(x)}$$

$$p(H_0|x) = \frac{(1-\lambda) P_0(x)}{\lambda P_1(x) + (1-\lambda) P_0(x)}$$

STATS 105

11/12

t-distribution

$$\bar{X}_1, \bar{X}_2, \dots, \bar{X}_n \sim N(\mu, \sigma^2)$$

unknown constant

large "n"

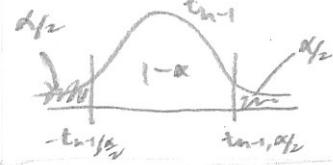
$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$(1-\alpha)CI: [\bar{X} - t_{n-1, \alpha/2} \cdot \frac{s}{\sqrt{n}}, \bar{X} + t_{n-1, \alpha/2} \cdot \frac{s}{\sqrt{n}}]$$

z transformation:

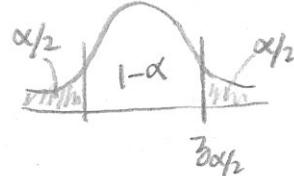
$$z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

$$T = \frac{\bar{X} - \mu}{s/\sqrt{n}} \sim t_{n-1} \xrightarrow{n \rightarrow \infty} N(0, 1)$$



1-alpha confidence interval:

$$[\bar{X} - z_{\alpha/2} \cdot \frac{s}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \cdot \frac{s}{\sqrt{n}}]$$



claim: $P(\text{mt CI}) = 1-\alpha$
 ↓
 fixed random

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

$$P(-z_{\alpha/2} \leq \frac{\bar{X}-\mu}{s/\sqrt{n}} \leq z_{\alpha/2}) = 1-\alpha \Rightarrow P\left(-z_{\alpha/2} \cdot \frac{s}{\sqrt{n}} \leq \bar{X} - \mu \leq z_{\alpha/2} \cdot \frac{s}{\sqrt{n}}\right) = 1-\alpha$$

$$P\left(\bar{X} - z_{\alpha/2} \cdot \frac{s}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \cdot \frac{s}{\sqrt{n}}\right) = 1-\alpha$$

What if σ^2 is unknown?

$$s^2 = \hat{\sigma}^2 = \frac{\sum_{i=1}^n (\bar{X}_i - \bar{X})^2}{n-1} \quad E(s^2) = \sigma^2 \Rightarrow s^2 \text{ is unbiased}$$

$$\sum_{i=1}^n (\bar{X}_i - \bar{X})^2 = \sum ((\bar{X}_i - \mu) - (\bar{X} - \mu))^2 \quad \bar{Y} = \frac{\sum Y_i}{n} = \frac{\sum (X_i - \mu)}{n} = \frac{\sum X_i}{n} - \mu$$

$$\hookrightarrow Y_i \quad \hookrightarrow \bar{Y} \quad | \quad E[Y_i^2] = E[(\bar{X}_i - \mu)^2] = \sigma^2$$

$$\begin{aligned} \textcircled{1} \quad \sum_{i=1}^n (\bar{Y}_i - \bar{Y})^2 &= \sum_{i=1}^n (Y_i^2 - 2\bar{Y}Y_i + \bar{Y}^2) \quad ; \quad E[\bar{Y}] = E[(\bar{X} - \mu)^2] = \text{Var}(\bar{X}) = \frac{\sigma^2}{n} \\ &= \sum_{i=1}^n Y_i^2 - 2\sum_{i=1}^n \bar{Y}Y_i + \sum_{i=1}^n \bar{Y}^2 = \sum Y_i^2 - 2\bar{Y}\sum Y_i + n\bar{Y}^2 = \sum Y_i^2 - 2\bar{Y} \cdot n\bar{Y} + n\bar{Y}^2 = \sum Y_i^2 - n\bar{Y}^2 \end{aligned}$$

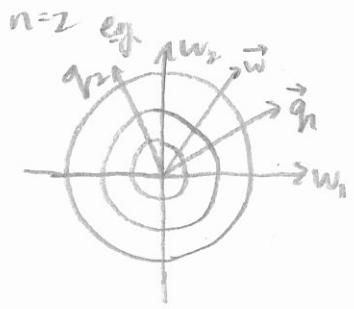
$$E[\textcircled{1}] = \sum E[Y_i^2] - nE[\bar{Y}]^2 = \sum \sigma^2 - \frac{n\sigma^2}{n} = (n-1)\sigma^2$$

$$\sum (\bar{X}_i - \bar{X})^2 \rightarrow \sum (\bar{Y}_i - \bar{Y})^2 = \sum Y_i^2 - n\bar{Y}^2 = \sum Y_i^2 - n\left(\frac{\sum Y_i}{n}\right)^2 = \sum Y_i^2 - \left(\frac{\sum Y_i}{\sqrt{n}}\right)^2$$

$$\sum \frac{(\bar{X}_i - \bar{X})^2}{\sigma^2} = \sum \left(\frac{Y_i}{\sigma}\right)^2 - \left(\sum \frac{Y_i}{\sqrt{n}}\right)^2, \quad w_i = \frac{Y_i}{\sigma} = \frac{X_i - \mu}{\sigma} \sim N(0, 1) = \sum w_i^2 - \left(\sum \frac{w_i}{\sqrt{n}}\right)^2$$

$$f(w_1, w_2, \dots, w_n) = \prod_{i=1}^n \left(\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{w_i^2}{2}\right) \right) = \frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n w_i^2\right)$$

$$\vec{w} = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix} \Rightarrow f(\vec{w}) = \frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{1}{2} \|\vec{w}\|^2\right)$$



new orthonormal system
(q_1, q_2, \dots, q_n)

$$\vec{w} \rightarrow \begin{pmatrix} v_1 = \langle \vec{w}, \vec{q}_1 \rangle \\ v_2 = \langle \vec{w}, \vec{q}_2 \rangle \\ \vdots \\ v_n = \langle \vec{w}, \vec{q}_n \rangle \end{pmatrix} \Rightarrow f(\vec{v}) = \frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{1}{2} \|\vec{v}\|^2\right)$$

$$\sum w_i^2 = \sum v_i^2 \Rightarrow \sum w_i^2 - \left(\sum \frac{w_i}{\sqrt{n}}\right)^2 \quad \text{Let } \vec{q}_1 = \begin{pmatrix} \frac{1}{\sqrt{n}} \\ \frac{1}{\sqrt{n}} \\ \vdots \\ \frac{1}{\sqrt{n}} \end{pmatrix} = \sum_{i=1}^n v_i^2 - \langle \vec{w}, \vec{q}_1 \rangle^2$$

$$= \sum_{i=1}^n v_i^2 - \bar{v}_1^2 = \sum_{i=2}^n v_i^2 \sim \chi_{n-1}^2$$

Definition: $\bar{z}_1, \bar{z}_2, \dots, \bar{z}_n \sim N(\mu, 1)$

$\sum \bar{z}_i^2 \sim \chi_n^2$ (chi squared)

$$T = \frac{\bar{x} - \mu}{S/\sqrt{n}} = \frac{(\bar{x} - \mu)/\sigma}{\frac{S/\sigma}{\sqrt{n}}} = \frac{\frac{\bar{x} - \mu}{\sigma}}{\left(\frac{\sum (\bar{x}_i - \bar{x})^2}{(n-1)(\sigma^2)}\right)^{1/2}} / \sqrt{n} = \frac{v_1/\sqrt{n}}{\sqrt{\sum_{i=2}^n v_i^2 / ((n-1) \cdot n)}} = \frac{v_1}{\sqrt{\sum_{i=2}^n v_i^2 / (n-1)}} \sim \frac{N(0, 1)}{\sqrt{\frac{\chi_{n-1}^2 / (n-1)}{n-1}}}$$

Two Sample Problem:

$$\bar{X}_1, \bar{X}_2, \dots, \bar{X}_m \sim N(\mu_X, \sigma^2) \quad \Delta = \bar{X}_1 - \bar{X}_2 \quad \bar{X} \sim N(\mu_X, \frac{\sigma^2}{m}), \bar{Y} \sim N(\mu_Y, \frac{\sigma^2}{m})$$

$$\bar{Y}_1, \bar{Y}_2, \dots, \bar{Y}_n \sim N(\mu_Y, \sigma^2) \quad \hat{\Delta} = \bar{X} - \bar{Y} \sim N(\mu_X - \mu_Y, \frac{\sigma^2}{n} + \frac{\sigma^2}{m})$$

$$\sim N(\Delta, \sigma^2(\frac{1}{n} + \frac{1}{m}))$$

$$Z = \frac{\bar{X} - \bar{Y}}{\sigma \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim N(0, 1)$$

$$1-\alpha \text{ CI for } \Delta : \hat{\Delta} \pm z_{\alpha/2} \sigma \sqrt{\frac{1}{n} + \frac{1}{m}}$$

$$T = \frac{\bar{X} - \bar{Y}}{S \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{n+m-2}$$

$$1-\alpha \text{ CI for } \Delta : \hat{\Delta} \pm t_{n+m-2, \alpha/2} S \sqrt{\frac{1}{n} + \frac{1}{m}}$$

$$S^2 = \frac{\sum (\bar{X}_i - \bar{X})^2 + \sum (\bar{Y}_i - \bar{Y})^2}{n+m-2}$$

Hypothesis Testing

Want to test whether a coin is fair.
Flip n times Observed 55 heads

Null hypothesis: $H_0: p = \frac{1}{2}$ (p is probability of head)

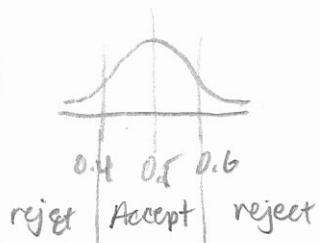
Alternative hypothesis: $H_1: p \neq \frac{1}{2}$

Collect Data: $X_1, X_2, \dots, X_n \sim \text{Ber}(p)$ $\hat{p} = \frac{\sum X_i}{n} \Rightarrow \hat{p}_{\text{observed}} = \frac{55}{100} = 0.55$
 $\sim N(p, \frac{p(1-p)}{n})$

Norm Distribution, reference distribution

$$\hat{p} \stackrel{H_0}{\sim} N(p = \frac{1}{2}, \frac{\frac{1}{2}(1-\frac{1}{2})}{100}) \sim N(\frac{1}{2}, \frac{1}{400}) \sim N(\frac{1}{2}, (0.05)^2)$$

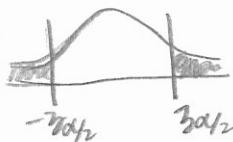
$$z(H_0): z = \frac{\hat{p} - \frac{1}{2}}{\sqrt{\frac{\frac{1}{2}(1-\frac{1}{2})}{100}}} = \frac{\hat{p} - 0.5}{0.05} \stackrel{H_0}{\sim} N(0, 1)$$



$$z_{\text{obs}} = \frac{0.55 - 0.5}{0.05} = 1 \quad \text{Compare } z_{\text{obs}} \text{ w/ } N(0, 1)$$

Type 1 Error: H_0 is true, but reject it.

$$\text{Prob(Type 1 error)} = \alpha \quad (\text{level of significance})$$



Type 2 Error: accept H_0 when H_0 is false
 H_1 is true

Prob(Type 2 error) Power of the test

Election Example:

Survey n people
 $X_i = \begin{cases} 1 & \text{vote for} \\ 0 & \text{not} \end{cases}$

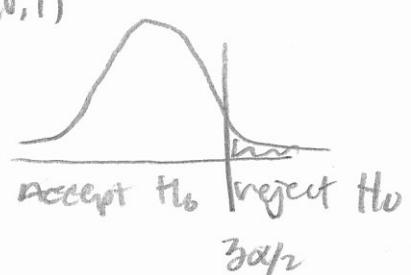
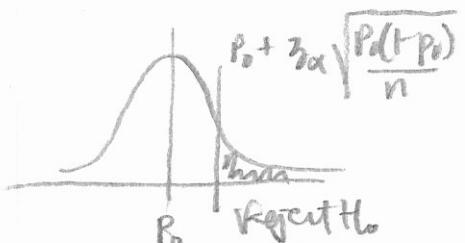
$X_1, X_2, \dots, X_n \sim \text{Ber}(p)$

$H_0: p = \frac{1}{2} = p_0$

$H_1: p > \frac{1}{2} = p_1$ (one-sided H_1)
 $(p \neq \frac{1}{2} \text{ two sided})$

$$\hat{p} \sim N\left(\frac{1}{2}, \frac{\frac{1}{2}(1-\frac{1}{2})}{n}\right) \sim N(p_0, \frac{p_0(1-p_0)}{n})$$

$$z(H_0): z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \stackrel{H_0}{\sim} N(0, 1)$$

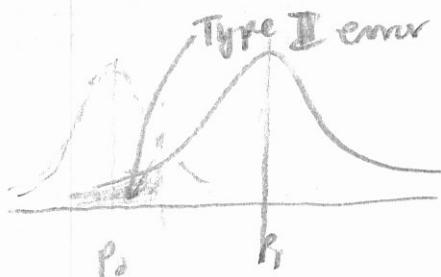


Type II error Accept H_0 when H_1 is true.

$$P(\text{Type II error}) > ?$$

$$H_0: p = p_0 \quad (0.5) \quad \hat{p} \stackrel{H_0}{\sim} N(p_0, \frac{p_0(1-p_0)}{n})$$

$$H_1: p = p_1 \quad (0.7) \quad \hat{p} \stackrel{H_1}{\sim} N(p_1, \frac{p_1(1-p_1)}{n})$$



STATS WS

11/26

1 sided, hand written, half sheet (NO Project)

Recall Bayes rule

$$\text{population: } p(\text{male}) = \lambda \quad [\underline{x} | z=1] \sim f_1(x)$$

$$p(\text{female}) = 1-\lambda \quad [\underline{x} | z=0] \sim f_0(x)$$

$$p(z=1 | \underline{x}=x) = \frac{\lambda f_1(x)}{\lambda f_1(x) + (1-\lambda) f_0(x)} = \frac{p(z=1) f(x | z=1)}{p(z=1) f(x | z=1) + p(z=0) f(x | z=0)} = \frac{p(1, x)}{p(1, x) + p(0, x)} = \frac{p(1, x)}{p(x)}$$

Count the population:

M = total # of people

$$\# \text{ of males} = M\lambda \rightarrow \# \text{ of males in } (x, x+\Delta x) = M\lambda f_1(x) \Delta x$$

$$\# \text{ of females in } (x, x+\Delta x) = M(1-\lambda) f_0(x) \Delta x$$

$$\# \text{ of people in } (x, x+\Delta x) = M(\lambda f_1(x) + (1-\lambda) f_0(x)) \Delta x$$

Among the slice of population in $(x, x+\Delta x)$, $= M f(x) \Delta x$

What is the proportion of males?

density

↳ density of entire population.

$$\frac{M\lambda f_1(x) \Delta x}{M(\lambda f_1(x) + (1-\lambda) f_0(x)) \Delta x} = p(\text{male} | \underline{x} \in (x, x+\Delta x))$$

↓
0
 $\underline{x} = x$

Bayes Test

$$H_0: \underline{x} \sim P_0(x)$$

$$H_1: \underline{x} \sim P_1(x)$$

$$\text{Prior prob: } p(H_0) = \lambda$$

$$\text{evidence: } [\underline{x} | H_0] \sim P_0(x)$$

$$[\underline{x} | H_1] \sim P_1(x)$$

$$\text{Posterior: } p(H_1 | \underline{x}=x) = \frac{\lambda P_1(x)}{\lambda P_1(x) + (1-\lambda) P_0(x)}$$

Bayes Decision Rule:

		Truth	
		H_0	H_1
Decision	Accept H_0	0	Loss II
	Reject H_0	Loss I	0

Risk: expected loss

① Risk of accepting $H_0 = P(H_0 | x) \cdot 0 + P(H_1 | x) \cdot \text{loss II}$ ② Risk of rejecting $H_0 = P(H_0 | x) \cdot \text{loss I} + P(H_1 | x) \cdot 0$ Reject H_0 if ② < ① or ① > ② or $\frac{①}{②} > 1$

$$\frac{\lambda P_1(x)}{\lambda P_1(x) + (1-\lambda) P_0(x)} \cdot \text{loss II} > 1 \Rightarrow \frac{P_1(x)}{P_0(x)} > \frac{(1-\lambda) \text{loss I}}{\lambda \text{loss II}}$$

likelihood ratio ↓

Reject H_0

$$\frac{P_1(x)}{P_0(x)} = \frac{\text{likelihood}(H_1)}{\text{likelihood}(H_0)} > \frac{\text{prior } Pr(H_0)}{\text{prior } Pr(H_1)} \cdot \frac{\text{loss of Type I}}{\text{loss of Type II}}$$

Generalize Likelihood Ratio

test whether die is fair.

$$H_0: p_1 = \frac{1}{6}, p_2 = \frac{1}{6}, \dots, p_6 = \frac{1}{6}$$

$$H_1: (p_1, p_2, \dots, p_6) \neq \left(\frac{1}{6}, \frac{1}{6}, \dots, \frac{1}{6}\right)$$

$$\text{Likelihood}(p_1, p_2, \dots, p_6) = p_1^{n_1} p_2^{n_2} \cdots p_6^{n_6}, \quad n_1 = \# \text{ of } 1s, \quad n_2 = \# \text{ of } 2s$$

$$\text{Suppose: } H_0: (p_1, p_2, p_3, \dots, p_6) = \left(\frac{1}{6}, \frac{1}{6}, \dots, \frac{1}{6}\right) \quad + n_6 = \# \text{ of } 6s.$$

$$H_1: (p_1, p_2, \dots, p_6) = (0.1, 0.1, 0.2, 0.2, 0.3, 0.1) \quad n = \text{total # of repetitions}$$

$$\text{Likelihood: } \frac{L(H_1)}{L(H_0)} = \frac{0.1^{n_1+n_2+n_6} 0.2^{n_3+n_4} 0.3^{n_5}}{\left(\frac{1}{6}\right)^{n_1+n_2+\dots+n_6}}$$

What if $H_1: (p_1, p_2, \dots, p_6) \neq \left(\frac{1}{6}, \frac{1}{6}, \dots, \frac{1}{6}\right)$

$$\frac{L(H_1)}{L(H_0)} = \frac{\hat{p}_1^{n_1} \hat{p}_2^{n_2} \cdots \hat{p}_6^{n_6}}{\left(\frac{1}{6}\right)^{n_1+n_2+\dots+n_6}} \quad \hat{p}_1 = \frac{n_1}{n}, \quad \hat{p}_2 = \frac{n_2}{n}, \dots, \hat{p}_6 = \frac{n_6}{n}$$

$$2 \log \frac{L(H_1)}{L(H_0)} \sim \chi^2_5 \quad (\hat{p}_1, \hat{p}_2, \dots, \hat{p}_6) = \left(\frac{n_1}{n}, \frac{n_2}{n}, \dots, \frac{n_6}{n}\right) \text{ MLE}$$

$$L(p_1, p_2, \dots, p_6) = p_1^{n_1} p_2^{n_2} \cdots p_6^{n_6}$$

$$l(p_1, \dots, p_6) = n_1 \log p_1 + \dots + n_6 \log p_6 = n_1 \log p_1 + \dots + n_5 \log p_5 + n_6 \log(1 - (p_1 + \dots + p_5))$$

$$\begin{aligned} \frac{\partial l}{\partial p_1} &= \frac{n_1}{p_1} - \frac{n_6}{1 - (p_1 + \dots + p_5)} \in p_6 = 0 \\ \frac{\partial l}{\partial p_2} &= \frac{n_2}{p_2} - \frac{n_6}{1 - (p_1 + \dots + p_5)} \in p_6 = 0 \end{aligned} \quad \left. \begin{array}{l} \frac{n_1}{p_1} = \frac{n_2}{p_2} = \dots = \frac{n_6}{p_6} = C \\ \frac{n_1}{C} = \frac{n_2}{C} = \dots = \frac{n_6}{C} = 1 \end{array} \right\} \quad \frac{n_1}{p_1} = \frac{n_2}{p_2} = \dots = \frac{n_6}{p_6} = C$$

$$\begin{cases} \hat{p}_1 = \frac{n_1}{C} \\ \hat{p}_2 = \frac{n_2}{C} \\ \vdots \\ \hat{p}_6 = \frac{n_6}{C} \end{cases} \quad \hat{p}_1 + \hat{p}_2 + \dots + \hat{p}_6 = 1 \quad \frac{n_1}{C} + \frac{n_2}{C} + \dots + \frac{n_6}{C} = 1 \quad C = n$$

$$\frac{\hat{p}_1^{n_1} \hat{p}_2^{n_2} \cdots \hat{p}_6^{n_6}}{\left(\frac{1}{6}\right)^{n_1+n_2+\dots+n_6}} = \frac{\max_{H_1} p_1^{n_1} p_2^{n_2} \cdots p_6^{n_6}}{\max_{H_0} p_1^{n_1} p_2^{n_2} \cdots p_6^{n_6}}$$

max out unknown in H₁

max out unknown in H₀

2. log like ratio $\stackrel{H_0}{\sim} \chi^2$ # of free parameters in H₁ - # of free parameters in H₀

5 0



overfitting
high cutoff value to avoid overfitting

STATS 105

12/03

Final: Similar to homework/lecture (6 Questions)

Regression

obs	predictor	Response	Simplest Regression
1	x_1	y_1	$y_i \approx \beta x_i$
2	x_2	y_2	more precisely
\vdots	\vdots	\vdots	
n	x_n	y_n	$y_i \stackrel{iid}{\sim} N(\beta x_i, \sigma^2)$
	training		$y_i = \beta x_i + \epsilon_i \quad \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2) \quad (1)$

(given)

Maximum Likelihood: Assume x_i are fixed & y_i random

$$L(\beta) = \prod_{i=1}^n f(y_i | x_i, \beta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \beta x_i)^2}{2\sigma^2}\right) = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta x_i)^2\right)$$

maximum likelihood \Rightarrow least squares

Find $\hat{\beta}$ by minimizing $\sum_{i=1}^n \underbrace{(y_i - \beta x_i)^2}_{\text{Residual}} = R(\beta)$



$$R(\beta) = \sum_{i=1}^n 2(y_i - \beta x_i)(-x_i) = 0 \Rightarrow \hat{\beta} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$$

Scatter plot:

Vector:
$$R(\beta) = |\vec{Y} - \vec{\beta} \vec{X}|^2$$

$$\langle \vec{X}, \vec{Y} - \vec{\beta} \vec{X} \rangle = 0 \Rightarrow \sum_{i=1}^n x_i (y_i - \beta x_i) = 0$$

Under hypothetical Repetition, y_i random, β random (different from repetition to repetition)

$$\beta = \frac{\sum x_i y_i}{\sum x_i^2} = \frac{\sum x_i (\beta_{\text{TRUE}} x_i + \epsilon_i)}{\sum x_i^2} = \frac{\sum x_i^2 \beta_{\text{TRUE}} + \sum x_i \epsilon_i}{\sum x_i^2} = \beta_{\text{TRUE}} + \frac{\sum x_i \epsilon_i}{\sum x_i^2}$$

↓ signal ↓ noise

$$\begin{aligned} E[\hat{\beta}] &= \beta_{\text{TRUE}} + \frac{\sum x_i E[\epsilon_i]}{\sum x_i^2} = \beta_{\text{TRUE}} \quad (\text{unbiased}) \\ \text{Var}[\hat{\beta}] &= \frac{\sum x_i^2 \text{Var}(\epsilon_i)}{(\sum x_i^2)^2} = \frac{\sigma^2}{\sum x_i^2} \end{aligned} \quad \left. \begin{aligned} \hat{\beta} &\sim N(\beta_{\text{TRUE}}, \frac{\sigma^2}{\sum x_i^2}) \\ \text{Sampling Distribution} \end{aligned} \right.$$

95% CI: $\hat{\beta} \pm z_{\alpha/2} \left(\frac{\sigma}{\sqrt{\sum x_i^2}} \right) \quad (\sigma^2 \text{ known})$
 (-1.96)

T-distribution (σ^2 unknown) by s^2

Testing: $H_0: \beta_{\text{TRUE}} = 0$, $H_1: \beta_{\text{TRUE}} \neq 0$.

$$\text{Under } H_0: \hat{\beta} \sim N(0, \frac{\sigma^2}{\sum x_i^2}) \Rightarrow z = \frac{\hat{\beta}}{\sqrt{\frac{\sigma^2}{\sum x_i^2}}} \sim N(0, 1)$$



($H_1: \beta > 0$)

decision rule: reject $z > z_\alpha$
Type I: α

Only if $\hat{\beta}$ (or z) large enough we believe $\beta_{\text{TRUE}} \neq 0$.
Avoid overfitting.

Simple Regression:

$$y_i \approx \alpha + \beta x_i \quad y_i = \alpha + \beta x_i + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

$$y_i \stackrel{iid}{\sim} N(\underbrace{\alpha + \beta x_i}_{\text{TRUE TRUE}}, \sigma^2)$$

Least Square: $R(\alpha, \beta) = \sum (y_i - (\alpha + x_i \beta))^2$

pretend β is known: $\hat{\alpha} \leftarrow \min \sum (y_i - \beta x_i - \bar{x}_i \cdot 1)^2 \Rightarrow \hat{\alpha} = \frac{\sum (y_i - \beta x_i) \times 1}{\sum 1^2 (=n)}$

$$R(\alpha, \beta) = R(\beta) = \sum (y_i - (\bar{y} - \bar{\beta} \bar{x} + \beta x_i))^2$$

$$= \sum ((y_i - \bar{y}) - \beta(x_i - \bar{x}))^2 = \sum (\tilde{y}_i - \beta \tilde{x}_i)^2$$

$$\hat{\beta} = \frac{\sum \tilde{x}_i \tilde{y}_i}{\sum \tilde{x}_i^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{n} \sum (x_i - \bar{x})^2} = \frac{\text{Sample Covariance}}{\text{Sample Var of } x}$$

$$\hat{\alpha} = \bar{y} - \bar{\beta} \bar{x} \Rightarrow \bar{y} = \hat{\alpha} + \hat{\beta} \bar{x} \quad (\bar{x}, \bar{y}) \text{ on regression line}$$

$$\sum (x_i - \bar{x}) = \sum x_i - n \bar{x} = 0$$

$$\hat{\beta} = \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x}) x_i} = \frac{\sum \tilde{x}_i (\alpha + \beta x_i + \epsilon_i)}{\sum \tilde{x}_i x_i} = \frac{\alpha \sum \tilde{x}_i + \beta \sum \tilde{x}_i x_i + \sum \tilde{x}_i \epsilon_i}{\sum \tilde{x}_i x_i} = \beta_T + \frac{\sum \tilde{x}_i \epsilon_i}{\sum \tilde{x}_i^2}$$

$$E[\hat{\beta}] = \beta_{\text{TRUE}} \quad \text{Var}[\hat{\beta}] = \frac{\sigma^2}{\sum \tilde{x}_i^2} \Rightarrow \hat{\beta} \sim \left(\beta_{\text{TRUE}}, \frac{\sigma^2}{\sum \tilde{x}_i^2} \right) \quad \begin{matrix} \text{sample} \\ \text{distribution} \end{matrix}$$

Multiple Regression:

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i$$

weight = intercept + β_1 height + β_2 height² + β_3 gender + ... + error

$$\beta_1 \quad 1 \quad x_{i2} \quad x_{i3} \quad x_{i4}$$

useful for predictions

linear in β 's. Can be nonlinear in variables

$$R(\beta) = \sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij})^2 = \sum_{i=1}^n (y_i - \vec{x}_i^T \vec{\beta})^2$$

$$\vec{x}_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix} \quad \vec{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix}$$

$$\frac{\partial R(\beta)}{\partial \beta_j} = \sum_{i=1}^n 2(y_i - \sum_{j=1}^p \beta_j x_{ij})(-x_{ij}) \Rightarrow \frac{\partial R}{\partial \beta} = \begin{pmatrix} \frac{\partial R}{\partial \beta_1} \\ \frac{\partial R}{\partial \beta_2} \\ \vdots \\ \frac{\partial R}{\partial \beta_p} \end{pmatrix} = -2 \sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij}) \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix}$$

$$\sum \vec{x}_i (y_i - \vec{x}_i^T \vec{\beta}) = 0$$

$$= \sum \vec{x}_i y_i - \sum \vec{x}_i \vec{x}_i^T \vec{\beta} = 0$$

$$= -2 \sum (y_i - \vec{x}_i^T \vec{\beta}) \vec{x}_i = 0$$

$$\sum \vec{x}_i \vec{x}_i^T \vec{\beta} = \sum \vec{x}_i y_i \Rightarrow \hat{\beta} = \left(\sum \vec{x}_i \vec{x}_i^T \right)^{-1} \sum \vec{x}_i y_i$$

Simple Regression:

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2}$$

$$\begin{array}{c} \text{Diagram showing } y \text{ vs } x_1 \text{ and } x_2. \\ \text{The regression line is } \hat{y} = \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2. \\ \text{The residual vector is } \langle x_j, y_i - \sum \hat{\beta}_j x_j \rangle = 0. \end{array}$$

Regression

obs	1	predictors	2	3	...	p	Response
1		x_{11}	x_{12}	\vdots	\vdots	x_{1p}	y_1
2		x_{21}	x_{22}	\vdots	\vdots	x_{2p}	y_2
3		x_{31}	x_{32}	\vdots	\vdots	x_{3p}	\vdots
n		x_{n1}	x_{n2}	\vdots	\vdots	x_{np}	y_n

Linear Regression

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i$$

linear in β 's.

May not be linear in original variables

Least Square Method:

find $(\beta_1, \beta_2, \dots, \beta_p)$ by minimizing $\sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij})^2 \Rightarrow (\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_p)$

Vector Notation

$$\beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix} \quad x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix}$$

$$y_i = x_i^\top \beta + \epsilon_i$$

$$R(\beta) = \sum_{i=1}^n (y_i - x_i^\top \beta)^2$$

$$\hat{\beta} = (\sum_{i=1}^n x_i x_i^\top)^{-1} \sum_{i=1}^n (x_i y_i)$$

Statisticians: understanding observed data

$$H_0: \beta_j = 0 \quad \& \quad H_1: \beta_j > 0 \quad (\text{to avoid overfitting})$$

Machine Learners: predicting testing data

Regularization (to avoid overfitting)

$$\min \left(\sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \text{penalty}(\beta_1, \dots, \beta_p) \right)$$

$$\text{penalty}(\beta_1, \dots, \beta_p) = \sum_{j=1}^p \beta_j^2$$

$$= \sum_{j=1}^p |\beta_j|$$

Lasso (Least Absolute shrinkage & selection operator)

Recommender System

Users	items	1	2	...	i	...	p
1							
2							
...							
n							

r_{ui}

Aspects/factors/Dimensions:
dimensions

item i : $\vec{a}_i = \begin{pmatrix} a_{i1} \\ a_{i2} \\ \vdots \\ a_{ik} \end{pmatrix}$

a_{ik} = appeal of item i in the k^{th} aspect

user u : $\vec{p}_u = \begin{pmatrix} p_{u1} \\ p_{u2} \\ \vdots \\ p_{uk} \end{pmatrix}$

p_{uk} = preference of user u in k^{th} aspect

$$r_{ui} = \sum_{k=1}^K p_{uk} \cdot a_{ik} + \epsilon_{ui}$$

$$r_{ui} = \mu + b_u + b_i + \sum_{k=1}^K p_{uk} a_{ik} + \epsilon_{ui}$$

$$\min_{\substack{\text{obs} \\ \text{ratings}}} \sum (r_{ui} - (\mu + b_u + b_i + \vec{p}_u^\top \vec{a}_i))^2 + \lambda_1 |\vec{p}_u|^2 + \lambda_2 |\vec{a}_i|^2$$

$$\Rightarrow \hat{\mu}, \hat{b}_u, \hat{b}_i, \hat{p}_u, \hat{a}_i$$

Simple Version:

$$\sum (r_{ui} - \vec{p}_u^\top \vec{a}_i)^2 \quad \text{Alternating: Given } \vec{p}_u, \text{ find } \vec{a}_i \\ \text{Given } \vec{a}_i, \text{ find } \vec{p}_u$$

Logistic Regression

$y_i \in \{0, 1\}$ - statistics

$\{-1, +1\}$ - machine learning

Assumption: $y_i \sim \text{Ber}(p_i)$; $\log \frac{p_i}{1-p_i} = \sum_{j=1}^p \beta_j x_{ij} = \vec{x}_i^\top \vec{\beta} \Leftrightarrow p_i = \frac{\exp(\vec{x}_i^\top \vec{\beta})}{1+\exp(\vec{x}_i^\top \vec{\beta})}$

maximum likelihood:

$$L(\vec{\beta}) = \prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i} = \prod_{i=1}^n \left(\frac{\exp(\vec{x}_i^\top \vec{\beta})}{1+\exp(\vec{x}_i^\top \vec{\beta})} \right)^{y_i} \left(\frac{1}{1+\exp(\vec{x}_i^\top \vec{\beta})} \right)^{1-y_i}$$

$$= \prod_{i=1}^n \frac{\exp(y_i \vec{x}_i^\top \vec{\beta})}{1+\exp(\vec{x}_i^\top \vec{\beta})} \Rightarrow l(\vec{\beta}) = \sum_{i=1}^n y_i \vec{x}_i^\top \vec{\beta} - \log(1+\exp(\vec{x}_i^\top \vec{\beta}))$$

$$l'(\vec{\beta}) = \left(\frac{\partial l}{\partial \beta_1}, \dots, \frac{\partial l}{\partial \beta_p} \right) = \sum y_i \vec{x}_i - \frac{\exp(\vec{x}_i^\top \vec{\beta})}{1+\exp(\vec{x}_i^\top \vec{\beta})} \cdot \vec{x}_i = \sum_{i=1}^n \vec{x}_i (y_i - p_i)$$

Learning Algorithm:

gradient ascent

Start from $\beta^{(t)}$

$$\text{iterate: } \beta^{(tn)} = \beta^{(t)} + \gamma \sum_{i=1}^n x_i (y_i - p_i^{(t)})$$

learning rate previous page

- adaboost (form of logistic regression)

X_i : weak classifiers

↓ boost

Strong classifiers

- support vector machine (maximum margin principle)

--- want a large Margin (M)