

References

Introduction to Linear Regression Analysis

- Douglas C. Montgomery

Applied Regression Analysis: a research tool

- John O. Rawlings

Brockwell, Peter J. - Intro. to TS and Forecasting

Midsem 30M Endsem 50M Project 20M

Regression Problem

Predicting dependent variable on the basis of independent variable(s).

- might be a categorical value / real value (of both the independent / dependent variable)

Polynomial regression also falls in the category of linear - regression.

Definition Conditional expectation of Y given X , is known as regression of Y on X . It can be denoted as $\hat{y} = g(x, \beta) = E[Y|X=x]$

- $g(x, \beta)$ can be a line, curve, plane, surface ...
- x can be stochastic / non-stochastic
- Y is always stochastic or a random variable.

Associated statistical problems with regression

- Interval Estimation
- Testing of Hypothesis
- Estimability of certain combination of regressors
Admissibility

[All combination of regressors won't be estimated.
given data, some combination might]

- Contrast
- Necessity of building regression : ANOVA
- Classification (if categorical prediction)

Associated mathematical problems with regression

- Notion of error
- How error space would interact with prediction space
- optimal subspace of regressors
- Reducing bias and variance of regression coeff.
- Rate of reduction of error if we incorporate more regressors

Vector Representation

1. Simple Linear Regression

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

↪ ↓ →
 dependent independent random error
 Data = { $(y_i, x_i) | i=1, 2, \dots, n$ }

$$\tilde{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad X = [\tilde{1} \ \tilde{x}] = \left[\begin{array}{c|c} \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} & \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \end{array} \right] \quad \epsilon = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$$

$$\tilde{Y} = X \beta + \epsilon$$

$(n \times 1)$ $(n \times 2)$ (2×1) $(n \times 1)$
 Problem in \mathbb{R}^n

Mathematics on \mathbb{R}^n

Vector Space $(V, +, \cdot)$ where $+ : V \times V \rightarrow V$; $\cdot : F \times V \rightarrow V$

$$V = \mathbb{R}^n; F = \mathbb{R}$$

$$lP_n : \left\{ \sum_{i=0}^n a_i x^i \mid a_i \in \mathbb{R}, x \in \mathbb{R} \right\}$$

$$\underline{a} = \begin{pmatrix} a_0 \\ \vdots \\ a_n \end{pmatrix} \quad \underline{x} = \begin{pmatrix} 1 \\ x \\ \vdots \\ x^n \end{pmatrix} \quad lP_n = \left\{ \underline{a}^T \underline{x} \mid \underline{a} \in \mathbb{R}^n; x \in \mathbb{R} \right\}$$

Since a polynomial can be identified by one variable irrespective of the domain, it still is a linear regression problem.

All the coefficient are identified uniquely by an $n \times 1$ vector (\underline{a} in the above case)

Subspace $S \subseteq V$ and it satisfies all the properties of V , then S is called subspace of V .

- check if $0 \in S$
- check if $x+ay \in S$?

Example $V = \mathbb{R}^2$

- $S_1 = \{(x, 0), x \in \mathbb{R}\}$
- $S_2 = \{(0, y), y \in \mathbb{R}\}$
- $S_3 = \{(x, mx), x \in \mathbb{R}\}$

generalizing, $S = \left\{ \underline{a}^T \underline{x} = a_1 x_1 + a_2 x_2 + \dots + a_n x_n = 0 \mid a_i \in \mathbb{R}, x_i \in \mathbb{R} \right\}_{i=1,2,\dots,n}$

$S \subseteq V$ and S is subspace.

- For $V = \mathbb{R}^3$, $S = \left\{ \underline{a}^T \underline{x} = a_1 x_1 + a_2 x_2 + a_3 x_3 = 0 \mid a_i \in \mathbb{R}, x_i \in \mathbb{R}, i=1,2,3 \right\}$
- $V = lP_n$ $S = lP_k$, $k \leq n$ S is a subspace of lP_n .
- 0 element is always a subspace (trivial subspace)**

- For $V = \mathbb{P}_n$, would $S = \{k\text{-degree polynomial}\}$ be a subspace?

Not necessarily. (only 0 elt. would be a subspace)

Ex. $v_1 = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}$ and $v_2 = \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix}$, then $v_1 + v_2 = \begin{pmatrix} 0 \\ 2 \\ 0 \end{pmatrix}$ is

not a 2-degree polynomial. (Not closed under addition)

given \tilde{x} , we predict \hat{y} if we know $\hat{\beta}_0$ and $\hat{\beta}_1$.

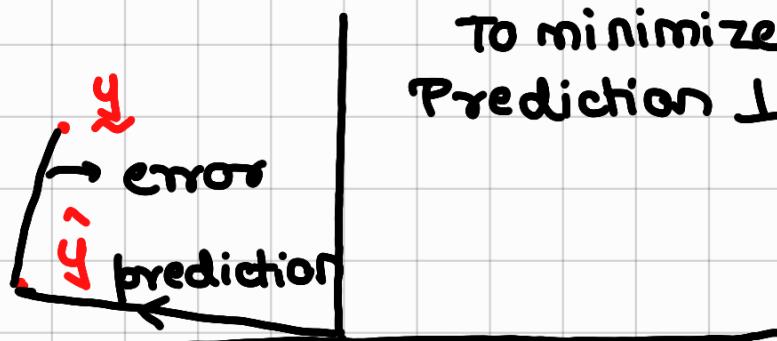
$$\tilde{Y} = X\beta + \epsilon$$

$$= \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \epsilon$$

$$\hat{y} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \hat{\beta}_0 + \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \hat{\beta}_1 = \underbrace{\frac{1}{n} \hat{\beta}_0}_{\text{prediction}} + \underbrace{\frac{x}{n} \hat{\beta}_1}_{\text{given from } \mathbb{R}^n}$$

We have 2 vectors $\tilde{1}$ and \tilde{x} . But their linear combination will not ever span full \mathbb{R}^n space.
($n \leq 1 ??$)

Entire regression problem is that you will try to predict \hat{y} but would not be able to span whole \mathbb{R}^n space.



To minimize error,
Prediction \perp error.

Recap of Linear independence, Span and Basis & dimensions.

Orthogonal Vectors

vectors $\underline{u}, \underline{v} \in \mathbb{R}^n$ are said to be orthogonal if
 $\langle \underline{u}, \underline{v} \rangle = \underline{u}^T \underline{v} = \sum_{i=1}^n u_i v_i = 0$

functions $u(t) : [a, b] \rightarrow \mathbb{R}$ and $v(t) : [a, b] \rightarrow \mathbb{R}$
are said to be orthogonal functions if
 $\int_a^b u(t)v(t) dt = 0$

Orthogonal Complement

Projection Matrix

Idempotent Matrix

A projection matrix is an idempotent matrix (Prove)
An idempotent matrix has eigenvalues 0 or 1 (Prove)
If $\{\underline{v}_1, \underline{v}_2, \dots, \underline{v}_k\}$ is an orthonormal basis the subspace
 $S \subseteq V$ then the orthogonal projection matrix of S
is $P_S = \sum_{i=1}^k \underline{v}_i \underline{v}_i^T$

Column Space

Properties:

- 1) $C(A+B) = C(A) + C(B)$
- 2) $C(AB) \subseteq C(A)$
- 3) $\dim(C(A)) = \text{Rank}(A)$
- 4) $C(AA^T) = C(A) \Rightarrow \text{Rank}(AA^T) = \text{Rank}(A)$ [Prove it]

Quadratic Form $\underline{v}^T A \underline{v} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} v_i v_j$

Definition of Metric (Σ should be symm. and psd)

Statistics on \mathbb{R}^n

$\underline{x} = (x_1, x_2, \dots, x_n)^T$ a random vector with finite expectation for each component

Existence of Expectation

$E(x)$ exists $\Rightarrow \int |x| f(x) dx < \infty$, or,
 $\sum_x |x| f(x) < \infty$

Dispersion matrix

$$\begin{aligned} D(x) &= ((\text{Cov}(x_i, x_j)))_{n \times n} \\ &= E[(x - E(x))(x - E(x))^T] = \text{Cov}(x, x) \\ &= \Sigma, \text{ or,} \\ D(x) &= \Sigma_x = ((\sigma_{ij}))_{n \times n} = ((\text{Cov}(x_i, x_j))) \end{aligned}$$

Note $\text{Cov}(U_p, V_q) = ((\text{Cov}(U_i, V_j)))_{p \times q}$,

$$E(x+b) = E(x)+b$$

$$D(x+b) = D(x)$$

$$\text{Cov}(x+b, y+c) = \text{Cov}(x, y)$$

$$\begin{aligned} \text{Cov}(x_i, x_j) &= E(x_i x_j) - E(x_i) E(x_j) \\ &= E(x_j x_i) - E(x_j) E(x_i) = \text{Cov}(x_j, x_i) \end{aligned}$$

→ commutative in nature

$$\text{Hence, } \Sigma_x^T = \Sigma_x$$

Important Results

x ← random vector with n -components

$$E(x) = \mu ; D(x) = \Sigma$$

$$1) E(I^T x) = I^T \mu, I \in \mathbb{R}^n \leftarrow \text{constant vector}$$

$$2) D(I^T x) = I^T \Sigma I$$

$$3) E(AX) = A\mu, A \in \mathbb{R}^{p \times n} \leftarrow \text{constant matrix}$$

$$4) D(AX) = A\Sigma A^T \text{ and } \text{Cov}(AX, BX) = A\Sigma B^T$$

$$5) \text{ If } \text{Cov}(U_p, V_q) = \Gamma \text{ then } \text{Cov}(AU, BV) = A\Gamma B^T$$

$$5) \Rightarrow 4) \Rightarrow 2) ; 3) \Rightarrow 1)$$

$$E(\tilde{\boldsymbol{l}}^T \tilde{\mathbf{x}}) = E \sum_{i=1}^n l_i x_i = \sum l_i E(x_i) = \tilde{\boldsymbol{l}}^T \mu$$

$$D(\tilde{\boldsymbol{l}}^T \tilde{\mathbf{x}}) = \text{Var}(\tilde{\boldsymbol{l}}^T \tilde{\mathbf{x}}) = \text{Cov}(\tilde{\boldsymbol{l}}^T \tilde{\mathbf{x}}, \tilde{\boldsymbol{l}}^T \tilde{\mathbf{x}})$$

$$= \sum_i \sum_j l_i \text{Cov}(x_i, x_j) l_j$$

$$= \tilde{\boldsymbol{l}}^T \sum \tilde{\boldsymbol{l}}$$

$D(\mathbf{x}) = \Sigma$ is symmetric

$D(\mathbf{x}) = \Sigma$ is positive semi-definite

To show $D(\tilde{\mathbf{x}}) = \Sigma$ is positive semi-definite.

$$\mathbf{y} = \tilde{\boldsymbol{l}}^T \tilde{\mathbf{x}} \in \mathbb{R} \quad \tilde{\boldsymbol{l}} \neq 0$$

We know $\text{Var}(\mathbf{y}) \geq 0$

$$\text{Var}(\mathbf{y}) = \text{Var}(\tilde{\boldsymbol{l}}^T \tilde{\mathbf{x}}) = \tilde{\boldsymbol{l}}^T \Sigma \tilde{\boldsymbol{l}} \geq 0 \quad \forall \tilde{\boldsymbol{l}} \neq 0$$

$\Rightarrow \Sigma$ is positive semi-definite.

Theorem $P((\mathbf{x} - \mu) \in C(\Sigma)) = 1$.

Proof

$$\mathbf{x} \in \mathbb{R}^n, E(\mathbf{x}) = \mu \in \mathbb{R}^n, D(\mathbf{x}) = \Sigma$$

To show: $P((\mathbf{x} - \mu) \in C(\Sigma)) = 1$.

It is equivalent to show if $\tilde{\boldsymbol{l}} \in (C(\Sigma))^{\perp}$, then $\tilde{\boldsymbol{l}} \perp (\tilde{\mathbf{x}} - \tilde{\mu})$.

$$\begin{aligned} \tilde{\boldsymbol{l}} \in (C(\Sigma))^{\perp} &\iff \tilde{\boldsymbol{l}}^T \Sigma = \tilde{\boldsymbol{o}}^T \iff \tilde{\boldsymbol{l}}^T \Sigma \tilde{\boldsymbol{l}} = \tilde{\boldsymbol{o}}^T \tilde{\boldsymbol{l}} = 0 \\ &\iff \text{Var}(\tilde{\boldsymbol{l}}^T \tilde{\mathbf{x}}) = 0 \\ &\iff \text{Var}(\tilde{\boldsymbol{l}}^T (\tilde{\mathbf{x}} - \tilde{\mu})) = 0 \end{aligned}$$

$$\text{Now, } E(\tilde{\mathbf{x}} - \tilde{\mu}) = \mathbf{0} \Rightarrow E(\tilde{\boldsymbol{l}}^T (\tilde{\mathbf{x}} - \tilde{\mu})) = 0.$$

We have

$$E(\tilde{\boldsymbol{l}}^T (\tilde{\mathbf{x}} - \tilde{\mu})) = 0$$

$\text{Var}(\tilde{\boldsymbol{l}}^T (\tilde{\mathbf{x}} - \tilde{\mu})) = 0$. This implies,

$$P(\tilde{\boldsymbol{l}}^T (\tilde{\mathbf{x}} - \tilde{\mu}) = 0) = 1 \Rightarrow P(\tilde{\boldsymbol{l}} \perp (\tilde{\mathbf{x}} - \tilde{\mu})) = 1.$$

Hence, proved.

Some cases to consider

$$A = I_n \quad E(\tilde{x}) = \mu_{\tilde{x}}, \quad D(\tilde{x}) = \Sigma$$

Case-I :

$$\begin{aligned} E(\tilde{x}^T A \tilde{x}) &= E(\tilde{x}^T \tilde{x}) = E\left(\sum_{i=1}^n x_i^2\right) \\ &= \sum_{i=1}^n (E(x_i^2)) \\ &= \sum_{i=1}^n (\mu_i^2 + \sigma_{ii}) \\ &= \sum_{i=1}^n \sigma_{ii} + \sum_{i=1}^n \mu_i^2 \end{aligned}$$

$$E(\tilde{x}^T A \tilde{x}) = \text{tr}(\Sigma I_n) + \mu_{\tilde{x}}^T I_n \mu_{\tilde{x}}$$

Case-II :

If $x_i \sim N(0, 1)$

$$E(\tilde{x}^T \tilde{x}) = E\left(\sum_{i=1}^n x_i^2\right) = n.$$

$$X_n^2 \sim \text{Gamma}(\alpha = n/2, \lambda = 1/2)$$

$$E[G] = \alpha/\lambda = n; \quad V(G) = \alpha/\lambda^2 = 2n.$$

$$\text{Hence, } E(\tilde{x}^T \tilde{x}) = n, \quad V(\tilde{x}^T \tilde{x}) = 2n.$$

$$\begin{aligned} E(\tilde{x}^T A \tilde{x}) &= E \sum_i \sum_j a_{ij} E(x_i x_j) \\ &= \sum_{i=1}^n \sum_{j=1}^n a_{ij} (\sigma_{ij} + \mu_i \mu_j) \\ &= \sum_i \sum_j a_{ij} \sigma_{ij} + \sum_i \sum_j a_{ij} \mu_i \mu_j \\ &= \text{tr}(\Sigma A) + \mu_{\tilde{x}}^T A \mu_{\tilde{x}} \end{aligned}$$

$$E(\tilde{x}^T A \tilde{x}) = E[\text{tr}(\tilde{x}^T A \tilde{x})]$$

(... $\tilde{x}^T A \tilde{x}$ is a scalar quantity)

$$= E[\text{tr}(A \tilde{x} \tilde{x}^T)]$$

(... $\text{tr}(AB) = \text{tr}(BA)$)

$$= \text{tr}[E(A(\tilde{x} \tilde{x}^T))]$$

$$= \text{tr}[A(\Sigma + \mu \mu^T)] = \text{tr}(A\Sigma + A\mu \mu^T)$$

$$\begin{aligned}
 &= \text{tr}(A\Sigma) + \text{tr}(\mu^T A \mu) \\
 &= \text{tr}(A\Sigma) + \mu^T A \mu \\
 &\quad (\dots \mu^T A \mu \text{ is a scalar quantity})
 \end{aligned}$$

Multivariate Normal Distribution

$$f(x) = \frac{\exp\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\}}{(2\pi)^{n/2} |\Sigma|^{1/2}}$$

for some $\mu \in \mathbb{R}^n$ and p.s.d. Σ

$$\text{Univariate Normal: } f(x) = \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

If Σ is a diagonal matrix,

$$f(x) = \frac{\exp\left\{-\frac{1}{2} \sum_{i=1}^n (x_i - \mu_i)^2 / \sigma_i^2\right\}}{(\sqrt{2\pi})^n (\prod_{i=1}^n \sigma_i)}$$

If $X \sim N(\mu, \Sigma)$ then $AX \sim N(A\mu, A\Sigma A^T)$

Ex: If (X, Y) follow Bivariate Normal, find $f(x, y)$, i.e., pdf of (X, Y) .

Ex: Show that $Y|X=x$ follows $N(\mu_y + \rho \frac{\sigma_y}{\sigma_x} (x - \mu_x), (1 - \rho^2) \sigma_y^2)$ if (X, Y) follows Bivariate Normal distn.

$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim N\left[\begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \sigma_x^2 & \rho \sigma_x \sigma_y \\ \rho \sigma_x \sigma_y & \sigma_y^2 \end{pmatrix}\right]$$

Mahalanobis distance

Notion of distance when there is a fluctuation
 $(X - \mu)^T \Sigma^{-1} (X - \mu)$
 $\sum_{i=1}^n (x_i - \mu_i)^2 / \sigma_i^2$ for diagonal Σ .

Suppose $\mathbf{z} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \sigma_x^2 & \rho\sigma_x\sigma_y \\ \rho\sigma_x\sigma_y & \sigma_y^2 \end{pmatrix}\right)$

$$\begin{aligned} \mathbf{l}_1 &= (1, 0)^T & \mathbf{l}_1^T \mathbf{z} &\sim N(\mu_x, \sigma_x^2) \\ \mathbf{l}_2 &= (0, 1)^T & \mathbf{l}_2^T \mathbf{z} &\sim N(\mu_y, \sigma_y^2) \end{aligned}$$

Chi-Squared distribution

If $x_i \sim N(0, 1)$, then $\sum_{i=1}^n x_i^2$ follow chi-square distn. with n degree of freedom

If $\mathbf{x} \sim N(\boldsymbol{\mu}, I_n)$, then $\mathbf{x}^T \mathbf{x}$ follows chi-square distn. with n degrees of freedom and **non-centrality parameter** $\boldsymbol{\mu}^T \boldsymbol{\mu}$

$$E(\mathbf{x}^T \mathbf{x}) = n + \boldsymbol{\mu}^T \boldsymbol{\mu}.$$

for $n=1$,

$$\mathbf{x} \sim N(0, 1)$$

$$\mathbf{x}^2 \sim \chi_1^2 \quad E(\mathbf{x}^2) = 1$$

$$\mathbf{y} \sim N(\boldsymbol{\mu}, 1), \text{ then } E(\mathbf{y}^2) = \boldsymbol{\mu}^2 + 1.$$

Theorem If $\mathbf{x} \sim N(\boldsymbol{\mu}, I_n)$ then $\mathbf{x}^T A \mathbf{x}$ has chi-sq. iff A is idempotent

Moreover, $\mathbf{x}^T A \mathbf{x} \sim \chi_{df=\text{Rank}(A), ncp=\boldsymbol{\mu}^T A \boldsymbol{\mu}}^2$

Cochran's Theorem

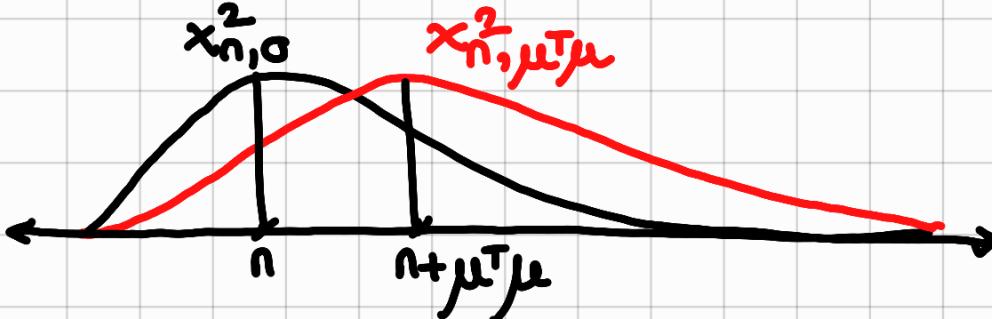
Let $\mathbf{x} \sim N(\boldsymbol{\mu}, I_n)$ and $\mathbf{x}^T A \mathbf{x} \equiv \sum_{i=1}^k \mathbf{x}^T A_i \mathbf{x}$ where A_i 's are symmetric and A is an idempotent matrix. Then $\mathbf{x}^T A_i \mathbf{x} \sim \chi_{\text{Rank}(A_i), \boldsymbol{\mu}^T A_i \boldsymbol{\mu}}$ and they are independent

$$Z \sim N(0, I_n)$$

$$X = Z + \mu, \mu \neq 0$$

$$Z^T I_n Z \sim \chi^2_{df=n, ncp=0}$$

$$X^T X \sim \chi^2_{df=n, ncp=\mu^T \mu > 0}$$



Ncp parameter has important significance while performing tail tests.

Construction of t-test and F-test

Let $x_i \stackrel{iid}{\sim} N(0, \sigma^2)$. Define $\hat{x} = \frac{1}{n} \sum_{i=1}^n x_i$ and $s^2 = \sum_{i=1}^n (x_i - \hat{x})^2$

$Z \sim N(0, 1)$ independent. Then, $\frac{Z}{\sqrt{Y/k}} \sim t_{k, ncp}$

$x_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$. Then $\tilde{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}, x \sim N\left[\mu \cdot \underbrace{\frac{1}{n} \cdot 1_T}_{\text{mean}}, \sigma^2 I_n\right]$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \tilde{x}^T \underbrace{\frac{1}{n} \cdot 1_T}_{\ell^T} \tilde{x} = \ell^T \tilde{x} \quad (\dots \ell = \frac{1}{n} \cdot 1^T)$$

$$\bar{x} \sim N(\ell^T \mu, \ell^T \Sigma \ell) = N(\mu, \sigma^2 \frac{1}{n})$$

As a consequence,

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

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

$$P = (I_n - \frac{1}{n} \cdot 1 \cdot 1^T)$$

Show that : i) $P^T = P$ ii) $P^2 = P$ iii) $\ell^T P = \tilde{\ell}^T$

$$\lambda_n \cdot 1^T (I_n - \lambda_n 1 1^T) = \lambda_n (1^T - \lambda_n n \cdot 1^T) = \lambda_n (1^T - 1^T)$$

$$S^2 \sim \chi_{df=\text{rank}(P), n-p=0}^2$$

Tutorial Class :

1. Prove that $C(AA^T) = C(A)$

We prove, i) $C(AA^T) \subseteq C(A)$
ii) $C(A) \subseteq C(AA^T)$

Let $\tilde{\ell} \in [C(AA^T)]^\perp$

$$\iff \tilde{\ell}^T AA^T = \tilde{\ell}^T$$

$$\iff \tilde{\ell}^T AA^T \tilde{\ell} = \tilde{\ell}^T \tilde{\ell} = 0$$

$$\iff (\tilde{\ell}^T A^T \tilde{\ell})^T (\tilde{\ell}^T A^T \tilde{\ell}) = 0$$

$$\iff A^T \tilde{\ell} = 0$$

$$\iff \tilde{\ell}^T A = \tilde{\ell}^T$$

$$\iff \tilde{\ell} \in [C(A)]^\perp$$

$\therefore [C(A)]^\perp \subseteq [C(AA^T)]^\perp$ and vice versa.

$$\Rightarrow C(AA^T) = C(A)$$

2 a) $V = \mathbb{R}^3$, $S = \{(x_1, x_2, x_3) : 2x_1 + x_2 + x_3 = 1\}$ No

b) $V = \mathbb{R}^2$, $S = \{(x_1, x_2) : x_1, x_2 \geq 0\}$

→ S over \mathbb{R} ? No

→ S over $\mathbb{R}^+ \cup \{0\}$? Yes

→ S over \mathbb{Z}_n ? Yes

Show that \bar{x} and S^2 independently distribute

$$x_i \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2) \Leftrightarrow \tilde{x} \sim N(\tilde{\mu}, \sigma^2 I_n)$$

$$\begin{aligned} \mathbf{1}^T &= \frac{1}{n} \mathbf{1}^T \\ \bar{x} &= \frac{1}{n} \mathbf{1}^T \tilde{x} \end{aligned}$$

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

$$S^2 = \frac{\tilde{x}^T P \tilde{x}}{n}$$

$$P = \left(I_n - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right)$$

$$P^T = P$$

$$P^2 = P$$

$$\mathbf{1}^T P = \mathbf{0}^T$$

Homework

$$\tilde{\mathbf{l}}^T \tilde{x} \sim N(\tilde{\mathbf{l}}^T \tilde{\mu}, \tilde{\mathbf{l}}^T \Sigma \tilde{\mathbf{l}})$$

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

$$\text{AS } x \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$$

$$\tilde{\gamma} = \left(\frac{x}{\sigma} \right) \sim N\left(\frac{\mu}{\sigma}, I_n \right)$$

$$S^2 = \frac{\tilde{x}^T P \tilde{x}}{n} = \sigma^2 \left(\frac{\tilde{x}}{\sigma} \right)^T P \left(\frac{\tilde{x}}{\sigma} \right) = \sigma^2 (\tilde{\gamma}^T P \tilde{\gamma})$$

$$\text{AS } x_i \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$$

$$y_i \stackrel{\text{iid}}{\sim} N(\mu_\sigma, 1)$$

$$\therefore \tilde{\gamma}^T P \tilde{\gamma} \sim \chi^2_{df = \text{rank}(P), ncp = (\mu/\sigma)^T P (\mu/\sigma)}$$

$$\text{Rank}(P) = (n-1)$$

$$(\mu^2/\sigma^2) \frac{1}{n} \mathbf{1}^T P \frac{1}{n} = \frac{1}{n} \mathbf{1}^T \frac{1}{n} = 0$$

$$I_n = (I_n - \frac{1}{n} \mathbf{1} \mathbf{1}^T) + \frac{1}{n} \mathbf{1} \mathbf{1}^T$$

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

$$S^2 = \sigma^2 \chi^2_{df = n-1, ncp = 0}$$

As a result, \bar{X} and S^2 are independent.

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

$$\frac{S^2}{\sigma^2} \sim \chi_{n-1}^2, n \neq 0 \quad \begin{matrix} \text{independent} \\ \swarrow \end{matrix}$$

t-statistic

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

$$x_i: \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2) \iff \tilde{x} \sim N(\mu, \sigma^2 I_n)$$

$$S^2 = \sigma^2 Y^T P Y$$

$$t = \frac{\sqrt{n}(\bar{X} - \mu)/\sigma}{\sqrt{\frac{(S^2/\sigma^2)/(n-1)}}}$$

$$t^2 = \frac{(\sqrt{n}(\bar{X} - \mu)/\sigma)^2 / 1}{(S^2/\sigma^2)/(n-1)} = \frac{Z^2 / 1}{\frac{S^2}{\sigma^2} / (n-1)} \sim F_{1, n-1}$$

Formulation of Hypothesis,

$$H_0: \mu = 5$$

$$H_1: \mu \neq 5$$

σ^2 unknown

Suppose we have the above hypothesis. What type of deformation would happen to F-distn. if the mean is shifted?

A non-centrality parameter would be introduced

in the numerator, and the expectation would increase. This forms a good measure to evaluate hypothesis



We construct t-statistic from the data.
used as **Upper tail test**

Upper Tail test

Under H_1 ,

$$\bar{X} \sim N(\mu - \mu_0, \sigma^2/n)$$

$$\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \sim N\left(\frac{\sqrt{n}}{\sigma}(\mu - \mu_0), 1\right)$$

References

Mathematical Statistics , Miller

Simple Linear Regression

$$y = \beta_0 + \beta_1 x + \epsilon$$

relation b/w response variable y and regressor variable x is linear in parameter, then it is called simple linear regression model.

Even $y = \beta_0 + \beta_1 e^x + \epsilon$ is a simple linear regression model.

But, $y = \frac{1}{\beta_0 + \beta_1 x} + \epsilon$, $y = \Phi(\beta_0 + \beta_1 x) + \epsilon$ is not.

Linear Model $y_i = \beta^T x_i + \epsilon_i$

Gauss-Markov Model $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$, $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$

Here, $\beta_0 \in \mathbb{R}$, $\beta_1 \in \mathbb{R}$, $\sigma > 0$ are unknown model parameters.

$$E(y_i) = \beta_0 + \beta_1 x_i$$

$$\text{Var}(y_i) = \sigma^2$$

$$\rightarrow y_i \sim N(\beta^T x_i, \sigma^2)$$

For simple linear regression,

$$y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

$\beta = (\beta_0, \beta_1)$ and σ^2 unknown

We define

$$\tilde{Y} = (Y_1, Y_2, \dots, Y_n)^T; \quad \tilde{\epsilon} = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)^T$$

$$\tilde{\beta} = (\beta_0, \beta_1)^T; \quad \tilde{x} = (x_1, x_2, \dots, x_n)^T$$

$$X = [I \quad \tilde{x}]$$

$$\tilde{Y} = X \tilde{\beta} + \tilde{\epsilon}$$

$$\tilde{Y} \sim N(X \tilde{\beta}, \sigma^2 I_n) \text{ as } \tilde{\epsilon} \sim N(0, \sigma^2 I_n)$$

It is necessary to define a model with the same metric throughout

Estimation of model parameters

least squared model

- Data: $\{(x_i, y_i) : i=1, 2, \dots, n\}$
- Model: $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ assumed
- Metric: square distance assumed

least squared condition

$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

$$S(\beta) = (\tilde{Y} - \tilde{X}\tilde{\beta})^T (\tilde{Y} - \tilde{X}\tilde{\beta})$$

$$= \tilde{Y}^T \tilde{Y} - \tilde{Y}^T \tilde{X} \tilde{\beta} - \tilde{\beta}^T \tilde{X}^T \tilde{Y} + \tilde{\beta}^T \tilde{X}^T \tilde{X} \tilde{\beta}$$

$$\frac{\partial S(\beta)}{\partial \beta} = 0$$

$$\left(\begin{array}{l} \frac{\partial \tilde{X}^T \tilde{A} \tilde{\beta}}{\partial \tilde{\beta}} = \tilde{\beta}^T (\tilde{A} + \tilde{A}^T) \\ \frac{\partial \tilde{A}^T \tilde{\beta}}{\partial \tilde{\beta}} = \frac{\partial \tilde{\beta}^T \tilde{A}}{\partial \tilde{\beta}} = \tilde{A}^T \end{array} \right)$$

This implies,

$$^{2x2} (\tilde{X}^T \tilde{X}) \tilde{\beta} = \tilde{X}^T \tilde{Y}$$

Assuming $\tilde{X}^T \tilde{X}$ is invertible,

$$\tilde{\beta} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \tilde{Y}$$

When will $\tilde{X}^T \tilde{X}$ not be invertible?

- When $\text{rank}(\tilde{X}^T \tilde{X}) < 2$, i.e., all x_i values are same.

Homework Find β_0, β_1 from the above result and compare it with the result which we got from the usual derivative.

$$\hat{\beta}_1 = S_{xy}/S_{xx}, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\text{Hence, } \hat{y}_{\text{new}} = \hat{\beta}_0 + \hat{\beta}_1 x_{\text{new}}$$

$$= \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_{\text{new}}$$

$$= \bar{y} + \hat{\beta}_1 (x_{\text{new}} - \bar{x})$$

passing through (\bar{x}, \bar{y})

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n (x_i - \bar{x})y_i$$

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

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} ; S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2$$

non-random

$$= \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{S_{xx}}$$

$$= \hat{l}^T \tilde{y} ; l_i = (x_i - \bar{x}) / S_{xx}$$

$$\hat{l} = (l_1 \ l_2 \ \dots \ l_n)^T$$

Linear Estimator

$$\tilde{y} \sim N(X\beta, \sigma^2 I_n)$$

$$\hat{\beta}_1 = \hat{l}^T \tilde{y} \sim N(l^T X \beta, \sigma^2 l^T I_n l)$$

$$E(\hat{\beta}_1) = l^T X \beta$$

$$= [0 \ 1]^T \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} = \beta_1$$

$$\text{var}(\hat{\beta}_1) = (l^T l) \sigma^2 = \sigma^2 \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{S_{xx}^2} = \sigma^2 / S_{xx}$$

$\hat{\beta}_1$ is an unbiased estimator of β_1 .

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = \frac{1}{n} 1^T \tilde{y} - \frac{\bar{x}}{S_{xx}} \sum_{i=1}^n (x_i - \bar{x}) y_i$$

$$= \frac{1}{n} 1^T \tilde{y} - \bar{x} \hat{l}^T \tilde{y}$$

$$= (\frac{1}{n} 1^T - \bar{x} \hat{l}^T) \tilde{y}$$

$$= \sum_{i=1}^n \left(\frac{1}{n} - \bar{x} \frac{(x_i - \bar{x})}{S_{xx}} \right) y_i$$

SLR

$$\hat{\beta}_1 \sim N(\beta_1, \sigma^2 / S_{xx}) ; \hat{\beta}_0 \sim N(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right))$$

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

$$\hat{y} = x \hat{\beta} = [x (x^T x)^{-1} x^T] \underline{y} = P_x \underline{y}$$

$$P_x^T = P_x \text{ and } P_x^2 = P_x \text{ (Homework)}$$

$$\epsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\Rightarrow y_i \stackrel{\text{independent}}{\sim} \mathcal{N}(\beta_0 + \beta_1 x_i, \sigma^2)$$

$$\underline{y} \sim \mathcal{N}(x \hat{\beta}, \sigma^2 I_n)$$

Estimated Error

$$e_i = y_i - \hat{y}_i \quad i=1, 2, \dots, n$$

$$\underline{e} = (e_1, e_2, \dots, e_n)^T$$

$$\underline{e} = (\underline{y} - \hat{\underline{y}}) = (I - P_x) \underline{y}$$

$$\underline{e} \sim \mathcal{N}(0, \sigma^2 (I_n - P_x)) \text{ no more iid}$$

[Using,
 $\underline{z} \sim \mathcal{N}(\mu, \Sigma) \Rightarrow A\underline{z} \sim \mathcal{N}(A\mu, A\Sigma A^T)$]

Squared estimate error

$$\sum_{i=1}^n e_i^2 = \underline{e}^T \underline{e}$$

$$= [(\underline{I}_n - P_x) \underline{y}]^T (\underline{I}_n - P_x) \underline{y}$$

$$= \underline{y}^T (\underline{I}_n - P_x) (\underline{I}_n - P_x) \underline{y}$$

$$= \underline{y}^T (\underline{I}_n - P_x) \underline{y}$$

$$\frac{\sum_{i=1}^n e_i^2}{\sigma^2} \sim \chi^2$$

$\text{df} = \text{rank}(\underline{I}_n - P_x)$
 $\text{ncp} = \frac{1}{\sigma^2} [(x \hat{\beta})^T (\underline{I}_n - P_x) x \hat{\beta}]$

$$e^T e \sim \chi^2_{(n-2)} \dots (I_n = P_x + (I_n - P_x) \therefore n = 2 + (n-2))$$

$$ncp = \frac{1}{\sigma^2} \beta^T X^T (I_n - P_x) X \beta$$

$$= \frac{1}{\sigma^2} \beta^T [X^T X - X^T P_x X] \beta$$

$P_x = X^T (X^T X)^{-1} X$
 orthogonal projection matrix for the column space
 of X .
 $\therefore P_x X = X$

Hence, $ncp = \frac{1}{\sigma^2} \beta^T O \beta = 0 \Rightarrow \frac{e^T e}{\sigma^2} \sim \chi^2_{(n-2), 0}$

$$E\left(\frac{e^T e}{\sigma^2}\right) = n-2 \Rightarrow E\left(\frac{e^T e}{n-2}\right) = \sigma^2$$

$$C(X) = \text{span}\{x_1, x_2, \dots, x_n\}$$

$X X^T$: projection matrix of $C(X)$.

$P_x = X (X^T X)^{-1} X^T$: orthogonal projection matrix

Prediction $\hat{y} = P_x y$

Estimated error $e = (I_n - P_x) y$

$$\begin{aligned}\hat{y}^T e &= y^T P_x (I_n - P_x) y \\ &= y^T (P_x - P_x) y \\ &= y^T O y = 0\end{aligned}$$

\hat{y} and e are independently distributed.

$$\text{Cov}(\hat{y}, e) = P_x (\sigma^2 I_n) (I_n - P_x)^T = 0_{n \times n}$$

\hat{y} and $e^T e = y^T (I_n - P_x) y$, are also independently distributed.

Homework: Verify that

$$\tilde{Y}^T (I_n - P_X) (I_n - P_X) \tilde{Y} = S_{yy} - \frac{S_{xy}^2}{S_{xx}}$$

One major advantage of this expression (RHS) is that we can readily estimate the error despite fitting the model.

Exercise: Find the MLE of $\hat{\beta}_0, \hat{\beta}_1, \sigma^2$.

Likelihood fn:

$$L(\beta) = \prod_{i=1}^n \exp \left\{ -\frac{1}{2} \left(\frac{y_i - \beta_0 - \beta_1 x_i}{\sigma} \right)^2 \right\} / \sigma \sqrt{2\pi}$$

$$= \left(\frac{1}{\sigma \sqrt{2\pi}} \right)^n \cdot \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \right\}$$

$$\hat{\beta}_{OLS} = \hat{\beta}_{MLE}; \quad \hat{\sigma}_{OLS} = \hat{\sigma}_{MLE}$$

Homework

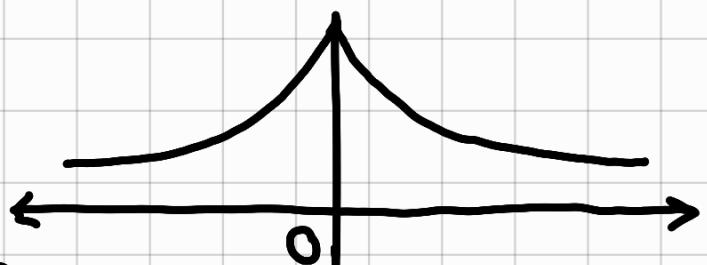
$$\hat{\sigma}_{LS}^2 = \frac{\mathbf{e}^T \mathbf{e}}{n-2} \quad \text{unbiased}$$

but

$$\hat{\sigma}_{MLE}^2 = \frac{\mathbf{e}^T \mathbf{e}}{n} \quad \text{biased}$$

Suppose,

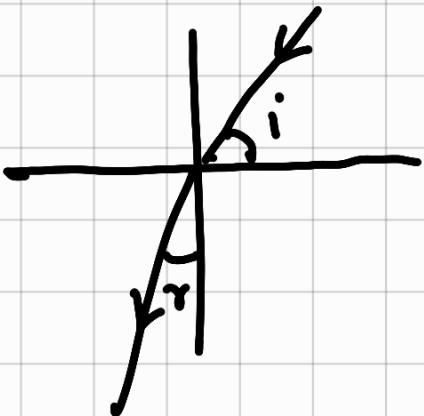
$$\epsilon_i \stackrel{iid}{\sim} f(\epsilon) = \frac{e^{-|\epsilon|/\lambda}}{2\lambda}$$



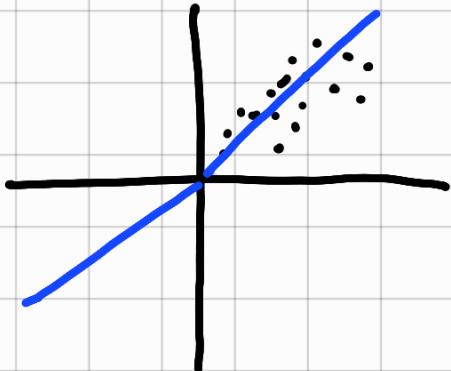
$$y_i \sim g(y) = \frac{1}{2\lambda} \cdot \exp \left\{ -\frac{|y_i - \beta_0 - \beta_1 x_i|}{\lambda} \right\}$$

Obtain the MLE of $\hat{\beta}_0, \hat{\beta}_1$ and λ .

Testing for Simple Linear Regression



$$\frac{\sin i}{\sin \gamma} = e.$$



$$H_0: \beta_0 = 0 (= b_0)$$

$$H_1: \beta_0 \neq 0 (= b_0)$$

$$\hat{\beta}_0 \sim N(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right))$$

Step-1

$$\Rightarrow \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \sim N(0, 1)$$

Step-2

$$\hat{\sigma}^2 = \frac{\mathbf{e}^\top \mathbf{e}}{n-2} = \frac{1}{n-2} \left(S_{yy} - \frac{S_{xy}^2}{S_{xx}} \right)$$

estimation of unknown parameters

Step-3

$$T = \frac{\hat{\beta}_0 - b_0}{\sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \sim t_{n-2, 0}$$

under H_0 .

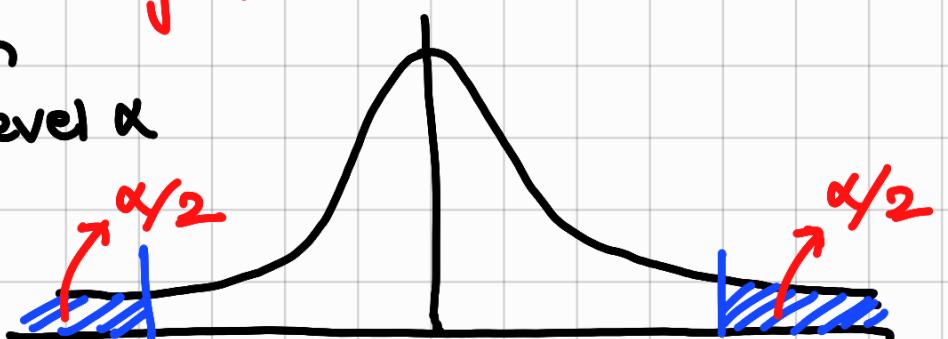
Step-4

Calculate the value T_{observed}

Step-5 Rejection condition

based on

H_1 , and level α



α is based on Type-1 error

$$P_{H_0}(\text{rejecting } H_0) \leq \alpha$$

... (probability of rejecting H_0 when H_0 is true)

$$\alpha = 0.01, 0.05 \text{ or } 0.1$$

$$T_{\alpha/2, n-2, 0}$$

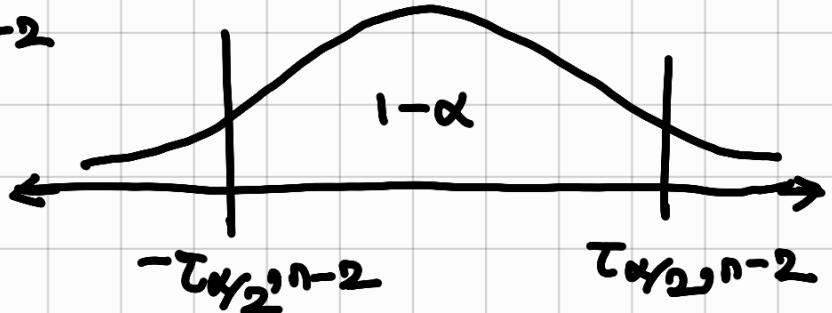
Step 6 Test Rule

We reject H_0 in favour of H_1 , at level α .

$$\text{if } |T_{\text{observed}}| > T_{\alpha/2, n-2}$$

95% confidence interval of β_0

$$\frac{\hat{\beta}_0 - \beta_0}{\sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \sim t_{n-2}$$



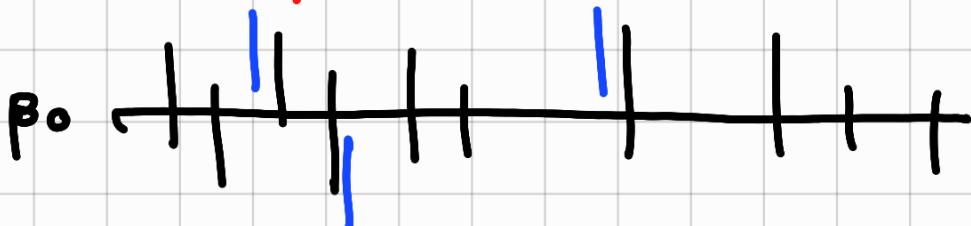
$$P(-T_{\alpha/2, n-2} < T < T_{\alpha/2, n-2}) = 1-\alpha$$

$$P\left(-T_{\alpha/2, n-2} < \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} < T_{\alpha/2, n-2}\right) = 1-\alpha$$

$$L(\underline{x}, \underline{y}) = \hat{\beta}_0 - T_{\alpha/2, n-2} \hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}$$

$$U(\underline{x}, \underline{y}) = \hat{\beta}_0 + T_{\alpha/2, n-2} \hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}$$

Interpretation of confidence interval



Homework Find the 95% CI of σ^2 .

$$D = \{(x_i, y_i) : i = 1, 2, \dots, n\}$$

$\hat{\beta}_0, \hat{\beta}_1$ obtained from data

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$$

$x_0 \notin \{x_1, x_2, \dots, x_n\}$
95% CI of $E(y_0)$.
for the same model

$$E(y_0) = \beta_0 + \beta_1 x_0$$

for parameter, we use term confidence interval
(.... CI of $E(y_0)$)

for variable, we use term prediction interval
(... PI of y_0)

$$y_0 \sim N(\beta_0 + \beta_1 x_0, \sigma^2)$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_0 \sim N(\beta_0 + \beta_1 x_0, \sigma^2(\hat{y}))$$

$$\text{var}(\hat{y}_0) = \text{var}(\hat{\beta}_0 + \hat{\beta}_1 x_0)$$

$$= \text{var}(\bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_0)$$

$$= \text{var}(\bar{y} + \hat{\beta}_1 (x_0 - \bar{x}))$$

$$= \text{var}(\bar{y}) + (x_0 - \bar{x})^2 \text{var}(\hat{\beta}_1) + \text{Cov}(\bar{y}, \hat{\beta}_1)$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \mathbf{1}^T \mathbf{y}$$

$$\hat{\beta}_1 = \frac{s_{xy}}{s_{xx}} = \sum_{i=1}^n \frac{(x_i - \bar{x})y_i}{s_{xx}} = \frac{\mathbf{l}^T \mathbf{y}}{s_{xx}}$$

$$\text{cov}(\bar{y}, \hat{\beta}_1) = \frac{1}{n} \mathbf{1}^T (\sigma^2 I_n) \mathbf{l} = \frac{\sigma^2}{n} \mathbf{1}^T \mathbf{l} = 0$$

$$l_i = (x_i - \bar{x}) / s_{xx}$$

$$y_0 \sim N(\beta_0 + \beta_1 x_0, \sigma^2)$$

$$y_0 = \beta_0 + \hat{\beta}_1 x_0 \sim N(\beta_0 + \beta_1 x_0, \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{s_{xx}} \right))$$

$$\frac{\hat{y}_0 - (\beta_0 + \beta_1 x_0)}{\hat{\sigma} \sqrt{\left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{s_{xx}} \right)}} \sim t_{n-2}$$

95% of CI of $E(y_0)$ is

$$\hat{y} \mp t_{\alpha/2, n-2} \hat{\sigma} \sqrt{\frac{1}{n} + \left(\frac{\bar{x} - x_0}{s_{xx}} \right)^2}$$

95% prediction interval of y_0

$$\text{independent } \begin{cases} y_0 \sim N(\beta_0 + \beta_1 x_0, \sigma^2) \\ y_0 \sim N(\beta_0 + \beta_1 x_0, \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{s_{xx}} \right)) \end{cases}$$

$$y_0 - \hat{y}_0 \sim N(0, \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{s_{xx}} \right))$$

Another case example for testing

$$(x, y) \sim \text{BVN}(\mu_x, \mu_y, \sigma_x^2, \sigma_y^2, \rho)$$

$$H_0: \rho = 0 \text{ v/s } H_1: \rho \neq 0$$

$$T = \frac{\gamma \sqrt{n-2}}{\sqrt{1-\gamma^2}} \sim t_{n-2}$$

$$\gamma = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}} = \frac{s_{xy}}{\sqrt{s_{xx}s_{yy}}}$$

$H_0: p = p_0 \neq 0$ v/s $p \neq p_0$

define $Z = \frac{1}{2} \log_e \left(\frac{1+\gamma}{1-\gamma} \right) = \tanh^{-1}(\gamma)$

$$\mu_Z = \frac{1}{2} \log_e \left(\frac{1+p_0}{1-p_0} \right) = \tanh^{-1}(p_0)$$

$$\hat{\sigma}_Z = (n-3)^{-1}$$

test statistic $W = \frac{Z - \mu_Z}{\hat{\sigma}_Z} \sim N(0, 1)$ as $n \rightarrow \infty$

Multiple linear regression

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \epsilon_i$$

$$\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

unknown parameters : $\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}$ and σ^2

Polynomial regression

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_k x_i^k + \epsilon_i$$

$$\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

Unknown : (β, σ^2)

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i1} x_{i2} + \beta_4 x_{i1}^2 + \beta_5 x_{i2}^2 + \epsilon$$

Unknown : (β, σ^2)

$$\tilde{y} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{pmatrix} \quad \tilde{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix} \quad \tilde{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}$$

$$X = \begin{bmatrix} 1 & c_1 & c_2 & \dots & c_k \end{bmatrix}$$

$$= \begin{bmatrix} 1 & x_{i1} & x_{i2} & \dots & x_{ik} \end{bmatrix} \rightarrow \text{i-th-row}$$

$$\tilde{y} = X \tilde{\beta} + \tilde{\epsilon} \rightarrow n \times 1$$

$$\begin{matrix} n \times 1 & \downarrow & (k+1) \times 1 \\ & \hookrightarrow & \end{matrix}$$

$$n \times (k+1)$$

Least squared condition

$$S(\beta) = (y - X\beta)^T (y - X\beta) \rightarrow \frac{\partial S(\beta)}{\partial \beta} = 0$$

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

If $|X^T X| \neq 0$,

$$\hat{y}_{\text{new}} = \hat{X}^T \hat{\beta} = \hat{x}_{\text{new}}^T \cdot \hat{\beta}$$

$$\hat{y} = \hat{X}^T \hat{\beta} = [\hat{x}(X^T X)^{-1} X^T] y = P_X y$$

$P_X \leftarrow$ orthogonal projection matrix of $C(X)$.

Estimated error vector

$$\tilde{\epsilon} = \tilde{y} - \hat{y} = (I - P_X) y$$

$$\hat{\beta} = (X^T X)^{-1} X^T y \text{ where } \tilde{y} \sim \mathcal{N}(\hat{X}\hat{\beta}, \sigma^2 I_n)$$

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

$\Rightarrow \hat{\beta}$ is an unbiased estimator.

Estimation of σ^2

$$\sigma^2 \left(\frac{e^T e}{\sigma^2} \right) = \left[(I_n - P_X) \frac{y}{\sigma} \right]^T \left[(I_n - P_X) \frac{y}{\sigma} \right]$$

$$= \sigma^2 \left[\frac{y^T}{\sigma} \cdot (I_n - P_X) \frac{y}{\sigma} \right]$$

$$\left[\begin{array}{l} \dots \frac{y}{\sigma} \sim N \left(\frac{x_B}{\sigma}, I_n \right) \\ \dots (I_n - P_X)^T = I_n - P_X \\ (I_n - P_X)^2 = I_n - P_X \end{array} \right]$$

$$\sigma^2 \left(\frac{e^T e}{\sigma^2} \right) \sim \sigma^2 \chi^2$$

$$df = \text{rank}(I - P_X) = n - 1$$

$$ncp = \left(\frac{x_B}{\sigma} \right)^T (I_n - P_X) \left(\frac{x_B}{\sigma} \right) = 0$$

$$E[e^T e] = \sigma^2(n - k + 1)$$

$$\Rightarrow E \left[\frac{e^T e}{n - k + 1} \right] = \sigma^2$$

Unbiased estimator of σ^2 is $\hat{\sigma}^2 = \frac{e^T e}{n - k - 1}$.

$$ncp: \left(\frac{x_B}{\sigma} \right)^T (I_n - P_X) \left(\frac{x_B}{\sigma} \right)$$

$$= \frac{1}{\sigma^2} B^T [x^T (I_n - P_X) x] B$$

$$= \frac{1}{\sigma^2} B^T [x^T x - x^T P_X x] B = \frac{1}{\sigma^2} B^T [x^T x - x^T x] B$$

Homework MLE of $\sigma^2 = \mathbf{e}^T \mathbf{e} / n$

Hence, it is not an unbiased estimator.

Testing the significance of feature(s)

1. $H_0: \beta_i = 0 (= b_i)$ Testing of linear combination of β_i 's, $i=1, \dots$
 $H_1: \beta_i \neq 0 (\neq b_i)$

2. $H_0: \beta_1 + 2\beta_2 = b_0 \iff H_0: \mathbf{l}^T \boldsymbol{\beta} = b_0$
 $H_1: \beta_1 + 2\beta_2 \neq b_0 \quad H_1: \mathbf{l}^T \boldsymbol{\beta} \neq b_0$

3. $H_0: \begin{cases} \mathbf{l}_1^T \boldsymbol{\beta} = b_1 \\ \mathbf{l}_2^T \boldsymbol{\beta} = b_2 \end{cases}$
 $H_1: H_0 \text{ is not true.}$

How to combine both the hypothesis?

Develop a methodology to simultaneously test both of them together

4. $H_0: \beta_1 = \beta_2 = \dots = \beta_K = 0$

$H_1: H_0 \text{ is not true}$

|||

More general case of above version

$H_0: \begin{cases} \beta_1 = 0 \\ \beta_1 - \beta_2 = 0 \\ \beta_1 - \beta_3 = 0 \\ \vdots \\ \beta_1 - \beta_K = 0 \end{cases}$

$H_1: H_0 \text{ is not true}$

→ This testing signifies if it is at all necessary to build a model this way or not.

→ Analysis of variance.

ANOVA regression model

Assumption: $\text{abs} |(\mathbf{X}^T \mathbf{X})| >> 0$

If $|X^T X| = 0$, then any $\hat{\beta}^T \beta$ may not be estimable
i.e. there will be no linear unbiased estimator of $\hat{\beta}^T \beta$.

i.e., $\nexists b^T y$ such that $E[b^T y] = \hat{\beta}^T \beta$.

We have

$$\begin{aligned} Y &\sim N(X\beta, \sigma^2 I_n) \\ \hat{\beta} &\sim N(\beta, \sigma^2 (X^T X)^{-1}) \\ \hat{\sigma}^2 &= \frac{e^T e}{n-k-1} \end{aligned}$$

$$(X^T X)^{-1} = C_{(k+1) \times (k+1)} = ((c_{ij}))$$

$$i, j = 0, 1, 2, \dots, k$$

Q1. $H_0: \beta_i = b_i$

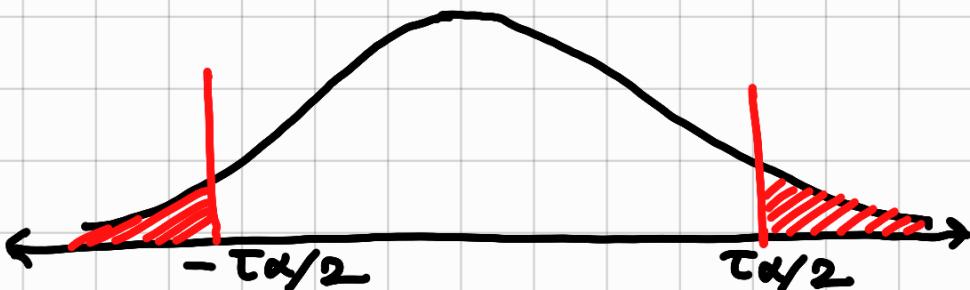
$$H_1: \beta_i \neq b_i$$

$$\begin{aligned} \hat{\beta}_i &\sim N(\beta_i, \sigma^2 c_{ii}) \\ \Rightarrow \frac{\hat{\beta}_i - \beta_i}{\sigma \sqrt{c_{ii}}} &\sim N(0, 1) \end{aligned}$$

Under H_0 ,

$$\begin{aligned} \frac{\hat{\beta}_i - b_i}{\sigma \sqrt{c_{ii}}} &\sim N(0, 1) \\ \Rightarrow T = \frac{\hat{\beta}_i - b_i}{\hat{\sigma} \sqrt{c_{ii}}} &\sim t_{n-k-1} \end{aligned}$$

Reject H_0 if
 $|T_{obs}| > T_{\alpha/2, n-k-1}$



$$Q2 \quad H_0: a\beta_i + b\beta_j = d$$

$$H_1: a\beta_i + b\beta_j \neq d$$

$$a\hat{\beta}_i + b\hat{\beta}_j \sim N(a\beta_i + b\beta_j, (a^2 c_{ii} + b^2 c_{jj} + 2ab c_{ij})\sigma^2)$$

Under H_0 :

$$\frac{(a\hat{\beta}_i + b\hat{\beta}_j - d)}{\sqrt{a^2 c_{ii} + b^2 c_{jj} + 2ab c_{ij}}} \sim N(0, 1)$$

$$\rightarrow T = \frac{a\hat{\beta}_i + b\hat{\beta}_j - d}{\sqrt{a^2 c_{ii} + b^2 c_{jj} + 2ab c_{ij}}} \stackrel{H_0}{\sim} t_{n-k-1}$$

We reject H_0 in favor of H_1 at level α if
 $|T_{\text{obs}}| > t_{\alpha/2, n-k-1}$

$$Q3. \quad H_0: \begin{cases} 2\beta_1 + 3\beta_2 = 1 \\ 4\beta_3 - 5\beta_6 = 7 \end{cases}$$

$H_1: H_0$ is not true

$$\begin{bmatrix} 0 & 2 & 3 & 0 & \dots & 0 \end{bmatrix} B = \begin{bmatrix} 1 \\ 7 \end{bmatrix}$$

$$AB = b$$

$$\hat{\beta} \sim N(\beta, C\sigma^2)$$

$$A\hat{\beta} \sim N(AB, ACAT\sigma^2)$$

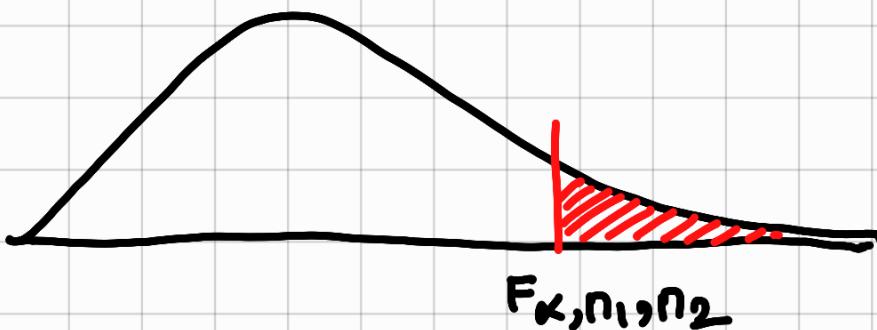
(... we can no longer make it into a t-statistic since this is a vector)

$$(A\hat{\beta} - b)^T \left[\frac{(ACAT)^{-1}}{\sigma^2} \right] (A\hat{\beta} - b) \stackrel{H_0}{\sim} \chi^2_{\text{rank}(ACAT)}$$

Note: $Z \sim N(\mu, \Sigma)$

$(Z - \mu)^T \Sigma^{-1} (Z - \mu)$ follows $\chi^2_{\text{rank}(\Sigma)}$

$$(\hat{A}\hat{\beta} - b)^T \frac{(AC\Delta^T)^{-1}}{\hat{\sigma}^2} (\hat{A}\hat{\beta} - b) \sim F_{\text{rank}(AC\Delta^T), n-k-1}$$



Reject H_0 if
 $F_{\text{obs}} > F_{k,n_1,n_2}$

Q4. ANOVA

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$H_1: H_0$ is not true

→ "do we even need to fit a model this way or not"

$$SST_{\text{Total}} = \sum_{i=1}^n (y_i - \bar{y})^2 = \mathbf{y}^T (I_n - \frac{1}{n} \mathbf{1} \mathbf{1}^T) \mathbf{y}$$

$$SSE_{\text{Error}} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \mathbf{y}^T (I_n - P_x) \mathbf{y}$$

$$SS_{\text{Regression}} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \mathbf{y}^T (P_x - \frac{1}{n} \mathbf{1} \mathbf{1}^T) \mathbf{y}$$

$$I_n - \frac{1}{n} \mathbf{1} \mathbf{1}^T = (I_n - P_x) + (P_x - \frac{1}{n} \mathbf{1} \mathbf{1}^T)$$

$$SST_{\text{Total}} = SSE_{\text{Error}} + SS_{\text{Regression}}$$

$$SSE_{\text{Error}} / \hat{\sigma}^2 \sim \chi^2_{n-k-1, n-1}$$

$$SS_{\text{Reg.}} / \hat{\sigma}^2 \sim \chi^2_{k, n-1}$$

... Cochran's thm.

$$SSTotal/\sigma^2 \sim \chi^2_{n-1, ncp=\lambda}$$

$$F = \frac{(SSReg/\sigma^2)/k}{(SSEerror/\sigma^2)/(n-k-1)} \sim F_{(k, \lambda), (n-k-1, \beta)}$$

If H_0 is true, $\lambda = 0$.
otherwise, $\lambda > 0$.



Exercise. Try ANOVA in the case of simple linear regression and the line passing through origin.

Rough Outline for the Proof of $\textcircled{*}$

$$\begin{aligned} \text{NCP of } (SSReg/\sigma^2) &= \left(\frac{x\beta}{\sigma}\right)^T \left(P_x - \frac{1}{n} \mathbf{1} \mathbf{1}^T\right) \left(\frac{x\beta}{\sigma}\right) \\ &= \frac{1}{\sigma^2} \beta^T \left[x^T \left(P_x - \frac{1}{n} \mathbf{1} \mathbf{1}^T\right) x \right] \beta \\ &= \frac{1}{\sigma^2} \beta^T \left[x^T x - \frac{1}{n} x^T \mathbf{1} \mathbf{1}^T x \right] \beta \end{aligned}$$

$$\beta = \begin{pmatrix} \beta_0 \\ \beta_R \end{pmatrix} \quad x = \begin{bmatrix} 1 : X_R \end{bmatrix}$$

Hence,

$$\text{NCP} = \frac{1}{\sigma^2} \beta^T \left[\left(\begin{array}{c|c} n & \mathbf{1}^T x_R \\ \hline x_R^T \mathbf{1} & x_R^T x_R \end{array} \right) - \left(\begin{array}{c|c} n & \mathbf{1}^T x_R \\ \hline x_R^T \mathbf{1} & \frac{1}{n} x_R^T \mathbf{1} \mathbf{1}^T x_R \end{array} \right) \right] \beta$$

$$= \frac{1}{\sigma^2} \begin{pmatrix} \beta_0 \\ \beta_R \end{pmatrix}^T \begin{pmatrix} 0 & \tilde{\Omega}^T \\ \tilde{\Omega} & A \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_R \end{pmatrix} = \frac{1}{\sigma^2} \underbrace{\beta_R^T A \beta_R}_{\text{not a function of } \beta_0}$$

$$A = X_R^T X_R - \frac{1}{n} X_R^T \underbrace{1}_{\sim} \underbrace{1^T}_{\sim} X_R,$$

Variance - covariance matrix of $\tilde{x}_1 \tilde{x}_2 \dots \tilde{x}_k$

What if $(X^T X)$ is not invertible? The rest of the part deals with this question.

Some definitions and theorems:

Linear Unbiased Estimator

$l^T y = \sum_i l_i y_i$ is said to be LUE of $\beta^T \beta$ if $E[l^T y] = \beta^T \beta \quad \forall \beta \in \mathbb{R}^{k+1}$.

Linear Zero function

$l^T y = \sum_i l_i y_i$ is said to be LZF if $E(l^T y) = 0$

- $\beta^T \beta$ is **estimable** if it has a LUE
- Best LUE (BLUE) is LUE $\hat{\beta}^T \beta$ with min. variance.

Theorem :

A linear function is BLUE of its expectation iff it is uncorrelated with all LZF.

$\hat{\beta}^T \beta$ is estimable iff $\exists l^T y$ such that $E(l^T y) = \hat{\beta}^T \beta \quad \forall \beta \in \mathbb{R}^{k+1}$

$$\Leftrightarrow l^T X \beta = \hat{\beta}^T \beta \quad \forall \beta \in \mathbb{R}^{k+1}$$

$$\boxed{l^T X = \hat{\beta}^T} \quad \boxed{\hat{\beta} \in C(C(X^T))}$$

$$\hat{\beta} = X^T l$$

$$\begin{aligned}
 \text{LZF} \quad & E(\tilde{l}^T y) = 0 \quad \forall \beta \in \mathbb{R}^{k+1} \\
 \Rightarrow & \tilde{l}^T \times \beta = 0 \quad \forall \beta \in \mathbb{R}^{k+1} \\
 \Rightarrow & \boxed{\tilde{l}^T x = 0^T} \iff \tilde{l} \in (C(x))^\perp \\
 & \Rightarrow \tilde{l} \in C(I_n - P_x)
 \end{aligned}$$

$$\begin{aligned}
 \tilde{l}^T x &= m^T (I_n - P_x) x \quad (l = (I_n - P_x)m) \\
 E[\tilde{l}^T y] &= \tilde{m}^T (I_n - P_x) x \beta \\
 &= \tilde{m}^T \beta
 \end{aligned}$$

Theorem :

A linear function is BLUE of its expectation iff it is uncorrelated with all LZF.

$$\text{given } E(\tilde{l}^T x) = b^T \beta \Rightarrow \tilde{l}^T x \beta = b^T \beta$$

$$\begin{aligned}
 \tilde{l}^T x &= \tilde{l}^T I_n y = \tilde{l}^T (P_x + I_n - P_x) y \\
 &= (\tilde{l}^T P_x y) + \underbrace{\tilde{l}^T (I_n - P_x) y}_{\text{LZF}}
 \end{aligned}$$

$$E(\tilde{l}^T P_x y) = l^T P_x x \beta = \tilde{l}^T x \beta = b^T \beta$$

$$\begin{aligned}
 \text{cov}(\tilde{l}^T P_x y, \tilde{l}^T (I_n - P_x) y) \\
 &= \sigma^2 (l^T P_x (I_n - P_x) l) \\
 &= \sigma^2 (\tilde{l}^T O \tilde{l}) \\
 &= 0.
 \end{aligned}$$

BLUE($\tilde{l}^T P_x y$)

See the prof's slides for recommended books.

Coefficient of determination

a) $R^2 = \frac{\text{Variation in } y \text{ explained}}{\text{Total variation in } y}$

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

$$= \frac{\text{SS Model}}{\text{SS Total}} = 1 - \frac{\text{SSE}_{\text{Error}}}{\text{SS Total}}$$

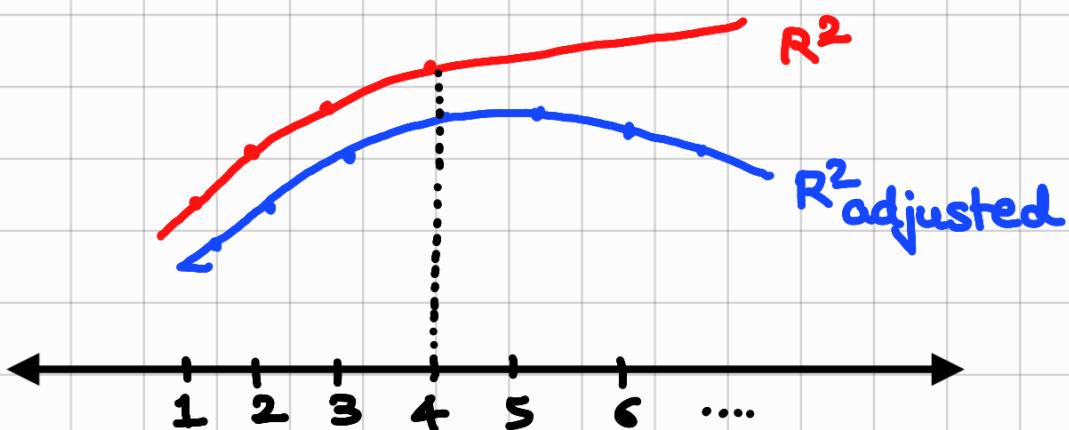
$$= 1 - \frac{\sum_{i=1}^n e_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

- $\rightarrow R^2 \in (0, 1)$
 \rightarrow increasing fn. of no. of variables

b) Adjusted R^2 :

$$R^2_{\text{adjusted}} = 1 - \frac{\sum_{i=1}^n e_i^2 / df}{\sum_{i=1}^n (y_i - \bar{y})^2 / df}$$

$$= 1 - \frac{\sum_{i=1}^n e_i^2 / (n-k-1)}{\sum_{i=1}^n (y_i - \bar{y})^2 / n-1} < R^2$$



\rightarrow find the maxima using local search heuristics.
 { even the courses of
 both the facads go
 hand-in-hand :)}

Polynomial Regression

→ analogous to multiple linear regression

$$y = X\beta + \epsilon$$

$$X = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)^T \text{ with } \tilde{x}_j = (1, x_{j1}, x_{j2}^2, \dots, x_{jk})^T$$

$$\beta = (\beta_0, \beta_1, \dots, \beta_k)^T$$

$k+2$ unknown model parameters: β and $\sigma^2 \gamma_0$.

$$Y \sim \mathcal{N}(X\beta, \sigma^2 I_n)$$

Problems of polynomial regression

→ Finding order of the model

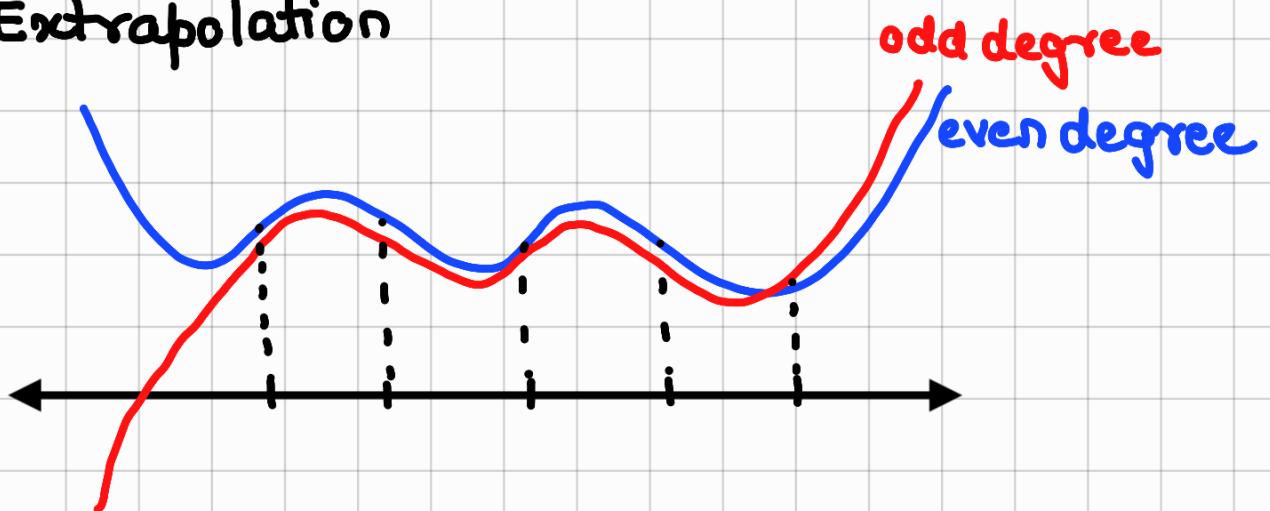
- forward selection

→ continue ↑ing the order and computed adjusted R^2 values.

→ find $\operatorname{argmax}(R^2_{\text{adj}})$

- backward elimination

→ Extrapolation



→ Ill-conditioning

$X^T X$ may be computationally singular

→ Hierarchy

not necessary that all degrees till k ($1, 2, \dots, k$) should be present to fit the model.

regression model won't do so.

Orthogonal Polynomial

$$u^T v = \sum_{i=1}^n u_i v_i = 0 \quad \text{orthogonal vectors}$$

$$\langle u(t), v(t) \rangle = 0$$

orthogonal functions

$$\text{e.g. } \int_0^1 u(t)v(t)dt = 0$$

given data $\{(x_i, y_i), i = 1, 2, \dots, n\}$

$$P_0(x_i) = 1 \quad \forall i = 1, \dots, n$$

$$\sum_{i=1}^n P_j(x_i) P_k(x_i) = P_j^T P_k = 0 \text{ for all } j \neq k.$$

$$\tilde{P}_j = \begin{pmatrix} P_j(x_1) \\ P_j(x_2) \\ \vdots \\ P_j(x_n) \end{pmatrix}$$

$$\tilde{P}_k = \begin{pmatrix} P_k(x_1) \\ P_k(x_2) \\ \vdots \\ P_k(x_n) \end{pmatrix}$$

$P_j(x) = \sum_{m=0}^j \theta_m x^m$ is the j th degree orthogonal polynomial

To fit a k -degree orthpoly to the data

$$y_i = \sum_{j=0}^k \alpha_j (P_j(x_i)) + \epsilon_i$$

$$\epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$

Unknown $\{(\alpha_0, \alpha_1, \dots, \alpha_k), \sigma^2\}$

$$X_0 = \begin{pmatrix} P_0(x_1) & P_1(x_1) & \cdots & P_k(x_1) \\ \vdots & \vdots & & \vdots \\ P_0(x_n) & P_1(x_n) & \cdots & P_k(x_n) \end{pmatrix}$$

Hence,

$$y = X_0 \alpha + \epsilon \quad \epsilon \sim N(0, \sigma^2 I_n) \Rightarrow y \sim N(X_0 \alpha, \sigma^2 I_n)$$

$$\hat{\alpha} = (x_0^T x_0)^{-1} x_0^T y$$

$$(x_0^T x_0) = \begin{bmatrix} \sum_{i=1}^n p_0(x_i)^2 & & & \\ & \sum_{i=1}^n p_1(x_i)^2 & & 0 \\ & & \ddots & \\ 0 & & & \sum_{i=1}^n p_k(x_i)^2 \end{bmatrix}$$

$$\hat{\alpha}_j = \frac{\sum_{i=1}^n p_j(x_i)y_i}{\sum_{i=1}^n p_j^2(x_i)}$$

$\hat{\alpha}_j$ is a linear estimator of $\hat{\alpha}_j = P_j^T y / P_j^T P_j = l^T y$

$$l_i = p_j(x_i) / (P_j^T P_j)$$

$$\begin{aligned} \hat{\alpha}_j &\sim N\left(\frac{P_j^T [x_0 \alpha]}{P_j^T P_j}, \sigma^2 \cdot \frac{(P_j^T I_n P_j)}{(P_j^T P_j)^2}\right) \\ &\equiv N\left(\alpha_j, \frac{\sigma^2}{P_j^T P_j}\right) \end{aligned}$$

SSE error

$$\begin{aligned} SSE &= \sum (I_n - P_{X_0}) y \\ &= y^T y - y^T P_{X_0} y \\ &= \sum_{i=1}^n y_i^2 - y^T x_0 \hat{\alpha} \end{aligned}$$

$$\begin{aligned} &\dots (P_{X_0} y = x_0 (x_0^T x_0)^{-1} x_0^T y \\ &= x_0 \hat{\alpha} \\ &= \sum_{i=1}^n y_i^2 - \sum_{j=0}^k \hat{\alpha}_j (\sum_{i=1}^n p_j(x_i) y_i) \end{aligned}$$

$$= \sum_{i=1}^n y_i^2 - (\hat{\alpha}_0 \sum_{i=1}^n p_0(x_i) y_i) - \sum_{j=1}^k \hat{\alpha}_j (\sum_{i=1}^n p_j(x_i) y_i)$$

(... $\sum_{i=1}^n p_0(x_i) y_i = n\bar{y}$)

$$\begin{aligned} SSE &= \sum_{i=1}^n y_i^2 - n\bar{y}^2 - \sum_{j=1}^k \hat{\alpha}_j (\sum_{i=1}^n p_j(x_i) y_i) \\ &= \sum_{i=1}^n (y_i - \bar{y})^2 - \sum_{j=1}^k \hat{\alpha}_j (\sum_{i=1}^n p_j(x_i) y_i) \end{aligned}$$

$$SSE = SST_{\text{Total}} - SS_{\text{Model}}$$

Multicollinearity

$$\tilde{y} = X\beta + \varepsilon ; \beta \in \mathbb{R}^{k+1}$$

$$\hat{\beta} = (X^\top X)^{-1} X^\top \tilde{y} .$$

If $|X^\top X| = 0$ or $|X^\top X| \approx 0 \leftarrow$ ill-condition

$$\hat{\beta} = \frac{\text{Adj}(X^\top X) X^\top \tilde{y}}{|X^\top X|}$$

$\Rightarrow \text{Var}(\hat{\beta}_j)$ may be unbounded.

$\hat{y}_0 = (1 \cdot \tilde{x}_0^\top) \hat{\beta}$ will also have large variance.

Why should this problem arise at all?

→ one regressor is a linear combination of the rest.

→ improper collection of the data may lead to multicollinearity.

→ Model building depending on the magnitude

of the regressor.

Consider the expected sq. error for $\hat{\beta}$.

$$\begin{aligned}
 & E\{(\hat{\beta} - \beta)^T (\hat{\beta} - \beta)\} \\
 &= E \sum_{j=0}^k (\hat{\beta}_j - \beta_j)^2 \\
 &= \sum_{j=0}^k E (\hat{\beta}_j - \beta_j)^2 = \sum_{j=0}^k \text{Var}(\hat{\beta}_j) \\
 &= \text{tr}\{\sigma^2 (x^T x)^{-1}\} \\
 & \quad \left[\dots \hat{\beta} \sim N(\beta, \sigma^2 (x^T x)^{-1}) \right]
 \end{aligned}$$

Note.

- 1) $x^T x$ is a symmetric matrix.
- 2) $(x^T x)^{-1}$ is also a symmetric matrix, if exists.
- 3) $x^T x$ is positive semidefinite matrix.
- 4) $x^T x$ has non-negative evs.

$$\lambda_0 > \lambda_1 > \lambda_2 > \dots \geq 0.$$

$$5) \lambda_{\max} = \max_z \frac{z^T (x^T x) z}{\|z\|^2 \neq 0}$$

$$\lambda_{\min} = \min_{\|z\|^2 \neq 0} \frac{z^T (x^T x) z}{z^T z}$$

$$6) \text{tr}(x^T x) = \sum \lambda_i$$

$$7) \text{tr}((x^T x)^{-1}) = \sum \lambda_i^{-1}$$

$$8) |x^T x| = \prod \lambda_i$$

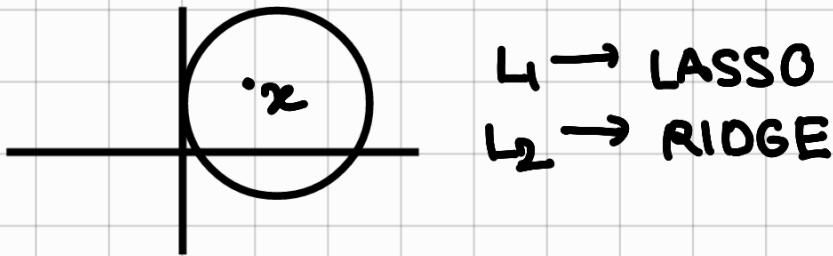
$$9) |(x^T x)^{-1}| = \prod \lambda_i^{-1} \uparrow \infty \text{ as some } \lambda_i \downarrow 0$$

Hence,

$$\begin{aligned}
 E\{(\hat{\beta} - \beta)^T (\hat{\beta} - \beta)\} &= \text{tr}\{\sigma^2 (x^T x)^{-1}\} \\
 &= \sigma^2 \sum_i (\lambda_i^{-1}) \\
 \hat{\beta} &\sim N(\beta, \sigma^2 (x^T x)^{-1})
 \end{aligned}$$

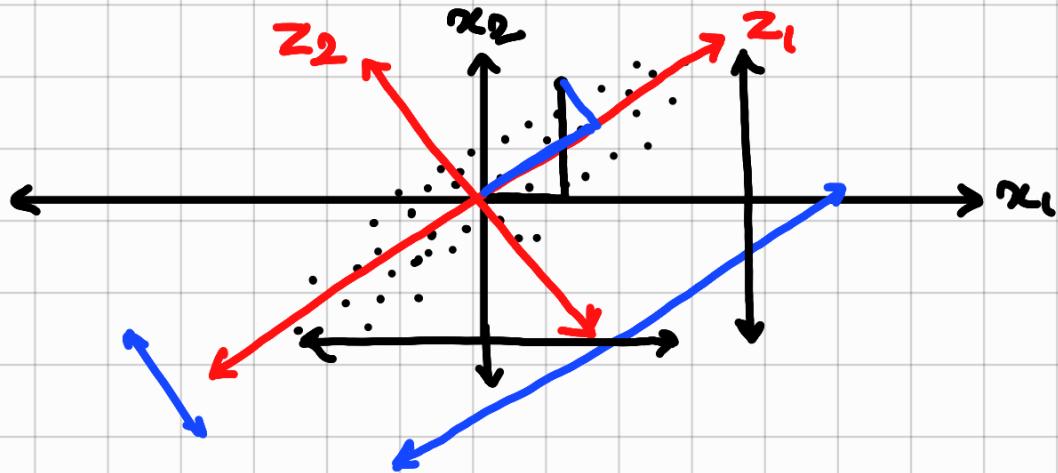
Possible Solutions

- Principal Component Regression
- restrict variance in euclidean norm / ℓ^2 -norm



- Principal Component Regression

SVD / PCA
Math Statistics



- one method of dimension reduction

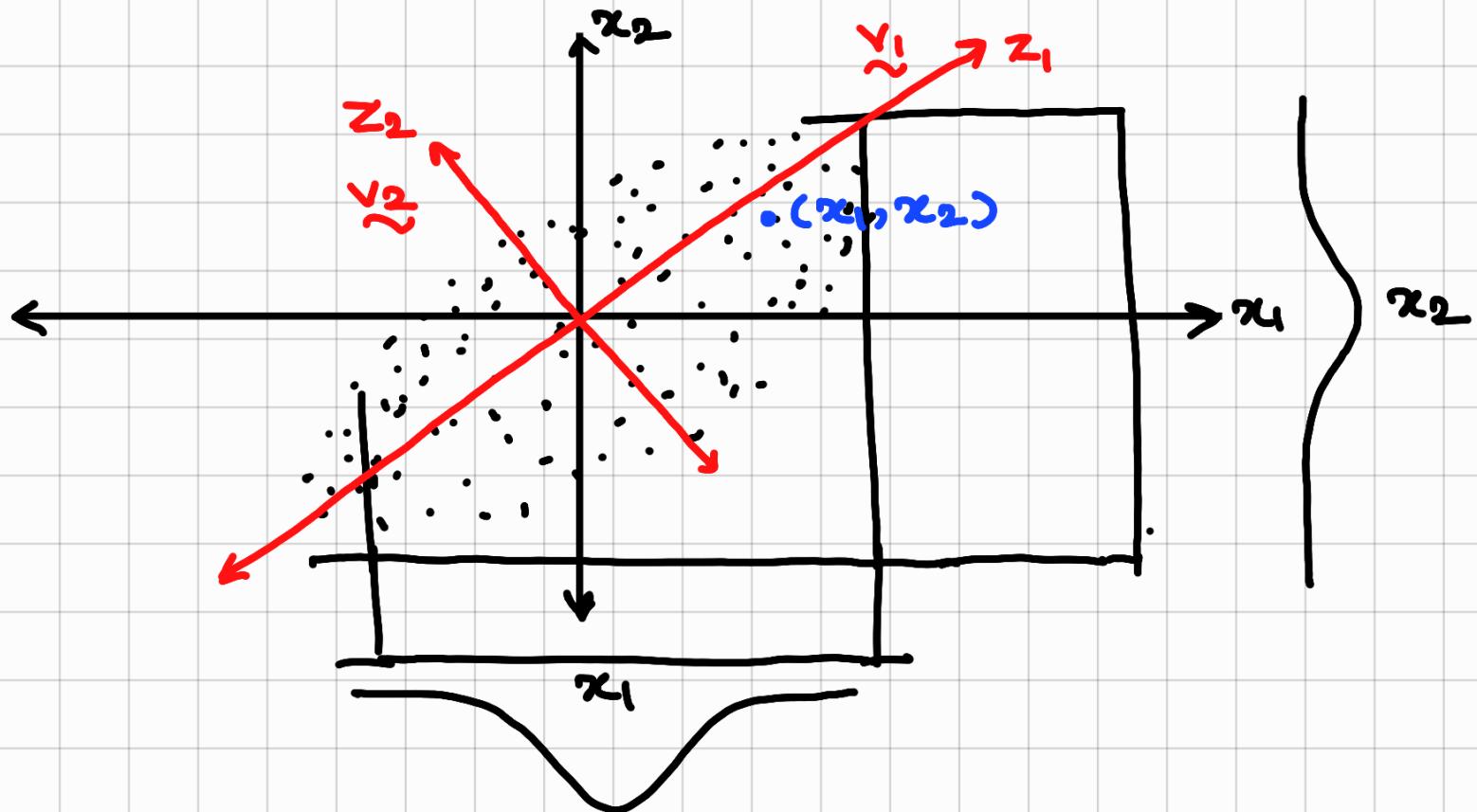
Tentatively, Midsem syllabus is upto and including Polynomial Regression.

Principal Component Regression

$$|x^T x| \approx 0$$

$$x^T x \text{ p.s.d}$$

$$\lambda_0 > \lambda_1 > \lambda_2 \dots \lambda_k \geq 0$$



$$e_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, e_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = x_1 \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} + x_2 \cdot \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$\langle \tilde{x}, e_1 \rangle \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \langle \tilde{x}_2, e_2 \rangle \cdot \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \tilde{x}$$

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \langle x, \tilde{v}_1 \rangle \tilde{v}_1 + \langle x, \tilde{v}_2 \rangle \tilde{v}_2$$

$$= z_1 \tilde{v}_1 + z_2 \tilde{v}_2$$

$$\lambda_0 > \lambda_1 > \lambda_2 \dots \dots > \lambda_k \geq 0$$

$$\begin{bmatrix} \tilde{v}_0 & \tilde{v}_1 & \tilde{v}_2 & \dots & \tilde{v}_k \end{bmatrix} = P$$

$$X^T X \rightarrow (\lambda_i, v_j)_{i=0,1,2,\dots,k}$$

$$X^T X = P D P^T$$

\downarrow \hookrightarrow orthogonal matrix
 diagonal matrix

$$P^T P = P P^T = I.$$

Example of orthogonal matrix : $(\cos \theta \ -\sin \theta)$
 orth. matrix for reflection $\rightarrow (\sin \theta \ \cos \theta)$

$$Y = X\beta + \epsilon$$

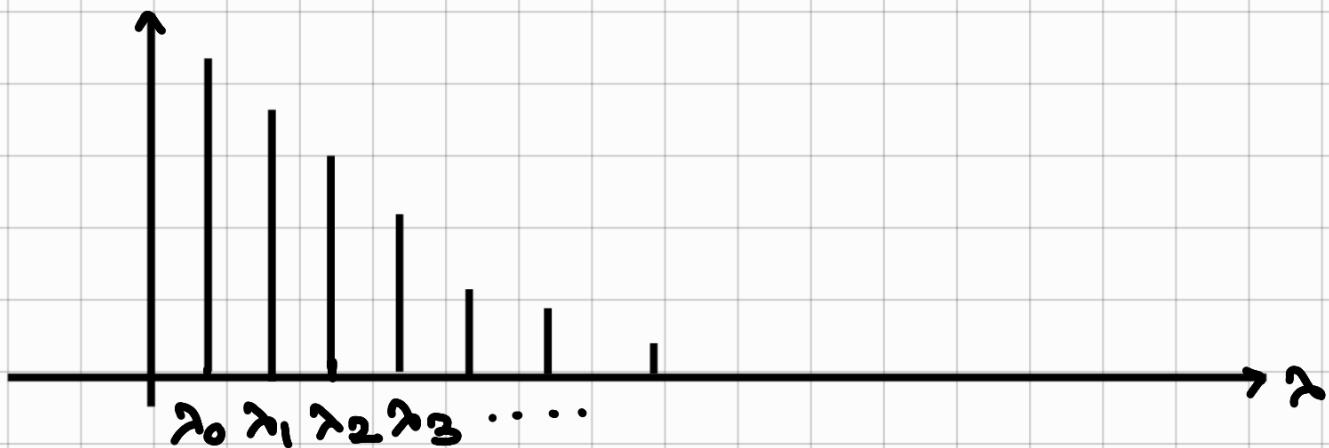
$$Y = X P P^T \tilde{B} + \tilde{\epsilon} \rightarrow \tilde{Y} = Z \tilde{\alpha} + \tilde{\epsilon}$$

$$\Rightarrow \hat{\tilde{\alpha}} = (Z^T Z)^{-1} Z^T \tilde{Y}$$

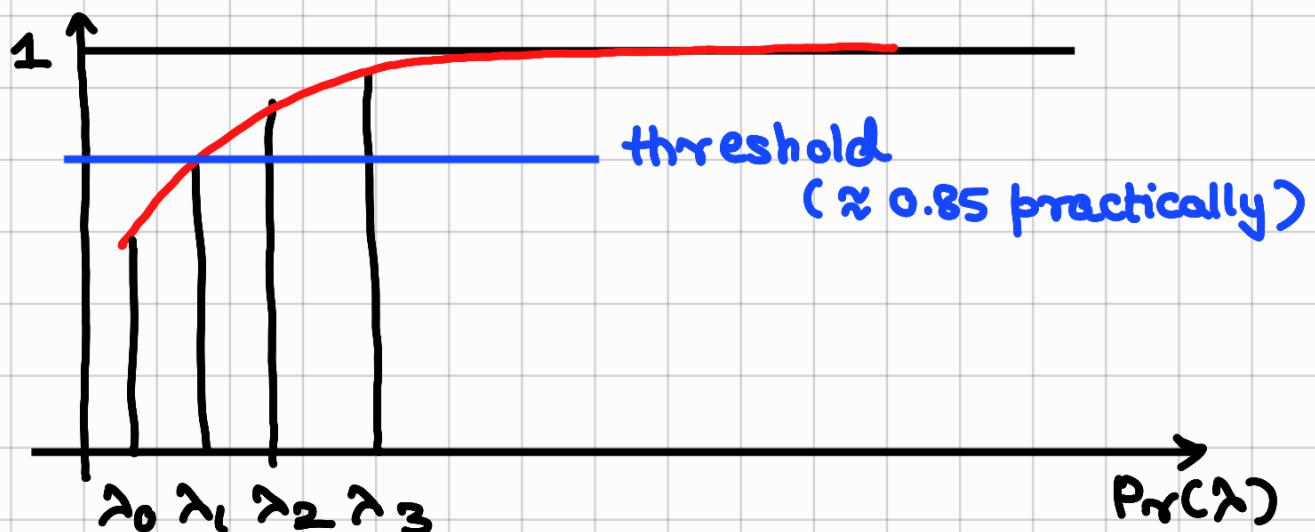
As a process,

$$\tilde{\alpha}_{(r)} = \begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_r \end{pmatrix}$$

(reducing the dimensions)



$$P_{\gamma}(\lambda) = \sum_{i=0}^t \lambda_i / \sum_{i=0}^k \lambda_i$$



$Z_{(r)}$ is the generated matrix from X and

$$[v_0 \ v_1 \ \dots \ v_r]$$

$$Z_{(r)} = X \underset{n \times (k+1)}{[v_0 \ v_1 \ \dots \ v_r]}$$

$$\hat{\alpha}_{(r)} = (Z_{(r)}^T Z_{(r)})^{-1} Z_{(r)}^T \mathbf{y}$$

Transform $\hat{\alpha}_{(r)}$ to β ,

$$\begin{aligned} \beta &= I\beta \\ \Rightarrow \tilde{\beta} &= P P^T \beta \\ \Rightarrow \tilde{\beta} &= P \cdot \beta \end{aligned}$$

$$\Rightarrow \hat{\beta}_{PC}^{(k+1)} = P \cdot \begin{pmatrix} \hat{\alpha}_{(r)} \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

Tentatively, Midsem syllabus is upto and including Polynomial Regression.

$$|x^T x| \approx 0$$

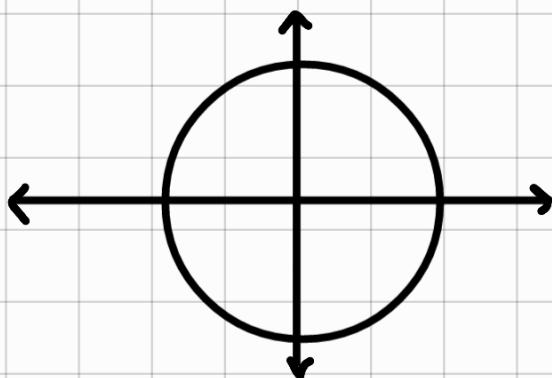
\rightarrow PCR

\rightarrow Shrinkage estimation or regularization

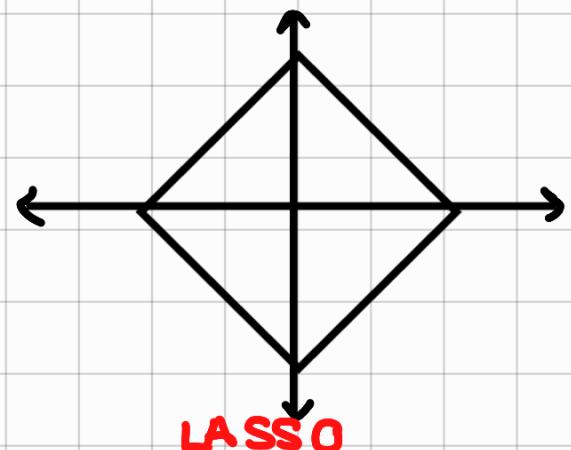
Shrinkage estimation / regularization

As $|x^T x| \approx 0$, the variance of $\hat{\beta}_{LS}$ is large. Hence, we want some bound on the norm of $\hat{\beta}$.

$$\begin{aligned} L_2 : \|\hat{\beta}\|_2^2 &\leq c_2 \quad \text{or} \quad L_1 : \|\hat{\beta}\|_1 \leq c_1 \\ \Leftrightarrow \sum_{i=0}^k |\beta_i|^2 &\leq c_2 \quad \Leftrightarrow \sum_{i=0}^k |\beta_i| < c_1 \end{aligned}$$



RIDGE



LASSO

Minimization of LS condition

$$S(\beta) = (\underline{y} - \underline{x}\beta)^T (\underline{y} - \underline{x}\beta) \text{ w.r.t. } \sum_{i=0}^k |\beta_i|^2 < c_2$$

Using Lagrangian multiplier

$$\hat{\beta}_R = \underset{\beta}{\operatorname{argmin}} (S(\beta) + \lambda(\beta^T \beta - c_2))$$

$$\hat{\beta}_R = (x^T x + \lambda I_{k+1})^{-1} x^T \underline{y}$$

Hoerl, Kennard, Baldwin (1975)

$$\lambda = (k+1) \hat{\sigma}_{LS} / \hat{\beta}_{LS}^T \hat{\beta}_{LS}$$

These $\hat{\sigma}_{LS}$, $\hat{\beta}_{LS}$ are generated from original model with/without using PCR depending on $(x^T x)^{-1}$.

Note :

$$1) \lambda \rightarrow \infty \Rightarrow \hat{\beta}_R = \underline{0}$$

$$2) \lambda \rightarrow 0 \Rightarrow \hat{\beta}_R = \hat{\beta}_{LS}$$

Is $\hat{\beta}_R$ an unbiased estimator of β ?

$$\begin{aligned} \hat{\beta}_R &= (x^T x + \lambda I_{k+1})^{-1} (x^T x) (x^T x)^{-1} x^T \underline{y} \\ &= \underbrace{[(x^T x + \lambda I_{k+1})^{-1} (x^T x)]}_{\text{not an identity matrix}} \hat{\beta}_{LS} \end{aligned}$$

$\hat{\beta}_R$ is not an unbiased estimator.

Notations :

$$W = (x^T x + \lambda I)^{-1}$$

$$S = (x^T x)$$

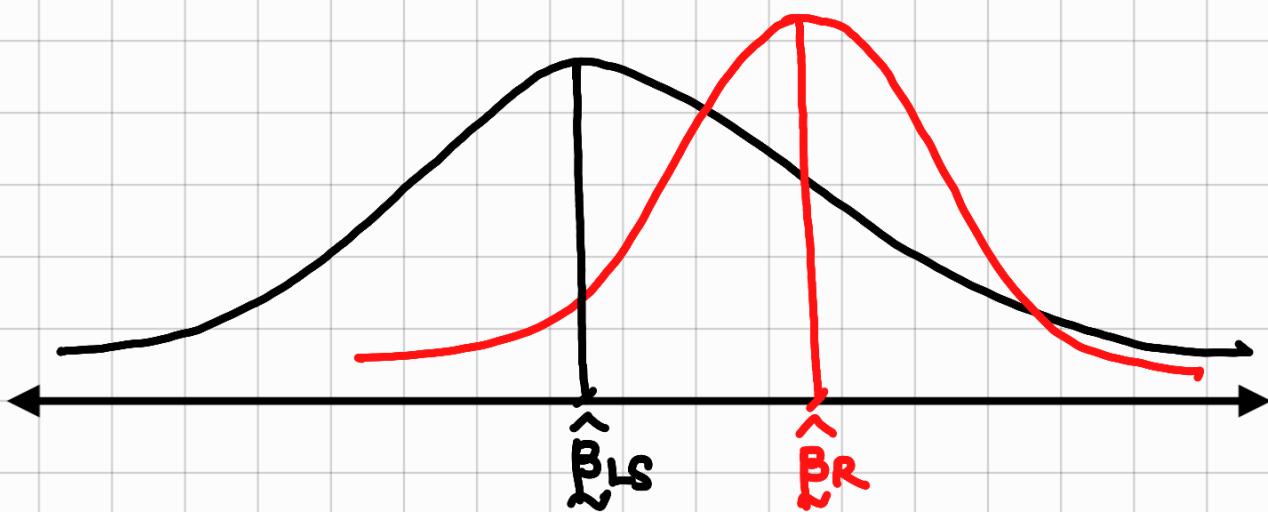
$$\therefore \hat{\beta}_R = W S \hat{\beta}_{LS}$$

Does RIDGE estimate have less variance?

$$\begin{aligned}
 D(\hat{\beta}_{LS}) - D(\hat{\beta}_R) &= \sigma^2(S^{-1} - WSS^{-1}SW) \\
 &= \sigma^2(S^{-1} - WSW) \\
 &= \sigma^2 W(W^T S^{-1} W - S)W \\
 &= \sigma^2 W((S + \lambda I)S^{-1}(S + \lambda I) - S)W \\
 &= \sigma^2 W(2\lambda I + \lambda^2 S^{-1})W
 \end{aligned}$$

↓ psd ↓ psd
 psd ↓ psd

Hence, $D(\hat{\beta}_{LS}) - D(\hat{\beta}_R)$ is a psd matrix.
This implies RIDGE estimate has less variance.



$$\text{Bias } E(\hat{\beta}_R) - \hat{\beta} = (WS - I)\hat{\beta}$$

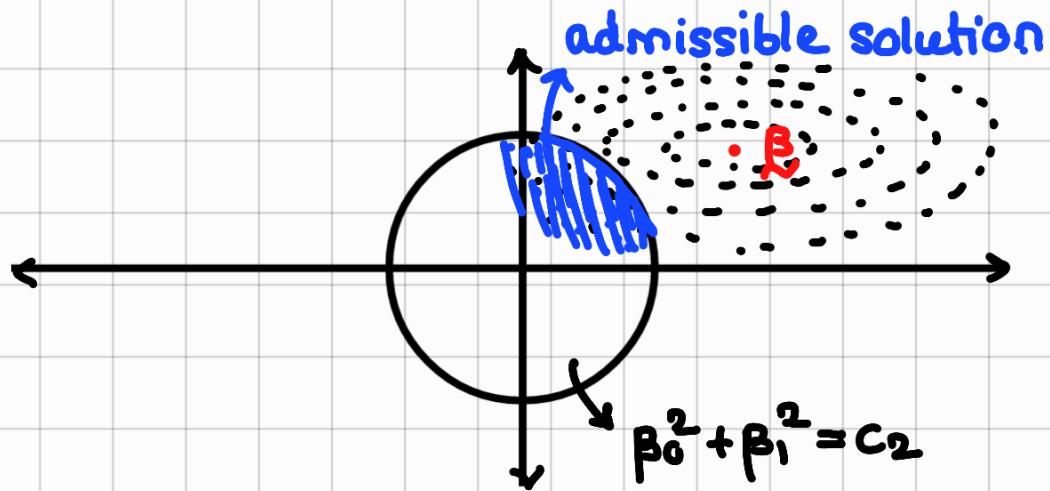
$$\text{MSE}(\hat{\beta}_R) = \text{tr}[D(\hat{\beta}_R)] + \hat{\beta}^T W^2 \hat{\beta} \lambda^2$$

$$= \sigma^2 \sum_{i=0}^k \frac{\lambda_i}{(\lambda_i + \lambda)^2} + \lambda^2 \sum_{i=0}^k \frac{\beta_i^2}{(\lambda_i + \lambda)^2}$$

$$= \sigma^2 \sum_{i=0}^k \frac{\lambda_i}{(\lambda_i + \lambda)^2} + \sum_{i=0}^k \frac{\beta_i^2}{(1 + \lambda_i/\lambda)^2}$$

λ_i 's are eigenvalues of $X^T X$.

Solution space for RIDGE Geometrical Intuition



Similar steps for LASSO (differentiation is a bit complex)

Variable Selection

$$\begin{aligned} Y &= X\beta + \varepsilon \\ &= [x_p | x_r] \left(\frac{\beta_p}{\beta_r} \right) + \varepsilon \end{aligned}$$

$(p+r=k+1)$

Data we received : (Y, x_p)

We need r more features to fit data into true model

True Model $Y = X\beta + \varepsilon$

Reduced Model $Y = x_p \hat{\beta}_p + \varepsilon$ reduced according to available data

Q1. Is $\hat{\beta}_p$ an unbiased estimator?

Q2. Does $\hat{\beta}_p$ have less variation compared to the true one?

Q3. Is σ_p^2 (estimate value of σ^2 in reduced model) an unbiased estimator of σ^2 ?

$$\text{Note: } \Sigma = \left(\begin{array}{c|c} \sum_{pp} & \sum_{pb} \\ \hline \sum_{pr} & \sum_{rr} \end{array} \right)$$

reduced model, \sum_{pp} is to be used.

In reduced model,

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

$$\begin{aligned} E(\hat{\beta}_p) &= (x_p^T x_p)^{-1} x_p^T E(y) \\ &= (x_p^T x_p)^{-1} x_p^T [x_p \beta_p + x_r \beta_r] \\ &= \beta_p + (x_p^T x_p)^{-1} (x_p^T x_r) \beta_r \end{aligned}$$

Hence, $\hat{\beta}_p$ is an unbiased estimator if $x_p^T x_r = 0$
 This would happen in case of PCA and orthogonal polynomials.

A1. Not true in general

$$\sigma_p^2 = \frac{y^T (I_n - x_p (x_p^T x_p)^{-1} x_p^T) y}{n-p}$$

$$\begin{aligned} E(\sigma_p^2) &= [\sigma^2(n-p) + \sigma^2(ncp)] / (n-p) \\ &= \sigma^2 + \sigma^2(ncp) / (n-p) \end{aligned}$$

$$\begin{aligned} ncp &= \frac{1}{\sigma^2} (x_B)^T (I_n - x_p (x_p^T x_p)^{-1} x_p^T) x_B \\ &= 0 + 0 + 0 + \frac{1}{\sigma^2} (\beta_r^T x_r^T (I_n - P x_p) x_r \beta_r) \\ &\quad \left. \begin{array}{l} \text{expand } x_B = x_p \beta_p + x_r \beta_r \\ \text{all } x_p \text{ terms get cancelled} \end{array} \right\} \end{aligned}$$

$$ncp = \frac{1}{\sigma^2} (\beta_r^T x_r^T (I_n - P x_p) x_r \beta_r)$$

A3. No (it is unbiased only when we capture all features)

Note: $\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$. Then,

$$\Sigma^{-1} = (\Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21})^{-1}, \text{ and}$$

$$(\Delta + BCD)^{-1} = A^{-1} - A^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1}$$

given in prof's slides

Project (Credits : Anubhab Mandal)

20/30 marks

Group size \rightarrow 1, 5 or 6

Report \rightarrow 5/6 pages ①

+ code (analysis
of dataset)
②

Justification of
maths + dataset
description using things
 taught in
 class)

+ PDF of ③
executed code
(Python or R)

(10 mins)

30 marks \rightarrow video recording ④
of project (Tams)

Upload all the stuff in a
drive folder

S All group members must submit project folder by the same name so that checked together

Submission → 15th April 2024

Generalized Linear Model

In linear model: $y_i = \beta^T x_i + \epsilon_i$, $i=1,2,\dots,n$

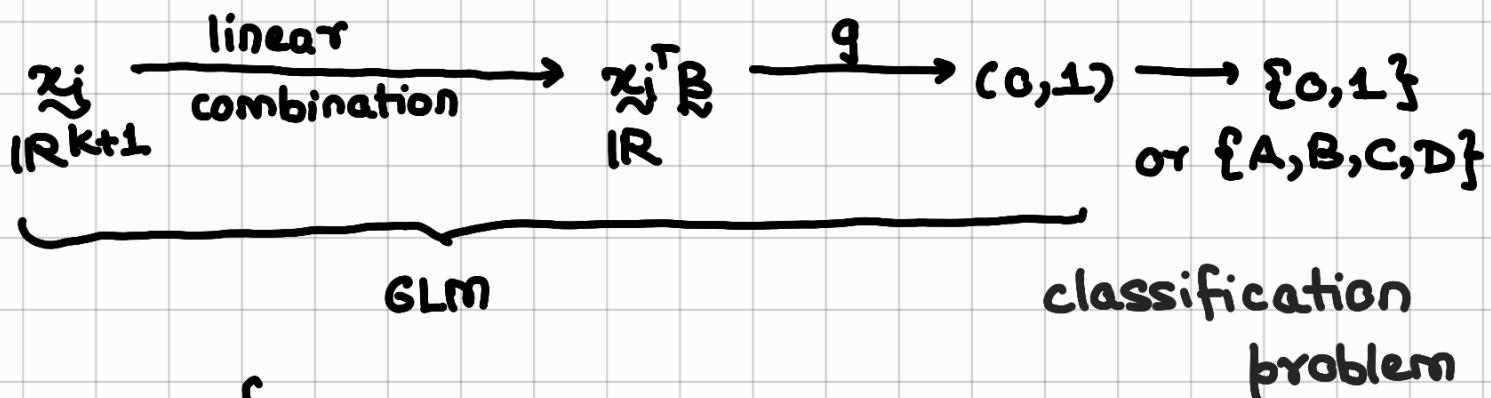
$$E(\epsilon_i) = 0 \quad V(\epsilon_i) = \sigma^2$$

$$\epsilon_i \stackrel{\text{ind}}{\sim} \mathcal{N}(0, \sigma^2)$$

$$E(y_i | x_i) = \beta^T x_i = I(\beta^T x_i) \quad (I(x) = x.)$$

For categorical variables as output, we need

$$E(y_i | x_i) = g(\beta^T x_i)$$



$$g(z) = \begin{cases} \frac{e^z}{1+e^z} & z \in \mathbb{R} \quad \text{logit-model} \\ \Phi(z) & z \in \mathbb{R} \quad \text{probit-model} \end{cases}$$

any cdf of a cont. random variable

Logistic Regression

$$P(y_i = 1) = 1 - P(y_i = 0) = \pi_i = E(y_i | x_i) = \frac{e^{x_i^T \beta}}{1 + e^{x_i^T \beta}}$$

$$\Leftrightarrow \log_e \left(\frac{\pi_i}{1 - \pi_i} \right) = x_i^T \beta \in \mathbb{R}$$

$\pi_i / (1 - \pi_i)$ is known as 'odd' of π_i

Two treatment problem; often defined as $\log\left(\frac{p_A}{1-p_A}\right) \rightarrow \text{log of odds ratio.} \approx \text{asymptotically follow normal distribution.}$

$$\psi = \log\left(\frac{p_A}{1-p_A}\right) ; \hat{\psi} \sim N(0, \sigma_y^2) \text{ for large } n.$$

Let $y \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$

$$f(y) = \prod_{i=1}^n p^{y_i} (1-p)^{1-y_i} \text{ joint pdf}$$

$$= p^{\sum_i y_i} (1-p)^{n - \sum_i y_i} = \left(\frac{p}{1-p}\right)^{\sum_i y_i} (1-p)^n$$

$$\log f(y) = (\sum_{i=1}^n y_i) (\log(p/(1-p))) + n \log(1-p) + \text{constant}$$

$$= D_1(y) D_2(p) + D_3(p) + D_4(y)$$

For any distribution (almost any), we are able to decompose $\log f(y)$ in this form.

$D_1(y) \rightarrow$ sufficient statistic (all info. of data)

$D_2(p) \rightarrow$ natural parameter structure

$D_3(p) \rightarrow$ efficient parameter estimation from data

$D_4(y) \rightarrow$ carries certain info. about ancillary part
depending on family of distribution. does
not depend upon parameter

Due to this natural structure, logistic model is preferred over other models.

$$\begin{aligned} \theta &= \log p/(1-p) \\ p &= e^\theta / (1+e^\theta) \end{aligned}$$

Now,

$$E(y_i|x_i) = \pi_i = \frac{e^{x_i^T \beta}}{1 + e^{x_i^T \beta}}$$

Likelihood fn. of (β)

$$L(\beta | \text{data}) = \prod_{i=1}^n \left\{ \pi_i^{y_i} (1 - \pi_i)^{1-y_i} \right\}$$

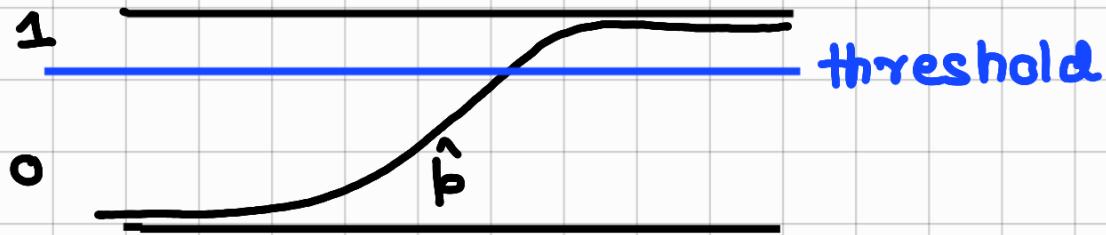
$$\log(L(\beta | y)) = \sum_{i=1}^n \left[y_i \log\left(\frac{\pi_i}{1-\pi_i}\right) \right] + \sum_{i=1}^n \log(1-\pi_i)$$

$$L(\beta) = \sum_{i=1}^n [y_i(x_i^T \beta)] - \sum_{i=1}^n \log(e^{x_i^T \beta})$$

$$\frac{\partial L(\beta)}{\partial \beta} = 0$$

we get $\hat{\beta}$ as an iterative solution

$$\hat{\beta} = \frac{e^{\hat{\beta} x}}{1 + e^{\hat{\beta} x}}$$



Error Analysis of MLR

$$\epsilon_i \sim \mathcal{N}(0, \sigma^2) \quad \text{assumption}$$

$$y_i = x_i^T \beta + \epsilon_i \quad \text{data}$$

$$e_i = y_i - \hat{y}_i \quad \text{estimated error}$$

$$H \equiv P_X$$

$$e = (y - \hat{y}) \sim \mathcal{N}(0, \sigma^2(I_n - P_X))$$

$$\hat{\sigma}^2 = \frac{e^T e}{n-k-1} \quad \begin{matrix} \text{unbiased estimate} \\ \text{of } \sigma^2 \end{matrix}$$

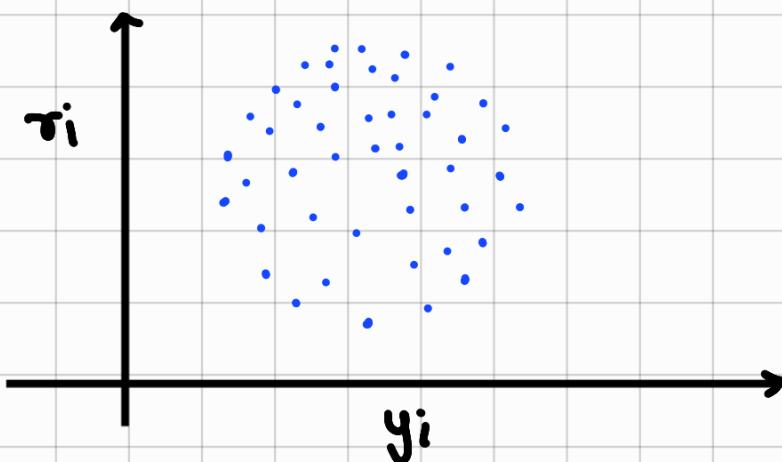
$$\text{cov}(e_i, e_j) = \begin{cases} \sigma^2(1-h_{ii}) & i=j \\ -\sigma^2 h_{ij} & i \neq j \end{cases}$$

$$MSE_{\text{error}} = \frac{SSE_{\text{error}}}{n-k-1}$$

Standardized Residual $d_i = e_i / \sqrt{MSE_{\text{error}}}$

Studentized Residual $r_i = e_i / \sqrt{MSE_{\text{error}}(1-h_{ii})}$

If all the analysis is correct, the pattern resembles :



Homework: Show that

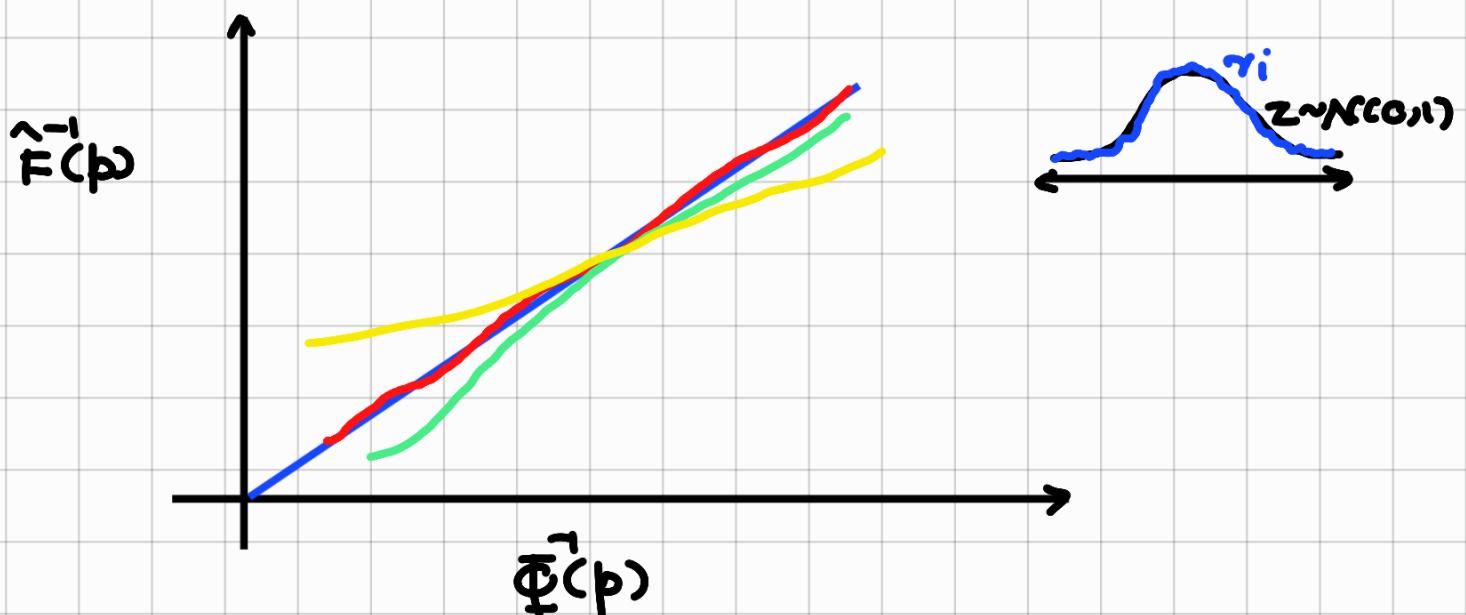
$$r_i = \frac{e_i}{\left(MSE_{\text{error}} \left[1 - \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{s_{xx}} \right) \right] \right)^{1/2}} \text{ for SLR.}$$

If $n \uparrow \infty$ and $k < \infty$, $\hat{\sigma}^2 = MSE_{\text{error}}$ will converge to σ^2 with probability 1.
 $\Rightarrow r_i \stackrel{a}{\sim} \mathcal{N}(0, 1)$

To validate this,

- Q-Q plot
- χ^2 -test
- KS test

Q-Q Plot



Tail part can most likely fluctuate. The middle part more or less follows the plot

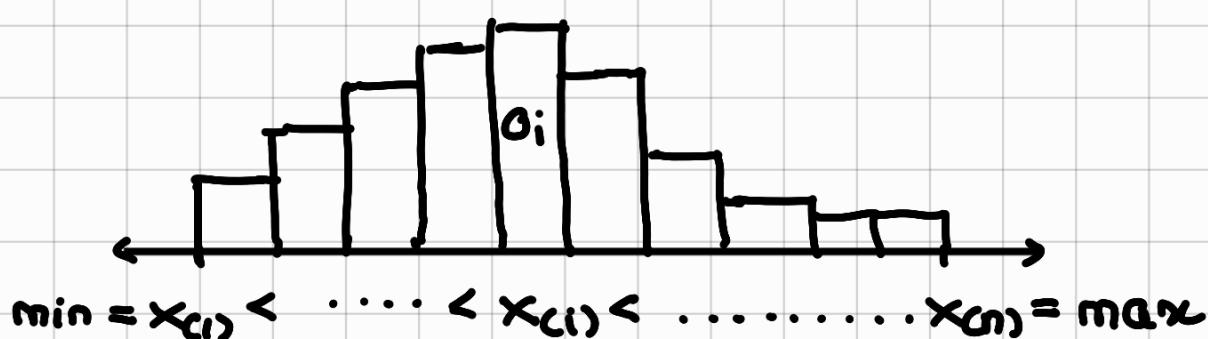
$\Phi^{-1}(p)$ cdf of $N(0,1)$
 \hat{F} is the empirical cdf of plot

χ^2 -test (goodness of fit)

pdf based test

$$\min = x_{(1)} < \dots < x_{(l)} < \dots \dots \dots x_{(n)} = \max$$

- arrange the data in an increasing order
- divide the range $(x_{(1)}, x_{(n)})$ into k -parts, such that each subinterval has more than 5 data



→ Let there are O_i many observations in the i th interval

→ E_i = expected no. of observations within i th interval

$$= n \times p_i$$

$$= (\text{total obs}) \times \text{prob. of } i\text{th interval}$$

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{k-1}$$

$df = k-1$ when parameters are known

df gets decremented with increasing no. of unknown parameters.

If we use $\hat{E}_i = n \times \hat{p}_i$, \hat{p}_i is the probability after estimating parameters from the data, then it will follow $\chi^2_{(k-1) - (\text{no. of parameters estimated})}$

KS Test

Assume that the target distribution specification is completely known.

$$Z_i \sim F(\text{completely known})$$

Z_i are continuous r.v.

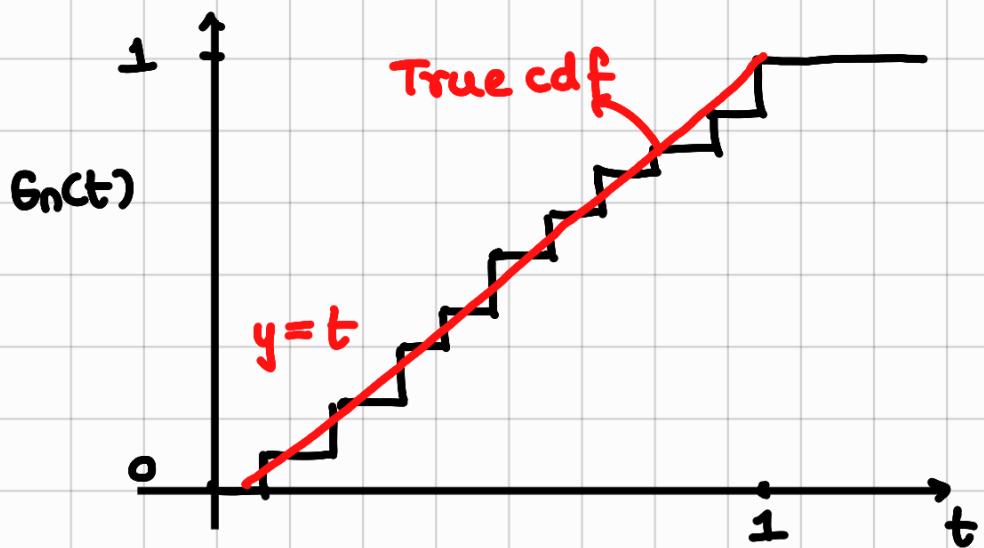
→ test only applicable to continuous r.v.

$$W_i = F(Z_i) \stackrel{iid}{\sim} U(0,1)$$

Empirical cdf of W is given as

$$G_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{\{w_i \leq t\}}$$

Plot of $G_n(t)$ v/s t looks like:



Test statistic : $\sup_{0 \leq t \leq 1} \sqrt{n} |G_n(t) - t| = D$

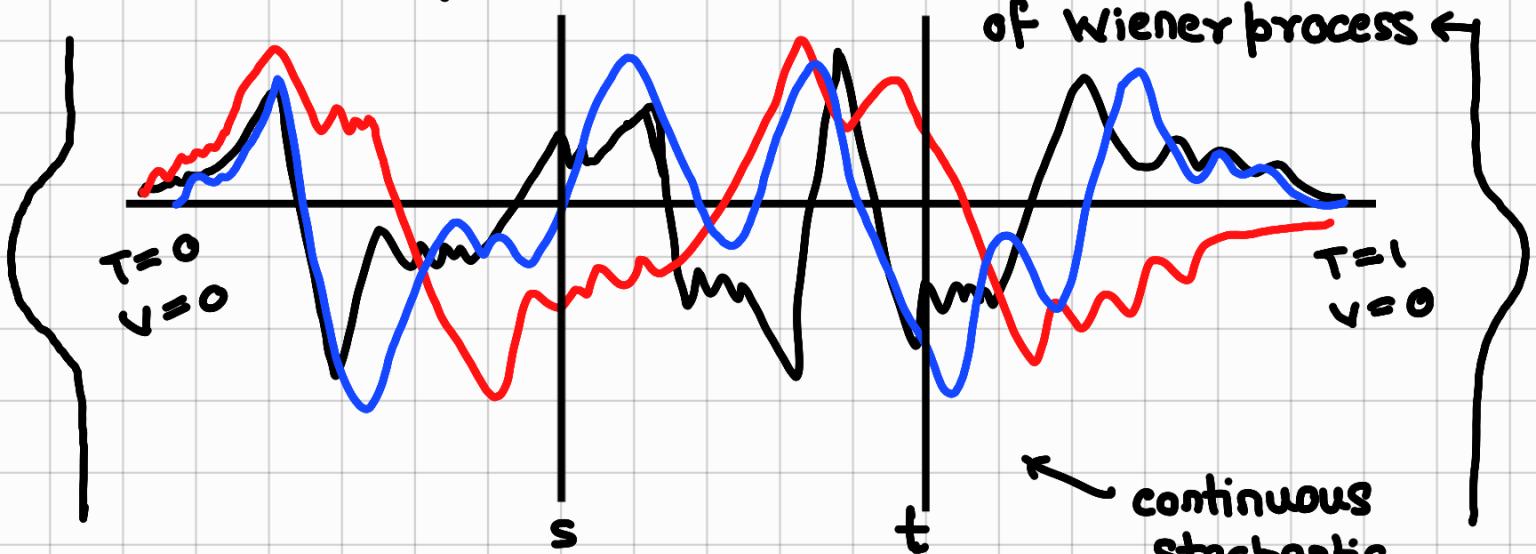
D is said to follow Kolmogorov distribution as $n \uparrow \infty$.

If $Z_i \sim F$ then $|G_n(t) - t| \rightarrow 0$ a.s. (almost surely)
 {Measure Theory PTSD}

Kolmogorov Distribution :

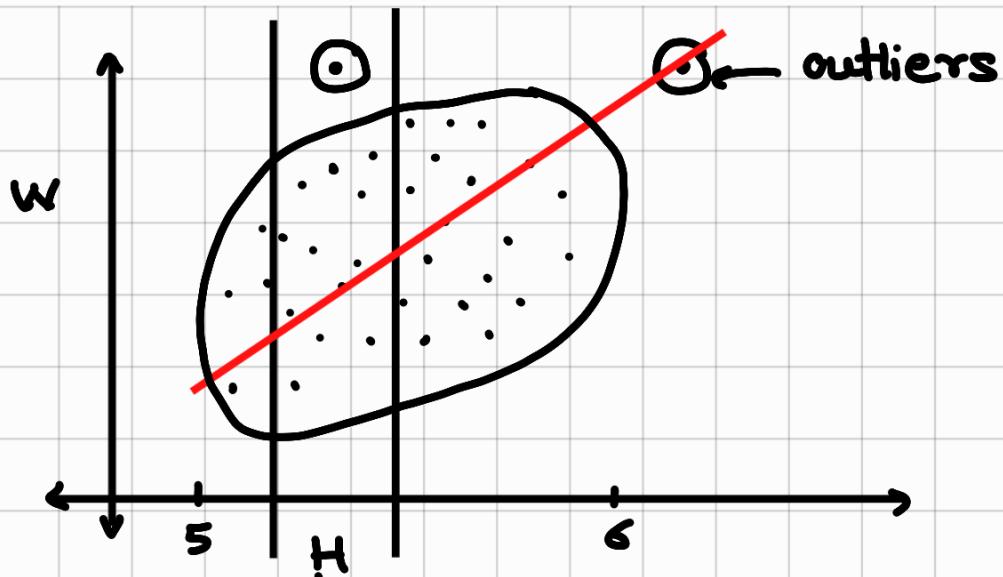
Brownian Bridge :

conditional distribution
of Wiener process ↪



$$\begin{pmatrix} B_0(s) \\ B_0(t) \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} s(1-s) & s(1-t) \\ s(1-t) & t(1-t) \end{pmatrix}\right)$$

$$K = \sup_{0 \leq t \leq 1} |B_0(t)| \leftarrow \text{r.v. of Kolmogorov distn.}$$



$$e_i = y_i - \hat{y}_i \\ = y_i - \tilde{x}_i^T \hat{\beta} \\ \left. \begin{array}{l} (y_1, \tilde{x}_1) \\ (y_2, \tilde{x}_2) \\ \vdots \\ (y_n, \tilde{x}_n) \end{array} \right\}_{n-1}$$

get $\hat{\beta}_{(1)}$
 $\hat{y}_{(1)} = \tilde{x}_1^T \hat{\beta}_{(1)}$
 $e_{(1)} = (y_1 - \hat{y}_{(1)}) = (y_1 - \tilde{x}_1^T \hat{\beta}_{(1)})$

$$\left. \begin{array}{l} (y_1, \tilde{x}_1) \\ \xleftarrow{(y_2, \tilde{x}_2)} \\ \vdots \\ (y_n, \tilde{x}_n) \end{array} \right\}_{n-1}$$

get $\hat{\beta}_{(2)}$
 $\hat{y}_{(2)} = \tilde{x}_2^T \hat{\beta}_{(2)}$
 $e_{(2)} = (y_2 - \hat{y}_{(2)}) = (y_2 - \tilde{x}_2^T \hat{\beta}_{(2)})$

Jack knife method

for $i=1$ to n

- 1) remove the observation
- 2) do the analysis for $(n-1)$ data
- 3) Compute the estimated error for i th data

$$e_{ci} = y_i - \hat{y}_{ci} = y_i - \hat{\beta}_{ci}^T x_i$$

Prediction Residual sum of square (PRESS)

$$\sum_{i=1}^n e_{ci}^2 = \sum_{i=1}^n (y_{ci} - y_i)^2$$

It can be shown that $e_{ci}/\sqrt{\text{Var}(e_{ci})} = e_i/\sqrt{\sigma^2(1-h_{ii})}$

Hence, fitting the whole n observations serves the purpose.

$$e_{ci} = e_i/(1-h_{ii}) \sim N(0, \sigma^2/(1-h_{ii}))$$

$$e_{ci}/\sqrt{\text{Var}(e_{ci})} = \frac{e_i/(1-h_{ii})}{\sqrt{\sigma^2/(1-h_{ii})}} = \frac{e_i}{\sqrt{\sigma^2(1-h_{ii})}} = t_i$$

$\hat{\sigma}^2 = \text{MSE}_{\text{Error}}$ may not be a good estimator when the i th observation has been removed

We can use

$$\begin{aligned}\hat{\sigma}^2 &= S_{ci}^2 = \frac{(n-k-1)\text{MSE}_{\text{Error}} - (e_i^2/(1-h_{ii}))}{n-k-2} \\ &= \frac{\sum_{i=1}^n e_i^2 - e_i^2/(1-h_{ii})}{n-k-2}\end{aligned}$$

$$T = \frac{e_{ci}}{S_{ci}^2/(1-h_{ii})} \sim t_{n-k-2} \text{ under } H_0$$

H_0 : i th observation is NOT an outlier

H_1 : i th observation is an outlier

$$\left[\underset{(k+1) \times (n-1)}{x_{ci}^T} \underset{(n-1) \times (k+1)}{x_{ci}} \right]^{-1} = \left[\underset{(k+1) \times n}{x^T x} - \underset{n \times (k+1)}{\underbrace{x_j}_{(k+1) \times 1}} \underset{1 \times (k+1)}{\underbrace{x_i^T}_{(k+1) \times 1}} \right]^{-1}$$

$$(A + \underset{n \times n}{\underbrace{uv^T}_{\text{rank 1}}})^{-1} = A^{-1} - \frac{A^{-1}uv^TA^{-1}}{1 + v^TA^{-1}u}$$

For related proofs, refer to Appendix C.7 and C.8 from the reference book { The prof is not accountable for the proof :)

Regression syllabus is completed.

Time-Series

Relation b/w exponential and geometric distribution:

$x \sim \exp(\lambda) \rightarrow \text{continuous}$



$$F(x) = 1 - e^{-\lambda x}$$

Consider $Y = [x]$

$$\begin{aligned} P(Y = r) &= P([x] = r) \\ &= P(r \leq x < r+1) \\ &= P(x < r+1) - P(x < r) \\ &= (1 - e^{-\lambda(r+1)}) - (1 - e^{-\lambda r}) \\ &= (e^{-\lambda})^r (1 - e^{-\lambda}) \end{aligned}$$

$\approx q^r p \leftarrow \text{geometric}(p = 1 - e^{-\lambda}) \rightarrow \text{discrete}$

This implies that the adjectives of time (discrete, continuous) depend upon the context of approximation of time, and not the data.

Hence, discrete/continuous are adjectives of time and not the data.

Example of time-series

e.g. 1 White Noise

A time series $\{W_t\}$ is said to follow white noise if

- $E(W_t) = 0$
- $V(W_t) = \sigma_w^2$
- and they are UNCORRELATED

$$W_t \sim WN(0, \sigma_w^2)$$

(WN stands for white noise)

- $x_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$
- $x_i \sim \begin{cases} N(0, 1) & \text{when } i \text{ is even} \\ \exp(i) - 1 & \text{when } i \text{ is odd} \end{cases}$

Homework : Construct uncorrelated but dependent WN for both the cases.

- (I) WN need not be normally distributed
- (II) WN need not be iid
- (III) iid sequence with zero mean and finite variance are always WN.
- (IV) WN is weakly stationary.
- (V) WN with normal distn. are strongly stationary.

e.g. 2 Binary process

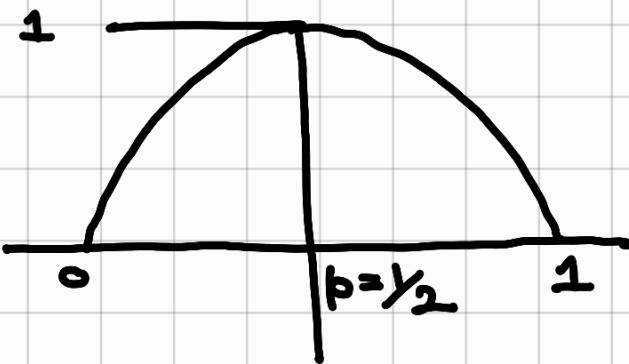
$$x_t = \begin{cases} +1 & \text{with prob. } p = \frac{1}{2} \\ -1 & \text{with prob. } 1-p = \frac{1}{2} \end{cases}$$

$$E(x_t) = 2p - 1$$

$$\text{Var}(x_t) = 4p(1-p)$$

$$z_t \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$$

$$x_t = 2z_t - 1$$



e.g. 3 Random walk on \mathbb{Z}



$$x_0 = 0$$

$$x_t = x_0 + \sum_{t=1}^T w_t$$

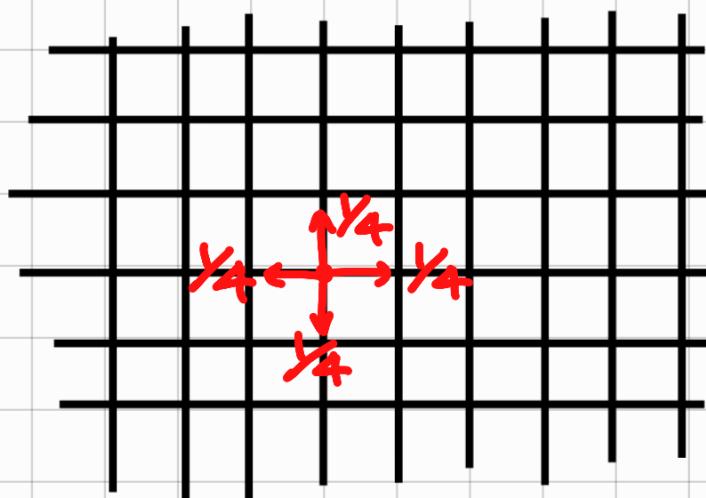
where w_t is a binary process with prob. (p).

$\binom{2n}{n} p^n (1-p)^n$: probability that it returns to 0 in $2n$ steps

Apply Stirling's approximation and check for $n \rightarrow \infty$.

(i) If ($p=y_2$) then sequence will eventually return to zero

(ii) Random walk on \mathbb{Z}^2



In \mathbb{Z}^2 with equal probability y_4 the sequence will eventually return to $(0,0)$.

Reference for proofs of above results:
Markov Chains by Norris

e.g.4. Random walk width drift

$$W_t \stackrel{\text{iid}}{\sim} E(W_t) = \delta$$

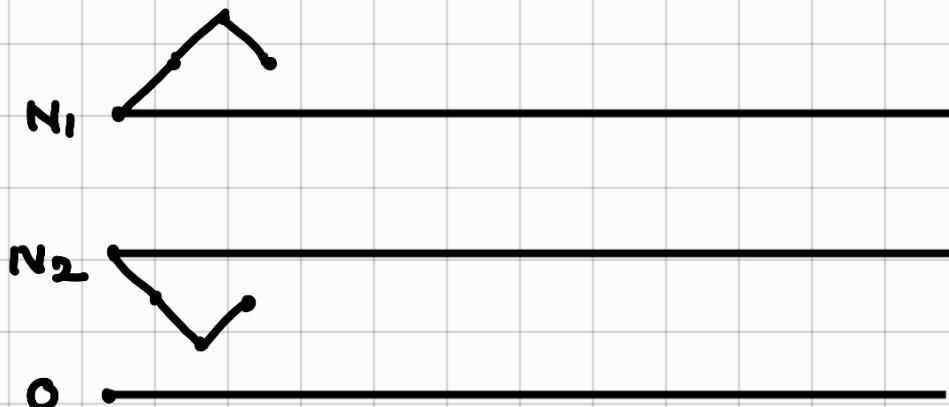
$$V(W_t) = \sigma_w^2$$

$$x_t = x_0 + \sum_{i=1}^t w_i$$

$$E(x_t) = \delta t$$

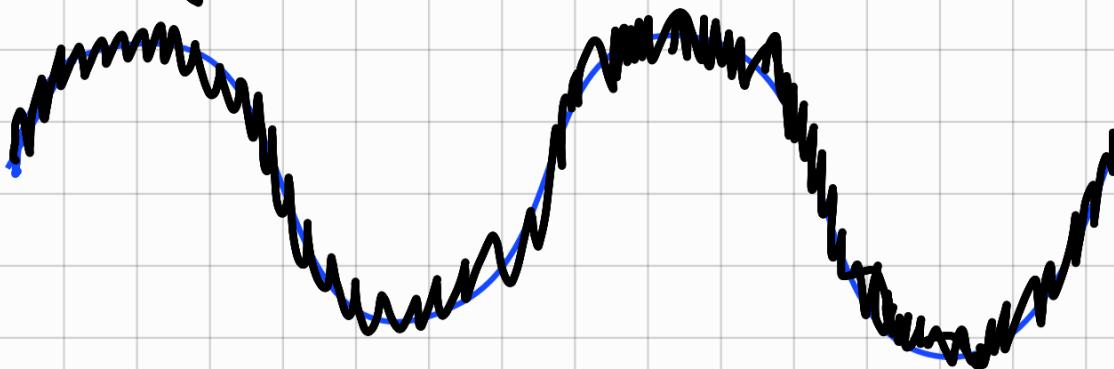
$$V(x_t) = \sigma^2 t$$

Gambler's ruin problem



First player which reaches O coins loses.

e.g.5 Signal with noise



$$x_t = A \sin(2\pi f t + \phi) + w_t$$

$$w_t \sim WN(0, \sigma_w^2)$$

e.g.6 Moving average prices (order one)

MACI process

$$\text{Let } w_t \sim WN(0, \sigma^2)$$

$$x_t = 1w_t + 8w_{t-1}$$

e.g. 7 Autoregressive process (order one)

AR(1)

$$x_t = \phi x_{t-1} + w_t \quad |\phi| < 1, \phi \neq 0$$

→ analogous / similar to the regression problem:

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

Wiener process (Brownian motion)

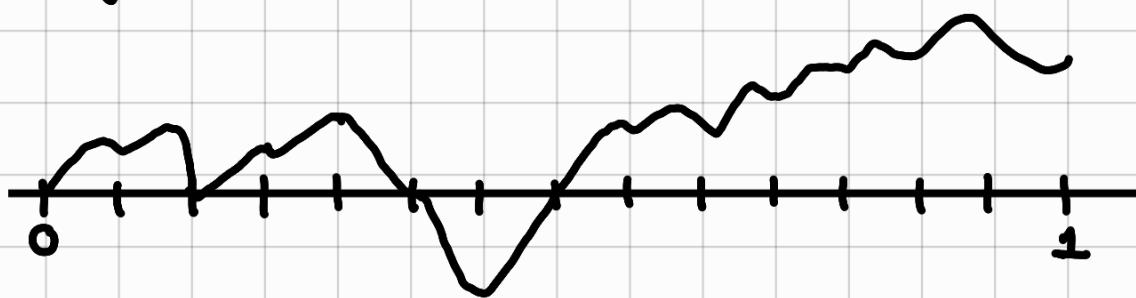
For definition, refer to the slides.

Generating data from Brownian motion (#)

$$T_1 = \sup_{0 < t < 1} W_t$$

$$T_2 = \int_0^1 W_t^2 dt$$

statistics used to measure qualities of brownian motion



$W_t \sim$ Brownian motion

$$y = x^2 \text{ on } (0,1)$$

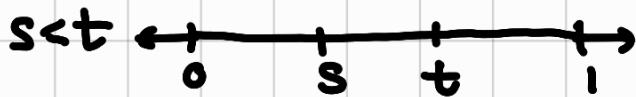
Remark: Let $\{w_i\}$ be a sequence of iid rvs with $E(w_i) = 0$ and $\text{Var}(w_i) = 1$ then $y_n(t) = \frac{1}{\sqrt{n}} \sum_{k=1}^{[nt]} w_k$ converges to Wiener process x_t on $[0,1]$ for large n .

$$y_n(t) = \frac{1}{\sqrt{n}} \sum_{k=1}^{[nt]} w_k$$

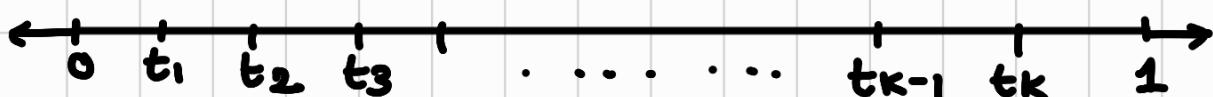
$$E(y_n(t)) = 0$$

$$\lim_{n \uparrow \infty} (\text{Var}(Y_n(t))) = \sqrt{\frac{[nt]}{n}} \left(\frac{1}{\sqrt{[nt]}} \sum_{k=1}^{[nt]} w_k \right) \xrightarrow{d} \sqrt{t} N(0, 1) \\ \equiv N(0, t)$$

Brownian Bridge :



$$\begin{pmatrix} Y_n(s) \\ Y_n(t) \end{pmatrix} \sim N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} s & \min(s, t) \\ \min(s, t) & t \end{pmatrix} \right)$$



$$\begin{pmatrix} Y_n(t_1) \\ Y_n(t_2) \\ \vdots \\ Y_n(t_k) \end{pmatrix} \sim N \left(\begin{pmatrix} 0 \\ \vdots \\ \vdots \end{pmatrix}, \begin{pmatrix} t_1 & \cdots & \min(t_1, t_j) & \cdots & t_k \\ \vdots & & \vdots & & \vdots \\ \cdots & \min(t_k, t_j) & \cdots & \cdots & t_k \end{pmatrix} \right)$$



Description of Brownian motion :

$$t = 0, X_0 = 0$$

$$t = t_1, X_{t_1} = X_0 + Z_1 \text{ where } Z_1 \sim N(0, t_1)$$

$$t = t_2, X_{t_2} = X_{t_1} + Z_2 \text{ where } Z_2 \sim N(0, t_2 - t_1)$$

$$\text{Var}(X_{t_2}) = 0 + t_1 + t_2 - t_1 = t_2$$

If we fix the initial and final time and we simulate a Brownian motion in b/w these, it forms a Brownian bridge.

Brownian Bridge ($B_0(t)$):

Let $B(t)$ be a brownian motion

$$B(t_1) = a \text{ and } B(t_2) = b, \quad t_1 < t_2$$

$$\text{mean} = a + \frac{t-t_1}{t_2-t_1} (b-a) \equiv \text{function of } t$$

$$\text{Var} = \frac{(t_2-t)(t-t_1)}{(t_2-t_1)}$$

$$\text{cov} = \frac{(t_2-t)(s-t_1)}{(t_2-t_1)} \text{ between } B_0(t) \text{ and } B_0(s)$$

$$B_0(t) = B(t) - \frac{t}{T} B(T)$$

Standard Brownian bridge on $[0,1]$

$$(t=0, B_0(t)=0) \text{ and } (t=1, B_0(t)=0) \quad \left\{ \begin{array}{l} B_0(t) = B(t) - tB(1) \\ 0 \leq t \leq 1 \end{array} \right.$$

$$E(B_0(t)) = 0 + \frac{t-0}{1} (0-0) = 0$$

$$\text{Var}(B_0(t)) = \frac{(1-t)(t-0)}{1} = t(1-t)$$

$$\text{Cov}(B_0(s), B_0(t)) = \frac{(1-t)(s)}{1} = s(1-t) \quad s < t$$

Brownian Bridge ($B_0(t)$)

→ KS Test for goodness of fit

→ Change point detection in temporal data

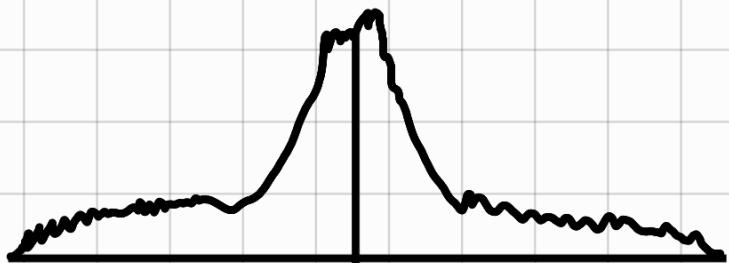
$$S_t := \frac{1}{\sqrt{n}} \sum_{i=1}^{[tn]} \frac{(x_i - \bar{x})}{\hat{\sigma}} ; \quad 0 \leq t \leq 1 \quad \text{CUSUM statistic}$$

↳ estimated process variance

$$S_t \Rightarrow B_0(t)$$

$$|S(t)|$$

$$T = \sup_{0 \leq t < 1} |S(t)|$$



→ these statistics are such that software cannot directly estimate.

→ we need to simulate the data to find the cut-off values.

$$T = \sup_{0 \leq t < 1} |S(t)| \rightarrow \sup_{0 \leq t \leq 1} |B_0(t)|$$

Hypothesis Testing :

$$H_0: E(x_1) = E(x_2) = \dots = E(x_n)$$

$$H_1: E(x_1) = E(x_2) = \dots = E(x_k) \neq E(x_{k+1}) \dots = E(x_n)$$

Formulation of variance and covariance :

$x_1, x_2, \dots, x_n \stackrel{iid}{\sim} F$ continuous distribution

$$Y = F(X) \sim U(0,1)$$

$$E(\sum_{i=1}^n I(y_i < t)) = nt, 0 \leq t \leq 1$$

$$V(\sum_{i=1}^n I(y_i < t)) = nt(1-t)$$

$$\text{Var}\left(\frac{1}{\sqrt{n}}(\sum_{i=1}^n I(y_i < t)) - \sqrt{n}t\right) = t(1-t)$$

$$H_n(t) = \sqrt{n} \left(\frac{1}{n} \sum I(y_i < t) - t \right)$$

$$\text{cov}(H_n(s), H_n(t)) \quad s < t$$

$$= \text{cov}\left(\frac{1}{\sqrt{n}} \sum_{i=1}^n I(y_i < t), \frac{1}{\sqrt{n}} \sum_{j=1}^n I(y_j < s)\right)$$

$$= \frac{1}{n} \sum_i \sum_j \text{cov}[I(y_i < t), I(y_j < s)]$$



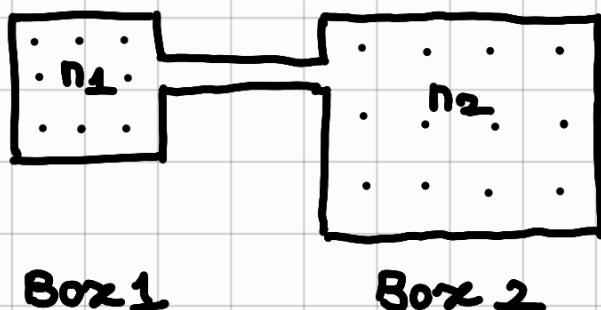
$$\text{cov} = \min(s, t) - st \\ = s(1-t)$$

$$\text{cov} = \frac{1}{n} \cdot n \left\{ \min(s, t) - st \right\} = s(1-t)$$

Stationary Time series

- strong stationarity
- weak stationarity

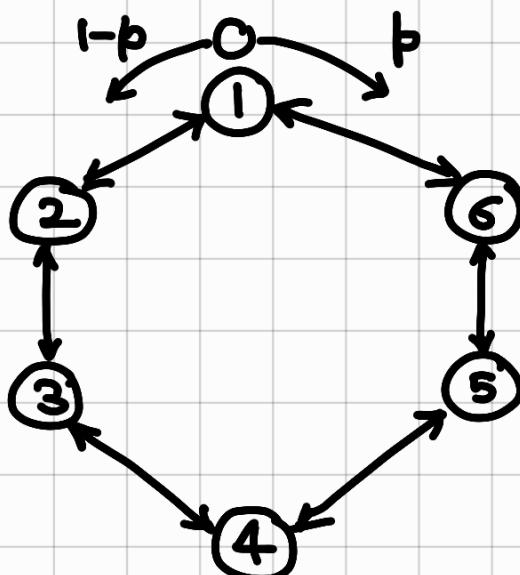
e.g. 1 $x_t \sim \text{bin}(n_1 + n_2, p_2)$ as $t \rightarrow \infty$



$$x_0 = n_1$$

x_t = # balls on box 1 at time t

e.g. 2



Limiting distribution of $\{x_t\}$ will follow uniform on $\{1, 2, 3, 4, 5, 6\}$. $P(x_t=1) = 1/6$.

x_t := location of the ball at time t .

Transition probability matrix :

	1	2	3	4	5	6
1	0	q	0	0	0	p
2	p	0	q	0	0	0
3	0	p	0	q	0	0
4	0	0	p	0	q	0
5	0	0	0	p	0	q
6	q	0	0	0	p	0

Strongly stationary timeseries

Let $\{x_t\}$ be a time series with a joint distribution fn. of (x_1, x_2, \dots, x_n) as

$$P(x_1 \leq a_1, x_2 \leq a_2, \dots, x_n \leq a_n) = F_n(a_1, a_2, \dots, a_n)$$

Now, $\forall n \in \mathbb{N}$, $\forall k \in \mathbb{Z}$, $\forall h \in \mathbb{Z}$ and $\forall a_i \in \mathbb{R}$,
If $P(x_{k+1} \leq a_1, x_{k+2} \leq a_2, \dots, x_{k+n} \leq a_n)$

$$= P(x_{k+h+1} \leq a_1, x_{k+h+2} \leq a_2, \dots, x_{k+h+n} \leq a_n)$$

$$= F(a_1, a_2, \dots, a_n),$$

then $\{x_t\}$ is strongly stationary.

Weakly stationary timeseries

A time series $\{x_t\}$ is said to be weakly stationary if

I) $\mu_t = E(x_t)$ is free from t.

II) $Cov(x_t, x_{t+h})$ is free from 't' but can be a function of 'h'.

→ this is similar to saying that this is upto second moment condition.

If a TS is strongly stationary with finite 2nd order moment, then it is weakly stationary.

Now, suppose $\{x_t\}$ is atleast weakly stationary.

$$\begin{aligned}\gamma_x(h) &= \text{Cov}(x_t, x_{t+h}) \\ &= E\{(x_t - \mu)(x_{t+h} - \mu)\} \\ &= E\{(x_{s-h} - \mu)(x_s - \mu)\} \\ &= \gamma_x(-h)\end{aligned}$$

... $t+h = s$ (say)

'h' is the lag.

Autocorrelation fn.

Auto-correlation coefficient or an 'atleast' weakly stationary timeseries is defined as

$$p_x(h) = \frac{\gamma_x(h)}{\gamma_x(0)}.$$

Significance/ Formulation :

$$\begin{aligned}\text{correlation coeff. } p_x(h) &= \frac{\text{cov}(x_{t+h}, x_t)}{\sqrt{\text{var}(x_{t+h}) \text{var}(x_t)}} \\ &= \frac{\gamma_x(h)}{\sqrt{\gamma_x(0) \gamma_x(0)}} \\ &= \frac{\gamma_x(h)}{\gamma_x(0)}.\end{aligned}$$

e.g.1 $\{x_t\}$ ^{WN or iid} $E(x_t) = 0, \text{var}(x_t) = \sigma^2$

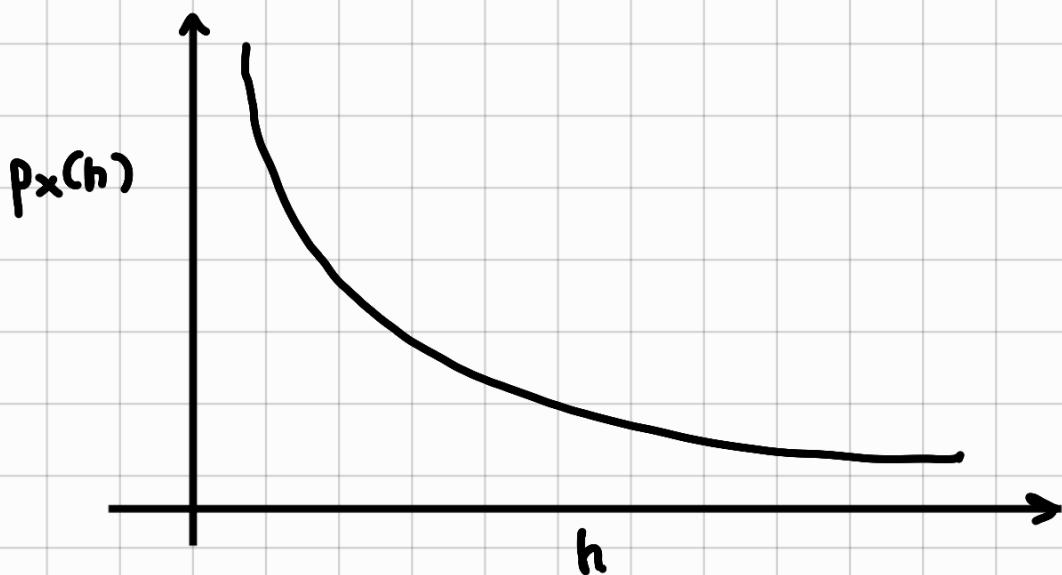
$$\gamma_x(h) = \begin{cases} \sigma^2 & h=0 \\ 0 & h \neq 0 \end{cases}$$

e.g.2 $s_t = \sum_{i=1}^t w_i$ $w_i \stackrel{iid}{\sim} N(0, \sigma^2)$

Random walk $\gamma_x(h) = \begin{cases} t\sigma^2 & h=0 \\ t\sigma^2 & h>0 \end{cases}$

$$\rho_X(h) = \frac{\text{cov}(S_t, S_{t+h})}{\sqrt{V(S_t)V(S_{t+h})}} = \frac{t\sigma^2}{\sqrt{t(t+h)}}$$

$$\rho_X(h) = \sqrt{\frac{t}{t+h}}$$



t is fixed.

e.g. 3 MA(1) Moving averages order one

$$Z_t = w_t + \theta w_{t-1} \text{ where } w_t \stackrel{iid}{\sim} N(0, \sigma^2)$$

$$E(Z_t) = 0$$

$$V(Z_t) = (1+\theta^2)\sigma^2$$

$$\gamma_Z(h) = \begin{cases} (1+\theta^2)\sigma^2 & h=0 \\ \theta\sigma^2 & h=\pm 1 \\ 0 & |h|>1 \end{cases}$$

e.g. 4 Autoregressive process of order one AR(1)

$$Z_t = \varphi Z_{t-1} + w_t \quad \left\{ \begin{array}{l} w_t \stackrel{iid}{\sim} N(0, \sigma^2) \\ |\varphi| < 1 \\ \varphi \neq 0 \\ Z_t \text{ is weakly stationary} \end{array} \right.$$

$$E(Z_t) = \varphi E(Z_{t-1}) + E(w_t)$$

$$E(z_t) = \varphi E(z_{t-1}) + 0$$

$$\Rightarrow E(z_t) = 0.$$

$$\begin{aligned} V(z_t) &= E(z_t^2) = E(\varphi^2 z_{t-1}^2 + w_t^2 + 2\varphi z_{t-1} w_t) \\ &= \varphi^2 E(z_{t-1}^2) + E(w_t^2) + 0 \\ &= \varphi^2 E(z_{t-1}^2) + \sigma^2 \end{aligned}$$

$$E(z_t^2) = \frac{\sigma^2}{1-\varphi^2}$$

$$V(z_t) = \frac{\sigma^2}{1-\varphi^2}$$

$$E(z_t) = 0$$

$$V(z_t) = \frac{\sigma^2}{1-\varphi^2} > \sigma^2$$

→ this is similar to BYN case.

$$(x, y) \sim \text{BYN}$$

$$V(y|x) = \sigma_y^2(1-p^2)$$

$$V(y) = \sigma_y^2$$

$$\gamma_z(h) = \begin{cases} (\frac{\sigma^2}{1-\varphi^2}) & h=0 \\ (\frac{\sigma^2}{1-\varphi^2}) \varphi^h & h \neq 0 \end{cases}$$

$$\gamma_z(h) = \text{cov}(z_{t+h}, z_t)$$

$$= \text{cov}(\varphi z_{t+h-1} + w_{t+h}, z_t)$$

$$= \varphi \text{cov}(z_{t+h-1}, z_t)$$

$$= \varphi^2 \text{cov}(z_{t+h-2}, z_t)$$

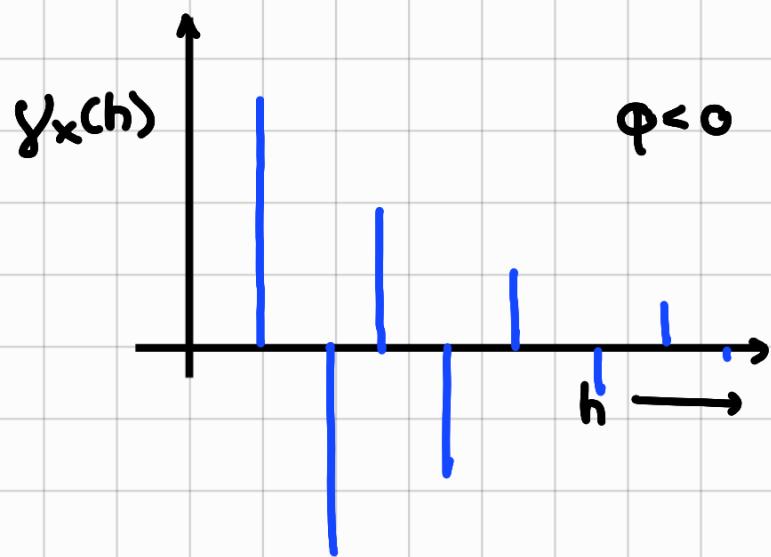
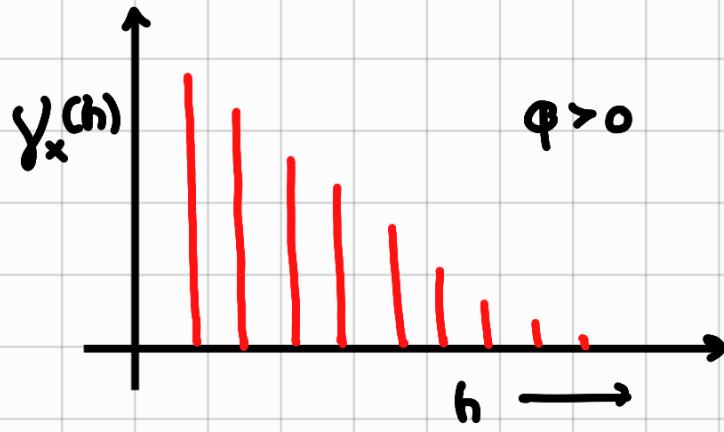
.....

$$= \varphi^h \text{cov}(z_t, z_t)$$

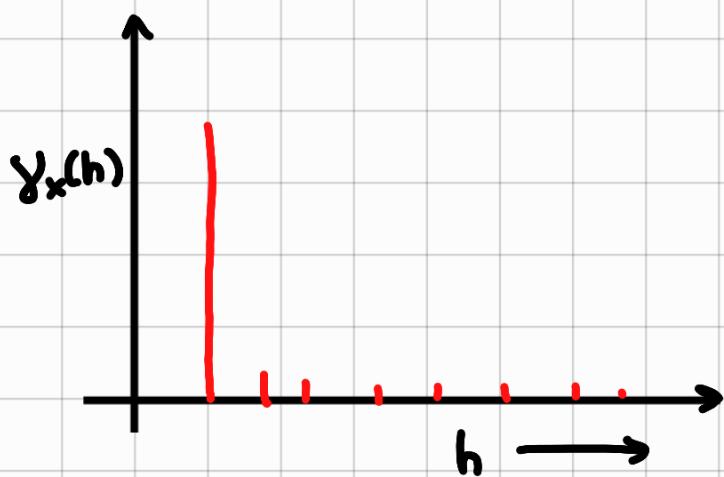
$$= \varphi^h \left(\frac{\sigma^2}{1-\varphi^2} \right) = \gamma_z(-h)$$

AR(1)

$$\gamma_x(h) = \gamma_x(-h) = \varphi^{|h|} \left(\frac{\sigma^2}{1-\varphi^2} \right)$$



MAC(1)



Auto-correlation fn.

Notations of operators

1. Backshift operator (B)

$$B^h x_t \equiv x_{t-h}$$

2. Difference operator (∇)

$$\nabla x_t \equiv x_t - x_{t-1} \equiv (I - B)x_t$$

$$\boxed{\nabla = I - B}$$

$$\nabla^h x_t = (I - B)^h x_t = \sum_{k=0}^h \binom{h}{k} (-1)^{h-k} (B^{h-k} x_t)$$

3. Seasonal Operator

$$\nabla_s = (I - B^s)$$

$$\nabla_s x_t = x_t - x_{t-s}$$

$$\text{Consider } f(x) = a_0 + a_1 x + a_2 x^2 = 1 + 2x + 3x^2$$

$$x=1, f(1) = 6$$

$$x=2, f(2) = 17$$

$$x=3, f(3) = 34$$

$$x=4, f(4) = 57$$

	∇	∇^2	∇^3
6			
17	11		
34	17	6	0
57	23	6	0

Linear process

A time series $\{x_t\}$ is a linear process if it can be represented as :

$$x_t = \mu + \sum_{j=-\infty}^{\infty} \psi_j z_{t-j} \text{ where } t \in \mathbb{Z},$$

$$= \mu + \sum_{j=-\infty}^{\infty} \psi_j B^j z_t$$

$$z_t \sim WN(0, \sigma^2),$$

$$\sum_{j=-\infty}^{\infty} |\psi_j| < \infty$$

$$= \mu + (\sum_{j=-\infty}^{\infty} \psi_j B^j) z_t$$

$$= \mu + \Psi(B) z_t$$

$$\Psi(B) = \sum_{j=-\infty}^{\infty} \psi_j B^j$$

$$x_t = \{\Psi(B) z_t\} + \mu$$

Note. $E(x)$ exists if $E(|x|) < \infty$, i.e.,

$$\sum_i (x_i) f(x_i) < \infty \text{ (or, } \int |x| f(x) dx < \infty)$$

Note. Absolutely summable series is always summable

Ex-1. $x_t \sim WN(0, \sigma^2)$

$$\psi_j = \begin{cases} 1 & \text{if } j=0 \\ 0 & \text{o.w.} \end{cases}$$

Ex 2. $x_t \sim \text{MA}(1)$

$$\psi_j = \begin{cases} 1 & \text{if } j=0 \\ 0 & \text{if } j=1 \\ 0 & \text{o.w.} \end{cases}$$

Ex 3. $x_t \sim \text{AR}(1)$

$$x_t = \phi x_{t-1} + z_t$$

$$\psi_j = \begin{cases} \phi^j & \text{if } j \geq 0 \\ 0 & \text{o.w.} \end{cases}$$

MA(q) process $q \in \mathbb{N}$

MA(q) process is a linear process where

$$\begin{cases} \psi_0 = 1 \\ \psi_j = \begin{cases} \theta_j & 1 \leq j \leq q \\ 0 & \text{o.w.} \end{cases} \end{cases}$$

$$\begin{aligned} x_t &= z_t + \sum_{j=1}^q \theta_j z_{t-j} \\ &= z_t + \left(\sum_{j=1}^q \theta_j B^j \right) z_t \\ &= (I + \sum_{j=1}^q \theta_j B^j) z_t \equiv \Theta_q(B) z_t. \end{aligned}$$

$z_t \sim \text{WN}(0, \sigma^2)$

$$\text{cov}(x_t, x_s) = \begin{cases} 0 & \text{if } |t-s| > q \\ \text{can be non-zero} & \text{o.w.} \end{cases}$$

If $z_t \stackrel{\text{iid}}{\sim}$ then, (x_t, x_j) are independent if $|t-s| > q$.

Auto-regressive Process of order $p \geq \text{AR}(p) p \in \mathbb{N}$

A process is said to be an AR(p) process if it can be represented as

$$x_t = \sum_{j=1}^p \varphi_j x_{t-j} + z_t$$

where $Z_t \sim WN(0, \sigma^2)$.

$$\Rightarrow (x_t) - (\sum_{j=1}^p \varphi_j B^j) x_t = Z_t.$$

$$\Rightarrow (I - \sum_{j=1}^p \varphi_j B^j) x_t = Z_t.$$

Let $\Phi_p(u) = (I - \sum_{j=1}^p \varphi_j B^j)$, then

$$\Phi_p(B) x_t = Z_t$$

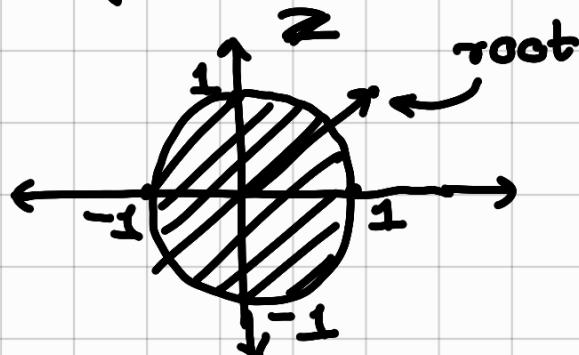
$$\Rightarrow x_t = \left(\frac{1}{\Phi_p(B)}\right) Z_t. \quad \rightarrow \textcircled{*}$$

$\frac{1}{\Phi_p(B)} = \frac{1}{1 - \sum_{j=1}^p \varphi_j B^j}$. This should satisfy

unit root condition for $\textcircled{*}$ to be feasible, that is, for eqn. $1 - q u \equiv 0$, root $u = \sqrt{q} > 1$.

(for our case, $q u \equiv \sum_{j=1}^p \varphi_j B^j$)

Also, $|q| < 1$, $q \neq 0$.



AR(1) process can be thought as MA(∞) process.

$$AR(1) = \lim_{q \rightarrow \infty} MA(q).$$

We can represent AR(1) process as MA(∞) process

$$x_t = \varphi x_{t-1} + Z_t \quad |\varphi| < 1 \quad \varphi \neq 0$$

$$\Rightarrow (I - \varphi B) x_t = Z_t$$

$$\Rightarrow x_t = \left(\frac{1}{1 - \varphi B}\right) Z_t$$

$$\Rightarrow x_t = \left(\sum_{j=0}^{\infty} (\varphi B)^j\right) Z_t$$

$$\Rightarrow x_t = \sum_{j=0}^{\infty} \varphi^j Z_{t-j}$$

$$x_t = (\sum_{j=0}^k \varphi_j z_{t-j}) + \varphi^{k+1} x_{t-(k+1)}$$

$$E(x_t - (\sum_{j=0}^k \varphi_j z_{t-j}))^2 = E\{(\varphi^{k+1} x_{t-(k+1)})^2\}$$

$$\begin{aligned}
 \lim_{k \rightarrow \infty} E(x_k - \sum_{j=0}^k \varphi^j z_{t-j}) &= \lim_{k \rightarrow \infty} E\{\varphi^{k+1} x_{t-(k+1)}\}^2 \\
 &= \lim_{k \rightarrow \infty} C(\varphi^{2k+2}) E(x_{t-(k+1)}^2) \quad \text{bounded} \\
 \left(\lim_{n \rightarrow \infty} E(|x_n - y|^2) \rightarrow 0 \Rightarrow x_n \xrightarrow{d} y. \right)
 \end{aligned}$$

$$\text{Hence, } x_t \stackrel{d}{=} \sum_{j=0}^{\infty} \varphi^j z_{t-j}$$

$$\lim_{n \uparrow \infty} (x_t - \sum_{j=0}^k \varphi^j z_{t-j})^2 = \lim_{k \uparrow \infty} \{ \varphi^{2k} E(x_{t-k}) \}$$

↓ AR(1) ↓ MA(k) ↓ bounded
0

$$y_k = \sum_{j=0}^k \varphi_j z_{t-j} \text{ such that } E|x_t - y_k(t)|^2 \rightarrow 0$$

$$\Rightarrow y_k(t) \xrightarrow{d} x_t \text{ as } k \uparrow \infty.$$

Note. x_1, x_2, \dots, x_n iid $E(x) = \mu$ $V(x) = \sigma^2$

$$T_k = \sqrt{k}(\bar{x} - \mu) / \sigma$$

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

Then $E(T_k - z)^2 \rightarrow 0 \Rightarrow T_k \xrightarrow{d} z$

Covariance of a Linear Process

Let $\{x_t\}$ be a weakly stationary process with $E(x_t) = 0$ and $\gamma_x(h) = \text{cov}(x_t, x_{t+h})$ exists.

If $\sum_{j=-\infty}^{\infty} |\psi_j| < \infty$, then defining

$y_t = \sum_{j=-\infty}^{\infty} \psi_j x_{t-j}$, we have

$$\rightarrow E(Y_t) = 0$$

$$\rightarrow \gamma_y(h) = \sum_{j=-\infty}^{\infty} \sum_{k=-\infty}^{\infty} \psi_j \psi_k \gamma_x(h+k-j)$$

Proof.

$$\begin{aligned} E|Y_t| &= E\left|\sum_{j=-\infty}^{\infty} \psi_j x_{t-j}\right| \leq \sum_{j=-\infty}^{\infty} |\psi_j| E|x_{t-j}| \\ &\leq M_1 \sum_{j=-\infty}^{\infty} |\psi_j| < \infty \end{aligned}$$

where $E|x_t| < M_1$.

Homework. Let Z be a random variable. Prove that
 $E(|Z|^\gamma) \leq \infty \implies E(|Z|^s) < \infty, \gamma \geq s$

Hence, $E|Y_t|$ exists.

$$E(Y_t) = E\left(\sum_{j=-\infty}^{\infty} \psi_j x_{t-j}\right) = \sum_{j=-\infty}^{\infty} \psi_j E(x_{t-j}) = 0.$$

$$\begin{aligned} E|Y_{t+h} Y_t| &= E\left|\sum_j \psi_j x_{t+h-j} \sum_k \psi_k x_{t-k}\right| \\ &= E\left|\sum_j \sum_k \psi_j \psi_k x_{t+h-j} x_{t-k}\right| \\ &\leq \sum_j \sum_k |\psi_j| |\psi_k| E(|x_{t+h-j}| |x_{t-k}|) \\ &\leq \sum_j \sum_k |\psi_j| |\psi_k| M_2 \end{aligned}$$

where $\max_{t,h \in \mathbb{Z}} E|x_{t+h-j} x_{t-k}| < M_2$.

So,

$$E|Y_{t+h} Y_t| \leq \sum_j \sum_k |\psi_j| |\psi_k| M_2 = M_2 (\sum_j |\psi_j| \sum_k |\psi_k|) < \infty.$$

$$\begin{aligned} \gamma_Y(h) &= \text{cov}(Y_{t+h} Y_t) \\ &= \text{cov}\left(\sum_j \psi_j x_{t+h-j} \sum_k \psi_k x_{t-k}\right) \\ &= E\left(\sum_j \sum_k \psi_j \psi_k x_{t+h-j} x_{t-k}\right) \\ &= \sum_j \sum_k \psi_j \psi_k E(x_{t+h-j} x_{t-k}) \\ &= \sum_j \sum_k \psi_j \psi_k \text{cov}(x_{t+h-j}, x_{t-k}) \\ &= \sum_j \sum_j \psi_j \psi_k \gamma_X(h-j) \end{aligned}$$

Homework. $x_t \sim WN(0, \sigma^2)$ and y_t is a linear process.

$$\text{Then, } \gamma_Y(h) = \sigma^2 \sum_{j=-\infty}^{\infty} \psi_j \psi_j h.$$

Homework. $y_t \sim AR(1)$. $\gamma_y(h) = \frac{\phi^h \sigma^2}{1-\phi^2}$.

ARMA process $\overbrace{ARMA(p,q)}$

A time series $\{x_t\}$ is said to follow ARMA (p,q) process if we have the representation

$$(x_t - \sum_{j=1}^p \phi_j x_{t-j}) = z_t + \sum_{j=1}^q \theta_j z_{t-j} \quad \forall t \in \mathbb{Z}, z \sim WN.$$

$$\begin{aligned} & \rightarrow (I - \sum_{j=1}^p \phi_j B^j) x_t = (I + \sum_{j=1}^q \theta_j B^j) z_t \\ & \rightarrow \Phi_p(B) x_t = \Theta_q(B) z_t \text{ where} \end{aligned}$$

$$\Phi_p(x) = 1 - \sum_{j=1}^p \phi_j x^j \text{ and } \Theta_q(x) = 1 + \sum_{j=1}^q \theta_j x^j$$

$\Phi_p(x)$ and $\Theta_q(x)$ should not have any common root.

For example, ARMA $(1,1)$

$$\begin{aligned} x_t - \phi x_{t-1} &= z_t + \theta z_{t-1} \\ \Rightarrow x_t &= \left(\frac{I + \theta B}{I - \phi B} \right) z_t \end{aligned}$$

$$\begin{cases} |\phi| < 1, \phi \neq 0 \\ \theta + \phi \neq 0 \end{cases}$$

$$\begin{aligned} x_t &= \left\{ (I + \theta B) \left(\sum_{j=0}^{\infty} (\phi B)^j \right) \right\} z_t = \Psi(B) z_t \\ &= (\psi_0 + \sum_{j=1}^{\infty} \psi_j B^j) z_t \end{aligned}$$

$$\psi_0 = 1$$

$$\psi_j = \begin{cases} (\theta + \phi) \phi^{j-1} & \forall j \geq 1 \end{cases}$$

Linear process representation

$$x_t = z_t + (\theta + \phi) \sum_{j=1}^{\infty} \phi^{j-1} z_{t-j}$$

$$\begin{aligned} y_x(\omega) &= \sigma^2 \sum_{j=0}^{\infty} \psi_j^2 = \sigma^2 [1 + (\theta + \phi)^2 \sum_{j=1}^{\infty} \phi^{j-2}] \\ &= \sigma^2 \left[1 + \frac{(\theta + \phi)^2}{1 - \phi^2} \right] \end{aligned}$$

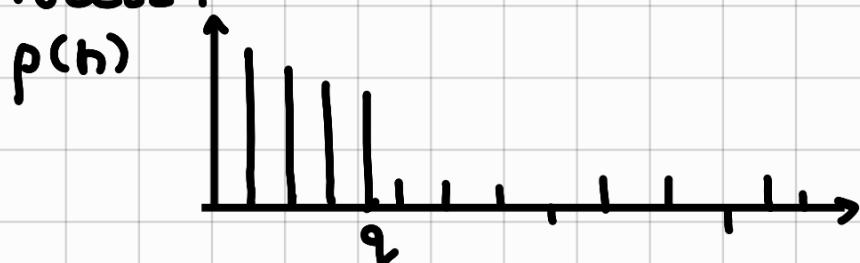
$$\gamma_x(1) = \sigma^2 \sum_{j=0}^{\infty} \psi_j \psi_{j+1} = \sigma^2 [(\theta+\phi)^2 + \frac{(\theta+\phi)^2 \phi}{1-\phi^2}]$$

$\gamma_x(h) = \phi^{h-1} \gamma_x(1)$. Homework.

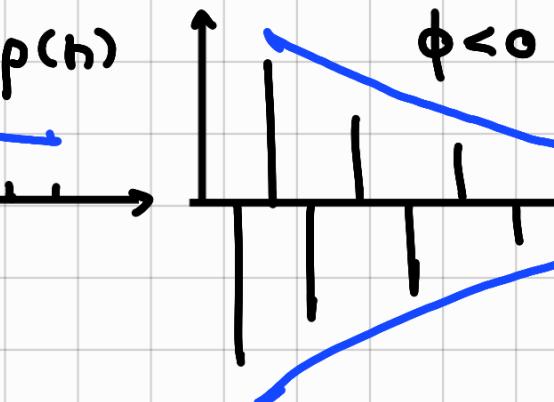
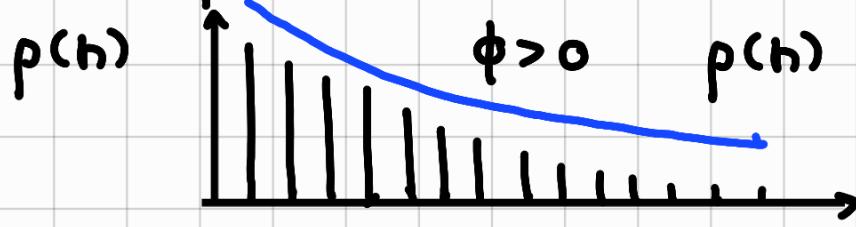
Autocorrelation coefficient of weakly stationary process.

$$\rho(h) = \frac{\gamma_x(h)}{\gamma_x(0)}$$

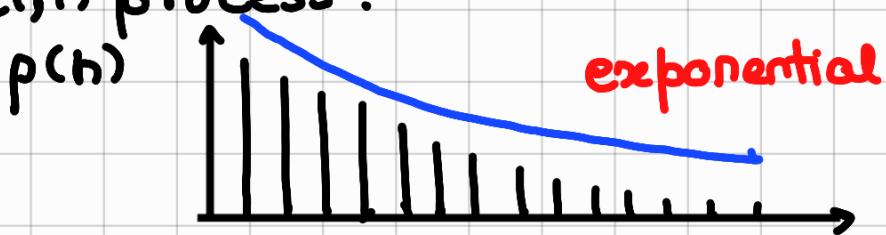
MA(q) process :



AR(1) process :



ARMA(1,1) process :



ACF can be used to distinguish b/w MA and (AR/ARMA)
But it cannot distinguish b/w AR & ARMA.
To do so, we need partial auto correlation coefficient.
(PACF)

Partial correlation Coefficient

given, $\begin{pmatrix} Y \\ Z \\ \vdots \\ X \end{pmatrix} \sim N \left(\begin{pmatrix} \mu_Y \\ \mu_Z \\ \vdots \\ \mu_X \end{pmatrix}, \begin{pmatrix} \sum_{YZ} & \sum_{YZX}^T \\ \sum_{YZX} & \sum_X \end{pmatrix} \right)$

$Y \in \mathbb{R}^1$

$Z \in \mathbb{R}^1$

$X \in \mathbb{R}^p$

$$E(Y) = \mu_Y \quad E(Z) = \mu_Z \quad E(X) = \mu_X$$

$$\Sigma_{YZ} = \begin{pmatrix} \sigma_{YY} & \sigma_{YZ} \\ \sigma_{YZ} & \sigma_{ZZ} \end{pmatrix}_{2 \times 2} \quad \Sigma_X = \sum_X = ((\sigma_{XiXj}))_{p \times p}$$

$$\sigma_{YX} = \text{cov}(Y, X) = \begin{pmatrix} \sigma_{YX_1} \\ \sigma_{YX_2} \\ \vdots \\ \sigma_{YX_p} \end{pmatrix}$$

$$\sigma_{ZX} = \text{cov}(Z, X) = \begin{pmatrix} \sigma_{ZX_1} \\ \sigma_{ZX_2} \\ \vdots \\ \sigma_{ZX_p} \end{pmatrix}$$

$$\Sigma_{YZX} = (\Sigma_{YX}, \Sigma_{ZX})_{p \times 2} \quad \Sigma_{YZX}^T = \begin{pmatrix} \Sigma_{YX}^T \\ \Sigma_{ZX}^T \end{pmatrix}_{2 \times p}$$

Partial correlation coefficient is the correlation b/w Y and Z after removing the effect of X and is denoted as $\rho_{YZ \cdot X}$

$$\rho_{YZ \cdot X} = \frac{\text{cov}(Y - E(Y|X), Z - E(Z|X))}{\sqrt{\text{var}(Y - E(Y|X)) \text{var}(Z - E(Z|X))}}$$

$$\Rightarrow \rho_{YZ \cdot X} = \frac{\text{cov}(e_{Y|X}, e_{Z|X})}{\sqrt{\text{var}(e_{Y|X}) \text{var}(e_{Z|X})}}$$

$$= \frac{\sigma_{yz} \cdot x}{\sqrt{\sigma_{yy} \cdot x \sigma_{zz} \cdot x}}$$

When we use it for time series we call it PACF.

conditional distribution of $(\begin{matrix} y \\ z \end{matrix}) \mid x = \tilde{x}$.

This also follows normal distribution,

$$N \left(\begin{pmatrix} \mu_y \\ \mu_z \end{pmatrix} + \sum_{yzx}^T \sum_x^{-1} (\tilde{x} - \mu_x), \sum_{yz} - \sum_{yzx}^T \sum_x^{-1} \sum_{yzx} \right)$$

Case : For Bivariate Normal,

$$\begin{pmatrix} u \\ v \end{pmatrix} \sim BVN \left(\begin{pmatrix} \mu_u \\ \mu_v \end{pmatrix}, \begin{pmatrix} \sigma_u^2 & \sigma_{uv} \\ \sigma_{uv} & \sigma_v^2 \end{pmatrix} \right)$$

$$v|u=x \sim N \left(\mu_v + \rho \frac{\sigma_v}{\sigma_u} (x - \mu_u), (1 - \rho^2) \sigma_v^2 \right)$$

$$\begin{aligned} E(v|u=x) &= \mu_v + \rho \frac{\sigma_u \sigma_v}{\sigma_u^2} (x - \mu_u) \\ &= \mu_v + \frac{\text{cov}(u, v)}{V(u)} (x - \mu_u) \\ &= \mu_v + \text{cov}(u, v) (V(u))^{-1} (x - \mu_u) \end{aligned}$$

$$\begin{aligned} V(v|u=x) &= \sigma_v^2 - \text{cov}(u, v) (V(u))^{-1} \text{cov}(u, v) \\ &= \sigma_v^2 - \rho \sigma_u \sigma_v (\sigma_u^2)^{-1} \rho \sigma_u \sigma_v \\ &= \sigma_v^2 - \rho^2 \sigma_v^2 \\ &= (1 - \rho^2) \sigma_v^2 \end{aligned}$$

Mean and the variance-covariance structure will remain same, if we do best linear prediction of $(\begin{matrix} y \\ z \end{matrix})$ based on

\mathbf{x} , even when $(\frac{Y}{X})$ are not jointly normally distributed
but 2nd order moments exists.