

Important limiting results:

Central Limit Theorem

If X_1, X_2, \dots, X_n is an iid sample with ANY distribution w/ mean μ and standard deviation σ , then:

$$\bar{X} \xrightarrow{D} N(\mu, \frac{\sigma^2}{n}), \text{ as } n \rightarrow \infty$$

(as the sample size increases, the sample mean behaves as if it is from a Normal distribution with mean μ and standard deviation (σ/\sqrt{n}))

Maximum Likelihood Estimate (MLE)

Likelihood function, $l_{\text{lik}}(\theta)$, corresponds to the probability of observing the particular data in our sample for various values θ , assuming the data came from some distribution w/ parameter θ .

Our goal is to find the value of θ that maximizes the likelihood function.

The Likelihood Function

$$l_{\text{lik}}(\theta) = P(X_1=x_1, X_2=x_2, \dots, X_n=x_n) = P(X_1=x_1)P(X_2=x_2)\dots P(X_n=x_n)$$

- * cause X_i are iid
- $= \prod_{i=1}^n f_\theta(x_i)$
- * since each X_i have density f_θ

Computing MLE

our goal is to compute the formula of the parameter θ that maximizes the likelihood function.

- Find likelihood $l_{\text{lik}}(\theta) = \prod_{i=1}^n f_\theta(x_i)$
- Calculate the log of the likelihood function (aka log likelihood)

$$l(\theta) = \log(l_{\text{lik}}(\theta)) = \sum \log(f_\theta(x_i))$$

- Differentiate in respect to θ , set to zero, and solve for θ .

Consistency:

An estimate of $\hat{\theta}_n$ of θ is consistent if $\hat{\theta}_n \xrightarrow{P} \theta$ as $n \rightarrow \infty$

Consistency of MLE:

Theorem: the MLE $\hat{\theta}_{\text{MLE},n}$ is a consistent estimator of the parameter, θ , that it is estimating, which means that:

$$\hat{\theta}_{\text{MLE},n} \xrightarrow{P} \theta \text{ as } n \rightarrow \infty$$

"actual value of the parameter estimate approaches the true value of the parameter"

Asymptotic Unbiasedness:

$$E[\hat{\theta}_{\text{MLE},n}] \rightarrow \theta \text{ as } n \rightarrow \infty$$

"the expected value of the parameter estimate approaches the true value of the parameter"

Asymptotic Normality of the MLE:

the MLE is asymptotically normal. If $\hat{\theta}_{\text{MLE},n}$ is the MLE estimate of a parameter θ whose true value is θ_0 , then as $n \rightarrow \infty$, we have that: $\hat{\theta}_{\text{MLE},n} \xrightarrow{D} N(\theta_0, \frac{1}{n I(\theta_0)})$ as $n \rightarrow \infty$

$I(\theta_0)$ is the Fisher information

$$\hookrightarrow -E \left[\frac{\partial^2}{\partial \theta^2} \log(f_\theta(x)) \right]_{\theta=\theta_0}$$

$I(\theta_0)$ measures how "peaked" $l(\theta)$ is around θ_0 . If this value is large, then it is easier to detect θ_0 . ∵ lower variance and vice versa (and should result the same distribution as CLT)

Cramer Rao Lower Bound

Theorem: if X_i are iid from a distribution with density f_θ , under smoothness conditions on f_θ , we have that:

If $\hat{\theta}$ is an unbiased estimator for θ , then: $\text{var}(\hat{\theta}) \geq \frac{1}{n I(\theta)}$ aka variance of MLE

This means the MLE has the LOWEST possible variance among unbiased estimators.

Efficiency:

Given two unbiased estimators $\hat{\theta}$ and $\tilde{\theta}$ of a parameter θ , the efficiency of $\hat{\theta}$ relative to $\tilde{\theta}$ is:

$$\text{eff}(\hat{\theta}, \tilde{\theta}) = \frac{\text{var}(\hat{\theta})}{\text{var}(\tilde{\theta})}$$

If $\text{eff}(\hat{\theta}, \tilde{\theta}) \leq 1$, then $\text{var}(\hat{\theta}) \geq \text{var}(\tilde{\theta})$, which implies that $\hat{\theta}$ is less efficient than $\tilde{\theta}$. Vice versa applies.

Efficient Estimators:

An unbiased estimator that achieves the Cramer Rao lower bound is called efficient. $\text{var}(\hat{\theta}) = \frac{1}{n I(\theta)}$ better than this

MLE is asymptotically efficient, but not really for finite samples.

Confidence Intervals:

95% confidence intervals computed from iid random samples contain the true population parameter 95% of the time.

Making CI

- if it follows CLT or MLE results:

$$P(\hat{\theta} - z_{\alpha/2} \sigma_{\hat{\theta}} \leq \hat{\theta} \leq \hat{\theta} + z_{\alpha/2} \sigma_{\hat{\theta}}) = 1 - \alpha$$

where α is our level 5% or 1% are most common.

Confidence Interval of MLE

$$\left[\hat{\theta}_{\text{MLE}} \pm \frac{z_{\alpha/2}}{\sqrt{n} I(\theta_0)} \right]$$

CI w/ unknown pop variance

If the X_i are iid and normally distributed with unknown population variance σ^2 , then an unbiased estimate of σ^2 is:

$$\sigma^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2 \quad (\text{sample variance})$$

And it turns out that:

$$\frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \sim t_{n-1} \quad t \text{ distribution with } n-1 \text{ degrees of freedom}$$

→ even if the data isn't normal, for large sample sizes, it still holds.

thus the CI is as follows:

$$\left[\bar{X} \pm t_{n-1, \alpha/2} \frac{\sigma}{\sqrt{n}} \right] \quad \text{where } t_{n-1, \alpha/2} \text{ is the value such that } P(T \leq t_{n-1, \alpha/2}) = 1 - \alpha$$

Hypothesis Test

Set up:

H_0 : accepted hypothesis

H_1 : alternate hypothesis

p-val: value to test at our significance level

null alternate p-value

$$H_0: \mu = \mu_0 \quad H_1: \mu > \mu_0 \quad P(Z \geq z|H_0) = 1 - \Phi(z)$$

$$H_0: \mu = \mu_0 \quad H_1: \mu < \mu_0 \quad P(Z \leq z|H_0) = \Phi(z)$$

$$H_0: \mu \neq \mu_0 \quad H_1: \mu \neq \mu_0 \quad P(|Z| \geq z|H_0) = 2(1 - \Phi(z))$$

One-Sample Z tests w/ known Variance

$$Z = \frac{\text{estimated val} - \text{null value}}{\text{SD of estimate}} = \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}} \sim N(0,1)$$

if the Z value looks unlikely to have come from a $N(0,1)$ distribution, that is evidence against the null hypothesis.

↳ due to known properties of the normal curve, we can find probability of the Z value happening as per the conditions of the alternate hypothesis.

↳ The p-value is not the probability that the null is false or is it the probability the null is true.

↳ just tells us the probability of the test statistic happening in the null distribution.

CLT and how it relates to testing:

If 95% CI for μ contains μ_0 , then we would not reject H_0 , at the $\alpha = 0.05$ level.

One Sample T-test w/ unknown variance

- if we don't have σ , we can use the sample SD w/ t-test instead!

- if the X_i are iid $N(\mu, \sigma^2)$, then:

$$T = \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}} \sim t_{n-1} \quad \text{assumption data is normal but tends to hold at } n > 30$$

Two Sample Z-test (var known)

X_1, \dots, X_n IID from pop w/ unknown μ_1 and known σ_1^2

Y_1, \dots, Y_m IID from pop w/ unknown μ_2 and known σ_2^2

$$H_0: \mu_1 = \mu_2 \rightarrow \text{gives us } \mu_1 - \mu_2 = 0$$

$$Z = \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m}}} \sim N(0,1)$$

Two Sample T-test (unknown and unequal σ^2)

AHA! Welch's test

same X and Y from above, but with unknown variance

$$T = \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m}}} \sim t_{df} \quad df = \frac{\left(\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m} \right)^2}{\frac{\sigma_1^4}{n^2} + \frac{\sigma_2^4}{m^2}}$$

Two-Sample t-test: var unknown but equal

AKA pooled variance

same conditions as above but $\sigma_1^2 = \sigma_2^2$

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

and we use pooled variance

$$\hat{\sigma}^2 = \frac{(n-1)\sigma_1^2 + (m-1)\sigma_2^2}{n+m-2} = \frac{\sum (X_i - \bar{X})^2 + \sum (Y_j - \bar{Y})^2}{n+m-2}$$

Paired Two Sample test (Z)

A paired test can be reduced to a single sample test that the mean difference is 0!

Define $D_i = X_i - Y_i$

$$H_0: \mu_D = 0 \quad Z = \frac{\bar{D}}{\sigma_D / \sqrt{n}}$$

$$H_1: \mu_D \neq 0 \quad \therefore \text{p-val} = 2(1 - \Phi(|Z|))$$

Paired Two Sample test (t)

$$T = \frac{\bar{D}}{\sigma_D / \sqrt{n}} \sim t_{n-1} \quad \text{but replace w/ tstat not } \Phi$$

Two Sample test for Proportions

$X_1, \dots, X_n \sim \text{Bernoulli}(p_1)$

$Y_1, \dots, Y_m \sim \text{Bernoulli}(p_2)$

estimate p from $H_0: p_1 = p_2$, $H_1: p_1 \neq p_2$ the sample

$$Z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1-\hat{p})} \left(\frac{1}{n} + \frac{1}{m} \right)} \sim N(0,1)$$

where \hat{p} is the pooled estimate of variance

$$\hat{p} = \frac{\sum X_i + \sum Y_i}{n+m}$$

* can use this test to test for independence between the data

All test above are PARAMETRIC tests

Mann Whitney Test

Non-normal data assumption

X_1, \dots, X_n are IID with unknown dist. F

Y_1, \dots, Y_m are IID with unknown dist. G

→ checks to see if there is a difference in the ranks of the TWO samples

$H_0: F = G \quad H_1: F \neq G$

First step is to rank all data, w/ 1 the lowest, increasing from there. Ties are split by giving .5 on the rank it would have been.

Second is to sum the ranks from each sample.

U-test Statistic:

computes the amount of overlap in the ranks of each sample

Smaller U means bigger difference

Smaller U test stat is MORE Significant

Using the R that is lower, you find the # of values smaller in the other sample, sample by sample, then sum. The ranks that were split are given .5.

ex.

A	6.25	Rank	R _A
A	5.50	4	R _B = 3+2 = 5
A	4.75	1	U = 1 + 1 = 2
B	5.25	3	
B	5.00	2	

But also a formula:

$$U_i = mn + \frac{n(n+1)}{2} - R_i$$

$$U_2 = mn + \frac{m(m+1)}{2} - R_2$$

$$U = \min(R_1, R_2)$$

We can use the normal curve to approximate.

$$E(U) = \frac{mn}{12}$$

$$Z = \frac{U - E(U)}{\sqrt{\text{Var}(U)}} \sim N(0,1)$$

Signed Test

typically on before/after data

X IID from some dist F

$H_0: F = G \quad H_1: F > G$

W ~ Binomial (n, .5)

Use normal approximation

$$E(W) = np \quad \text{Var}(W) = np(1-p)$$

$$Z = \frac{W - E(W)}{\sqrt{\text{Var}(W)}} \sim N(0,1)$$

→ this test has generally low power

Wilcoxon Signed Rank test

General Procedure

1. Remove observations w/ no difference

2. Compute R_i = rank of |D_i|

3. Calculate $W_i = R_i \times \text{sign}(D_i)$

4. Compute test statistic

$$W_+ = \sum_{i: D_i > 0} W_i \quad \text{in other words, sum the positive rank values}$$

If H_0 is true, then:

$$E(W_+) = \frac{n(n+1)}{4}$$

$$\text{Var}(W_+) = \frac{n(n+1)(2n+1)}{24}$$

$$Z = \frac{W_+ - E(W_+)}{\sqrt{\text{Var}(W_+)}} \sim N(0,1)$$

use normal dist to calc p-val

Example

Before	After	diff	abs	rank	signed rank
125	118	7	7	3	3
132	134	-2	2	1	-1
138	130	8	8	4	4
120	124	-4	4	2	-2

In this case, $W_+ = 3+4 = 7$

$$E(W_+) = 4(4+1)/4 = 5 \quad \text{Var}(W_+) = \frac{4(4+1)(24+1)}{24} = \frac{45}{24} = 1.875$$

$$Z = \frac{7 - 5}{\sqrt{1.875}} = \frac{2}{\sqrt{1.875}} = 0.69$$

test stat

Bonferroni Correction

for multiple testing

If we are conducting K test w/ $\alpha = 0.05$

$H_0: M_1 = M_2 = \dots = M_K$ Assume null is true in ALL cases

prob of rejection of ANY test = 5%. but prob of at least one one incorrect rejection of the null

$H_0: M_1 = M_2 = \dots = M_K$ \vdots $H_0: M_K = M_{K+1}$

The family wise error rate (FWER) is the probability of at least one incorrect rejection among K test.

Our goal is to set 0.05, and currently its 99.9% if $K=100$.

→ stricter and lower significance level for each test.

If you have K hypothesis tests, reject H_0 if the p-value is less than α/K , rather than α

our new α becomes $\alpha^{0.05} = 0.0005$

Bias:

An estimate is considered unbiased if $E(\hat{\theta}) - \theta = 0$

that is:

$$E(\hat{\theta}) = \theta$$

Analysis of Variance (ANOVA)

Suppose we have a sample J observations from each of I populations (groups)

$$G_1, G_2, \dots, G_I \quad H_0: \mu_1 = \mu_2 = \dots = \mu_I$$

$H_1: \text{at least one has different mean}$

Sums of Squares:

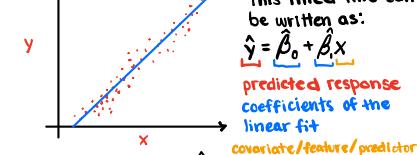
We can capture different types of var using sum of squares

Total Sum of Squares: compares each observation to overall mean

$$SS_T = \sum_{i=1}^J \sum_{j=1}^n ($$

Linear Regression (Continued)

Linear Relationships



The "hat", e.g. $\hat{y}, \hat{\beta}$, means that the quantity is a number computed (estimated) from the data.

Least Absolute Deviation (LAD)

estimating a linear relationship by taking the line that has the lowest average/total prediction error.

choosing a line based on the absolute value error is called LAD

L1 Loss function

$$\begin{aligned} \hat{\beta}_0 &= \arg\min_{\beta_0} \sum |y_i - (\beta_0 + \beta_1 x_i)| && \text{L1 loss is not differentiable} \\ \hat{\beta}_1 &= \arg\min_{\beta_1} \sum |y_i - (\beta_0 + \beta_1 x_i)| && \text{at the minimum} \end{aligned}$$

i.e. the LAD is the line such that the values of $\hat{\beta}_0$ and $\hat{\beta}_1$ are chosen as so to minimize the absolute difference between the observed and predicted response

Least Squares Algorithm

the L2 loss is the sum of the areas of the squares whose edges are the vertical distances from the point to the line
 "Squared" Prediction Err = $(y - \hat{y})^2$

L2 loss function

$$\begin{aligned} \hat{\beta}_0 &= \arg\min_{\beta_0} \sum (y_i - (\beta_0 + \beta_1 x_i))^2 && \text{note L2 loss is differentiable} \\ \hat{\beta}_1 &= \arg\min_{\beta_1} \sum (y_i - (\beta_0 + \beta_1 x_i))^2 && \text{at the minimum} \end{aligned}$$

i.e. the LS line is the line such that the values of $\hat{\beta}_0$ and $\hat{\beta}_1$ are chosen so as to minimize the squared difference between the observed and predicted responses

good thing about least squares is that we can derive a formula for the values that minimize loss for $\hat{\beta}_0$ and $\hat{\beta}_1$.

$$\hat{\beta}_0 = \bar{y} - \bar{\beta}_1 \bar{x} \quad \hat{\beta}_1 = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

Comparing LAD w/ LS

outliers - LS is more sensitive than LAD

Training and Test Set

before you begin, you need to split your data into:

Training Set: data to train your algorithm

Test Set: data to evaluate your algorithm

Commonly Split 60/40, 70/30, 80/20

Mean Squared Error (MSE)

value of L2 loss when you plug in the estimate parameters

$$MSE = \frac{1}{n} \sum (y_i - \hat{y}_i)^2 = \frac{1}{n} \sum (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2$$

Mean Absolute Deviation (MAD)

value of the L1 loss when you plug in the estimated parameters

$$MAD = \frac{1}{n} \sum |y_i - \hat{y}_i| = \frac{1}{n} \sum |y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)|$$

R² value (Coefficient of Determination)

R² value is squared correlation between observed and predicted response

$$R^2 = \frac{P_{y,\hat{y}}}{P_y} \quad 0 < R^2 < 1$$

corresponds to the proportion of variability in the observed response that can be "explained" by the covariate

$$R^2 = \frac{\hat{\sigma}_e^2 - \hat{\sigma}_\beta^2}{\hat{\sigma}_e^2} \quad \text{where } \hat{\epsilon}_i = y_i - \hat{y}_i$$

If $\hat{\sigma}_e^2$ is very small, then \hat{y}_i is very similar to y_i on average, which means x conveys a lot of information about y

Inference for $\hat{\beta}_0$ and $\hat{\beta}_1$

Linear relationships assumptions

ϵ_i is random thus y_i is random

ϵ_i are IID

$E(\epsilon_i) = 0$, $\text{var}(\epsilon_i) = \sigma^2$ for all i

x_i is fixed

β_0, β_1 are parameters that we need to estimate

Bias and Variance of LS estimates of β_0 and β_1

$\hat{\beta}_1$ is unbiased, that is, $E(\hat{\beta}_1) = \beta_1$,

proof: $E(\hat{\beta}_1) = E\left[\frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}\right] = \frac{\sum (x_i - \bar{x})E(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$

$\therefore E(y_i - \bar{y}) = E(\beta_0 + \beta_1 x_i + \epsilon_i - (\beta_0 + \beta_1 \bar{x} + \bar{\epsilon}))$

$= \beta_1(x_i - \bar{x}) + E(\epsilon_i - \bar{\epsilon})$

plug in $= \beta_1(x_i - \bar{x})$ $\therefore \hat{\beta}_1$ is unbiased

$$E(\hat{\beta}_1) = \frac{\beta_1 \sum (x_i - \bar{x})}{\sum (x_i - \bar{x})} = \beta_1$$

$\hat{\beta}_0$ is also unbiased!

Proof: $E(\hat{\beta}_0) = E[\bar{y} - \hat{\beta}_1 \bar{x}]$

$$= E[(\beta_0 + \beta_1 \bar{x} + \bar{\epsilon}) - \beta_1 \bar{x}]$$

$$= \beta_0 + \beta_1 \bar{x} + E(\bar{\epsilon}) - \beta_1 \bar{x}$$

$= \beta_0 \therefore \hat{\beta}_0$ is unbiased!

Variance of $\hat{\beta}_0$ and $\hat{\beta}_1$

The variance and covariance of the LS estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ are given by:

$$\text{var}(\hat{\beta}_0) = \frac{\sigma^2 \sum x_i^2}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\text{var}(\hat{\beta}_1) = \frac{n \sigma^2}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\text{cov}(\hat{\beta}_0, \hat{\beta}_1) = \frac{-\sigma^2 \sum x_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$\hat{\beta}_0$ and $\hat{\beta}_1$ are asymptotically normal

since the LS estimates (w/ normality assumption) are the MLE, this means that they are asymptotically normal

$$\begin{aligned} \hat{\beta}_0 &\sim N\left(\beta_0, \frac{\sigma^2 \sum x_i^2}{n \sum x_i^2 - (\sum x_i)^2}\right) \\ \hat{\beta}_1 &\sim N\left(\beta_1, \frac{n \sigma^2}{n \sum x_i^2 - (\sum x_i)^2}\right) \end{aligned}$$

Hypothesis testing and CI for β_0 and β_1

$$H_0: \beta_j = 0 \quad H_1: \beta_j \neq 0$$

Test Statistic: $t = \frac{\hat{\beta}_j}{\text{se}_{\hat{\beta}_j}}$ where $T \approx t_{n-p}$
 p is # of parameters

$$\begin{aligned} \hat{\beta}_0 &= \sqrt{\frac{\sigma^2 \sum x_i^2}{n \sum x_i^2 - (\sum x_i)^2}} \\ \hat{\beta}_1 &= \sqrt{\frac{n \sigma^2}{n \sum x_i^2 - (\sum x_i)^2}} \end{aligned}$$

and $\hat{\sigma}^2 = \frac{\text{RSS}}{n-p} = \frac{\sum r_i^2}{n-p}$

$\hat{\sigma}^2 = \frac{\sigma^2}{n-p}$

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

Confidence Interval for β_j :

$$[\hat{\beta}_j - t_{n-p, \alpha/2} \hat{\sigma}_{\hat{\beta}_j}, \hat{\beta}_j + t_{n-p, \alpha/2} \hat{\sigma}_{\hat{\beta}_j}]$$

Residual Plots for assessing inference assumptions

Evaluating homoskedasticity

Homoskedasticity - variance of error with each observation is identical and does not depend on x

Heteroskedasticity - the variance of error associated with each observation is different and may depend on x

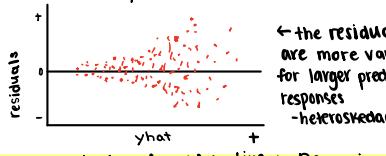
We don't observe ϵ , but perhaps we can get a sense for whether these assumptions are reasonable using the residuals

$$r_i = y_i - \hat{y}_i - \beta_0 - \beta_1 x_i$$

Visualizing Predictive Performance

residual plot

plots the residuals ($y_i - \hat{y}_i$) against the fitted values, \hat{y}_i (or covariate values x_i)



Matrix notation of Multiple Linear Regression

Matrix notation

Obs 1: $y_1 = \beta_0 x_1 + \beta_1 x_{1,1} + \beta_2 x_{1,2} + \dots + \beta_{p-1} x_{1,p-1} + \epsilon_1$

Obs 2: $y_2 = \beta_0 x_1 + \beta_1 x_{2,1} + \beta_2 x_{2,2} + \dots + \beta_{p-1} x_{2,p-1} + \epsilon_2$

Obs n: $y_n = \beta_0 x_1 + \beta_1 x_{n,1} + \beta_2 x_{n,2} + \dots + \beta_{p-1} x_{n,p-1} + \epsilon_n$

Also can be written as:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{1,1} & x_{1,2} & \dots & x_{1,p-1} \\ 1 & x_{2,1} & x_{2,2} & \dots & x_{2,p-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & x_{n,2} & \dots & x_{n,p-1} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Requirements on error vector

$$E(\epsilon) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{Cov}(\epsilon) = \sigma^2 I = \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix}$$

so $\text{Cov}(\epsilon) = \text{Cov}(X\beta + \epsilon) = \text{Cov}(\epsilon) = \sigma^2 I$

To generate a prediction, we need to estimate the vector β , and then compute

$$\hat{y} = X\hat{\beta}$$

LS estimate of Multiple LR

$S(\beta_0, \beta_1, \dots, \beta_{p-1}) = \sum (y_i - (\beta_0 + \beta_1 x_{i,1} + \dots + \beta_{p-1} x_{i,p-1}))^2$

is the same as finding the vector that minimizes

$$S(\beta) = \|y - X\beta\|^2 = (y - X\beta)^T (y - X\beta)$$

this is just a #!

$$\|X\|^2 = \sqrt{X^T X}$$

LS Estimate of β is:

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

$\hat{\beta}$ is unbiased

Proof:

$$E(\hat{\beta}) = E((X^T X)^{-1} X^T y) = (X^T X)^{-1} X^T E(y)$$

$$= (X^T X)^{-1} X^T E(X\beta + \epsilon) = (X^T X)^{-1} X^T X\beta = \beta$$

$\hat{\beta}$ variance-covariance matrix

LS estimate, $\hat{\beta}$ has var/cov matrix:

$$\text{Cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

Proof:

$$\text{Cov}(\hat{\beta}) = \text{Cov}((X^T X)^{-1} X^T y)$$

$$= (X^T X)^{-1} X^T \text{Cov}(y) X (X^T X)^{-1}$$

$$= (X^T X)^{-1} X^T \text{Cov}(\epsilon) X (X^T X)^{-1}$$

$$= (X^T X)^{-1} X^T (\sigma^2 I) X (X^T X)^{-1}$$

$$= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}$$

Theoretical Distribution of $\hat{\beta}$

$$\hat{\beta} \sim N(\beta, \sigma^2 (X^T X)^{-1})$$

each element is $\hat{\beta}_j \sim N(\beta_j, \sigma^2 C_{jj})$

where $C = (X^T X)^{-1}$

this assumes we know σ^2

Hypothesis Testing for β_j

$H_0: \beta_j = 0$ against $H_1: \beta_j \neq 0$

Test statistic: $t = \frac{\hat{\beta}_j}{\text{se}_{\hat{\beta}_j}}$ where $T \approx t_{n-p}$

(T is the # of parameters)

Confidence Interval for β_j

$$[\hat{\beta}_j \pm t_{n-p, \alpha/2} \hat{\sigma}_{\hat{\beta}_j}]$$

where $\hat{\sigma}_{\hat{\beta}_j}^2 = \hat{\sigma}^2 C_{jj}$

and $\hat{\sigma}^2 = \frac{\text{RSS}}{n-p} = \frac{\sum r_i^2}{n-p} = \frac{\|y - X\hat{\beta}\|^2}{n-p}$

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

Correlation + Covariance

correlation is "scaled" covariance

$$-1 \leq \text{corr}(x, y) \leq 1$$

population: $\text{corr}(x, y) = \frac{E((x - \mu_x)(y - \mu_y))}{\sigma_x \sigma_y} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$

data sample:

$$\text{corr}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

OLS Proof:

$$\text{LS coefficient proof } (\hat{\beta}_0)$$

$$S(\hat{\beta}_0, \hat{\beta}_1) = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - \beta_0 - \beta_1 x_i)^2$$

In order to find the values of $\hat{\beta}_0$ and $\hat{\beta}_1$ that minimize $S(\hat{\beta}_0, \hat{\beta}_1)$, we need to set

$$\frac{\partial S}{\partial \hat{\beta}_0} \text{ and } \frac{\partial S}{\partial \hat{\beta}_1}$$

to zero, and solve for $\hat{\beta}_0$ and $\hat{\beta}_1$, respectively.

$$\frac{\partial S}{\partial \hat{\beta}_0} = -2 \sum (y_i - \hat{y}_i - \beta_0) x_i$$

So $\frac{\partial S}{\partial \hat{\beta}_0} = 0$ implies that $0 = \sum (y_i - \hat{y}_i - \beta_0) x_i$

$$\Rightarrow n \hat{\beta}_0 = \sum y_i x_i - \sum \hat{y}_i x_i = \sum y_i x_i - \sum \hat{y}_i x_i$$

$$Substituting \hat{\beta}_0 with \hat{\beta}_0 = \bar{y} - \bar{\beta}_1 \bar{x}: \quad 0 = \sum y_i x_i - (\bar{y} - \bar{\beta}_1 \bar{x}) \sum x_i - \bar{\beta}_1 \sum x_i^2$$

$$Rearranging to make \hat{\beta}_1 the subject: \quad \hat{\beta}_1 = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

$$S(\hat{\beta}_0, \hat{\beta}_1) = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - (\bar{y} - \bar{\beta}_1 x_i))^2$$

→ find $\frac{\partial S}{\partial \hat{\beta}_0}$ and $\frac{\partial S}{\partial \hat{\beta}_1}$ then set = 0 to get proof. Use $\hat{\beta}_0$ in β , proof

Remembering Linear Algebra

if $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ then $A^{-1} = \frac{1}{\det(A)} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$

Det(A) = $a \cdot d - b \cdot c = ad - bc$

If all of these are matrices:

$$A+B = B+A$$

$$A+(B+C) = (A+B)+C$$

$$A+O = A$$

$$A(BC) = (AB)C$$

$$A(B+C) = AB + AC$$

$$IA = AI = A$$

Let r and s be scalars

$$r(sA) = (rs)A$$

$$(r+s)A = rA + sA$$

$$A(rB) = r(AB) = (rA)B$$

$$(A^T)^T = A$$

$$(A+B)^T = A^T + B^T$$

$$(AB)^T = B^T A^T$$

$$(rA)^T = rA^T$$

Probability Properties

$$0 \leq P(A) \leq 1$$

$$P(A) = 1 - P(A') \rightarrow 1 = P(A) + P(A')$$

if A and B are mutually exclusive $P(A \cap B) = 0$

$$P(A \cup B) = P(A) + P$$

