

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

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

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

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

### Multicollinearity

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

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

If  $|X^T X| = 0$  or  $|X^T X| \approx 0$  ← ill-condition

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

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

$\hat{y}_0 = (1 \cdot \tilde{x}_0^T) \hat{\beta}$  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.

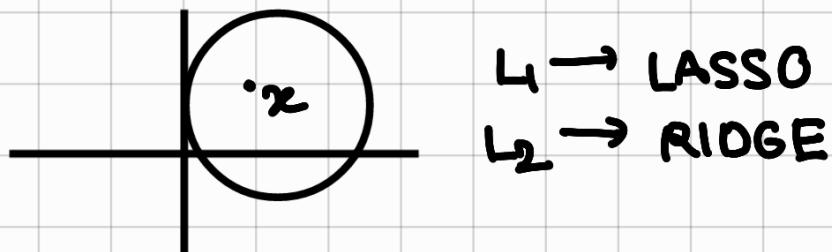
- 1)  $x^T x$  is a symmetric matrix.
- 2)  $(x^T x)^{-1}$  is also a symmetric matrix, if exists.
- 3)  $x^T x$  is positive semidefinite matrix.
- 4)  $x^T x$  has non-negative evs.  
 $\lambda_0 > \lambda_1 > \lambda_2 > \dots \geq 0$ .
- 5)  $\lambda_{\max} = \max_{\|z\| \neq 0} \frac{z^T (x^T x) z}{z^T z}$   
 $\lambda_{\min} = \min_{\|z\| \neq 0} \frac{z^T (x^T x) z}{z^T z}$
- 6)  $\text{tr}(x^T x) = \sum \lambda_i$
- 7)  $\text{tr}((x^T x)^{-1}) = \sum 1/\lambda_i$
- 8)  $|x^T x| = \prod_i \lambda_i$
- 9)  $|(x^T x)^{-1}| = \prod_i 1/\lambda_i \uparrow \infty \text{ as some } \lambda_i \downarrow 0$

Hence,

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

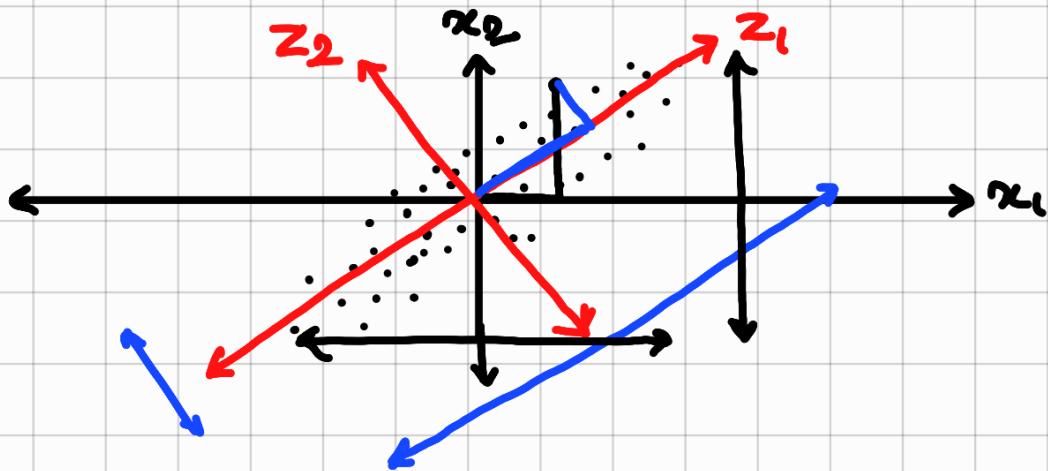
## Possible Solutions

- Principal Component Regression
- restrict variance in euclidean norm /  $\ell^2$ -norm



- Principal Component Regression

SVD  
math / PCA  
Statistics



- one method of dimension reduction

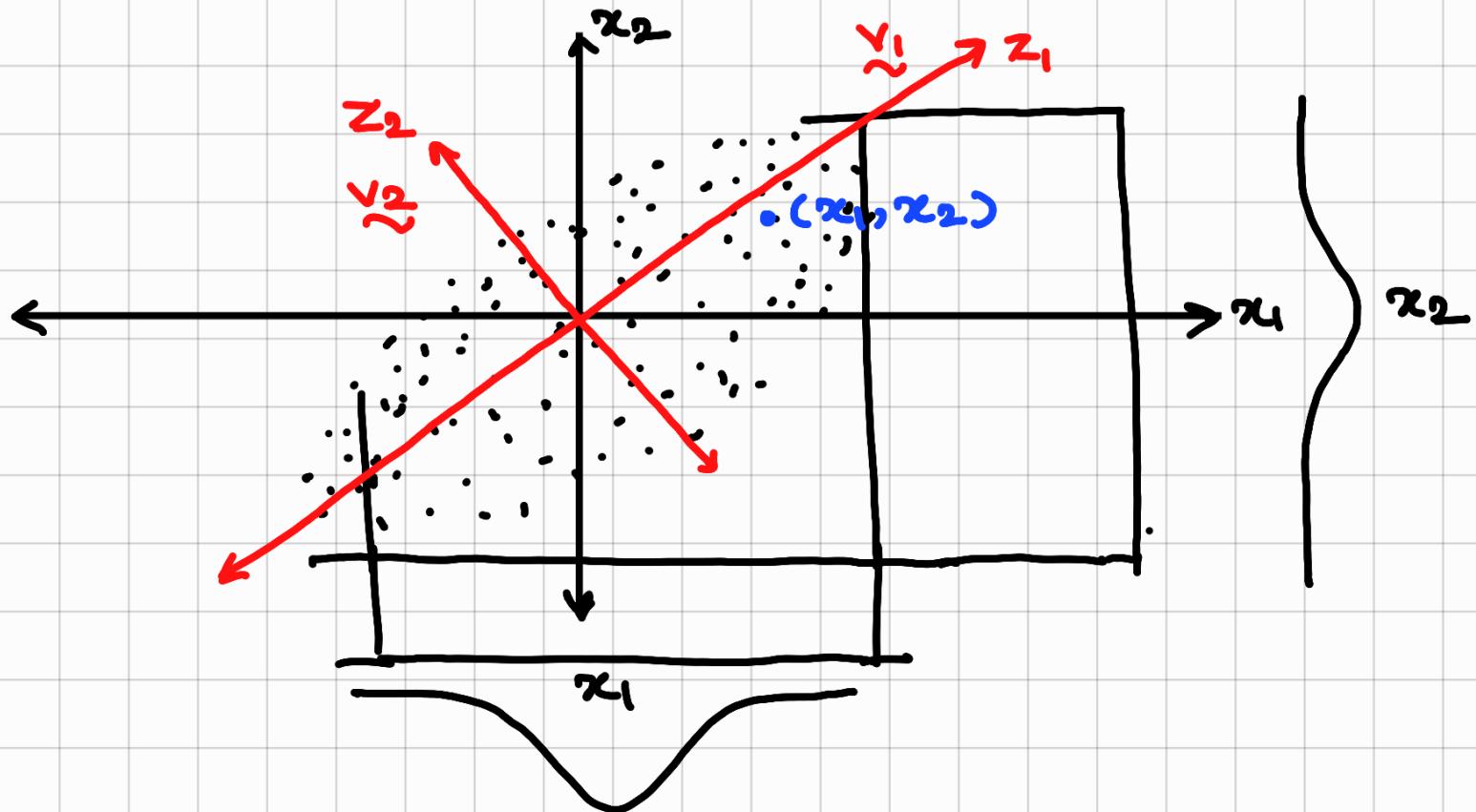
Tentatively, Midsem syllabus is upto and including Polynomial Regression.

## Principal Component Regression

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

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

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



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

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

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

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \langle x, v_1 \rangle \tilde{v}_1 + \langle x, v_2 \rangle \tilde{v}_2$$

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

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

$$[ \tilde{v}_0 \quad \tilde{v}_1 \quad \tilde{v}_2 \dots \dots \quad \tilde{v}_k ] = P$$

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

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

$\downarrow$   $\hookrightarrow$  orthogonal matrix  
diagonal matrix

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

Example of orthogonal matrix:  
 orth. matrix for reflection  $\rightarrow \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix}$

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

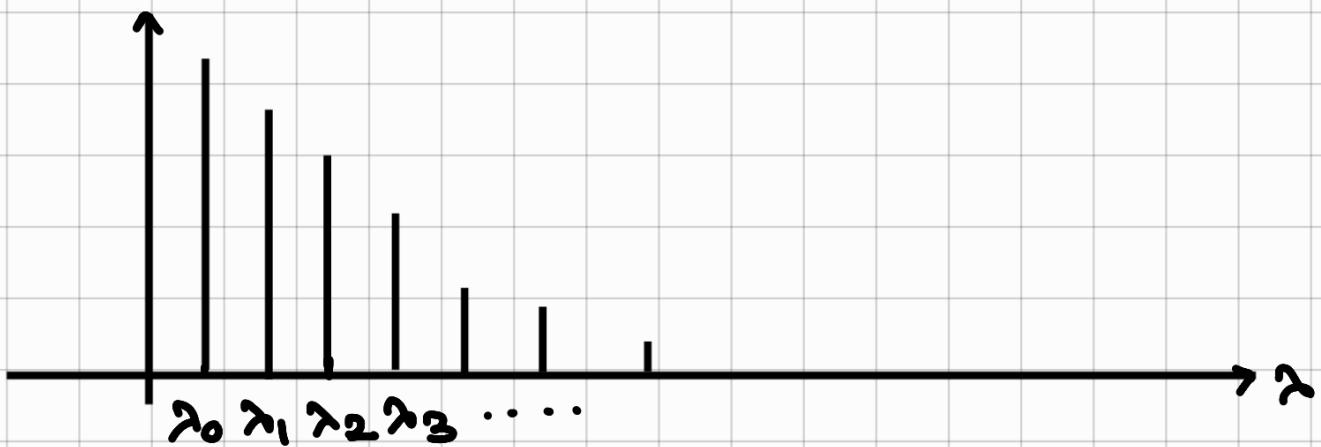
$$Y = X P P^T \tilde{B} + \epsilon \rightarrow \tilde{Y} = Z \tilde{\alpha} + \tilde{\epsilon}$$

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

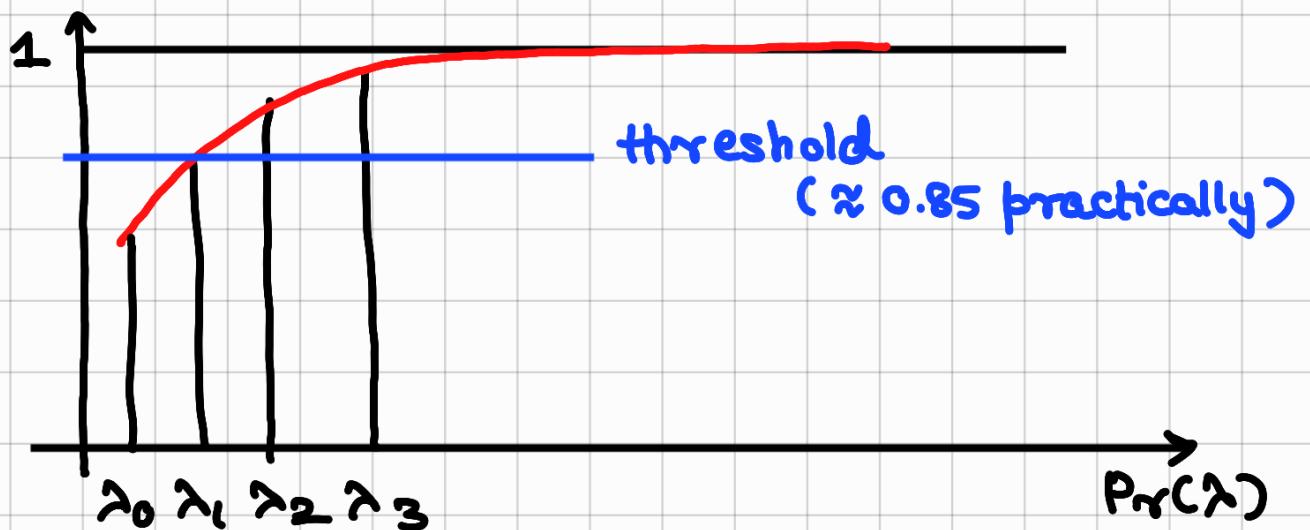
As a process,

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

(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

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

$$Z_{(r)} = X \begin{bmatrix} v_0 & v_1 & \dots & v_r \end{bmatrix}_{n \times (k+1)}$$

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

Transform  $\hat{\alpha}_{(r)}$  to  $\beta$ ,

$$\begin{aligned}\beta &= I\beta \\ \Rightarrow \beta &= P P^T \beta \\ \Rightarrow \beta &= P \cdot \hat{\alpha}_{(r)}\end{aligned}$$

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

Tentatively, Midsem syllabus is upto and including Polynomial Regression.

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

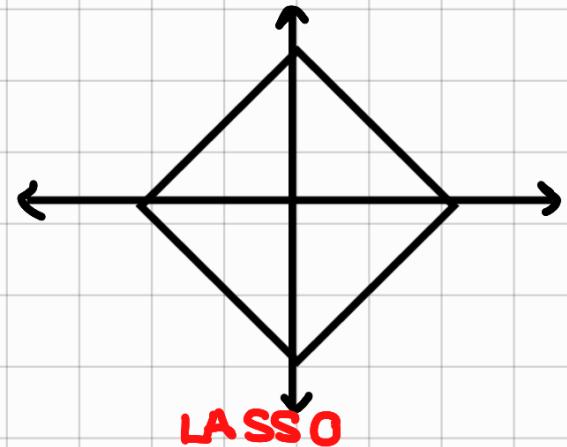
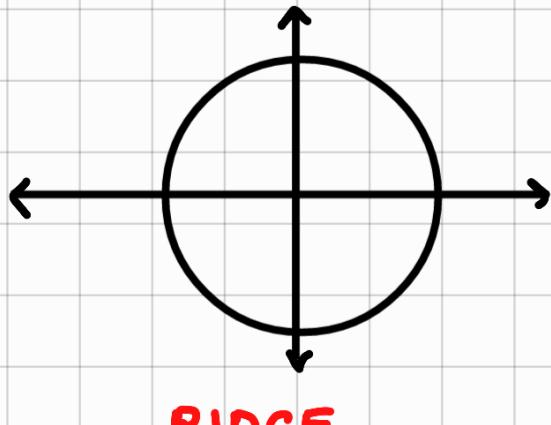
$\rightarrow$  PCR

$\rightarrow$  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  $\hat{\beta}$ .

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



## 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

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

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

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  $\beta$ ?

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

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

Notations :

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

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

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

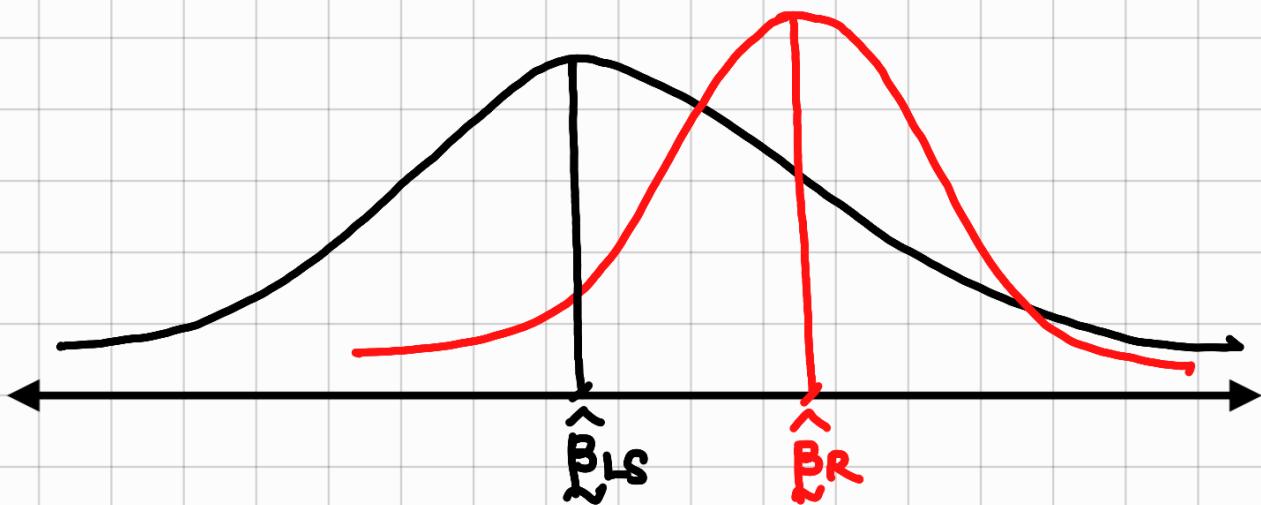
Does RIDGE estimate have less variance?

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

↓                    ↓                    ↓  
 psd                psd                psd

Hence,  $D(\hat{\beta}_{LS}) - D(\hat{\beta}_R)$  is a psd matrix.

This implies RIDGE estimate has less variance.



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

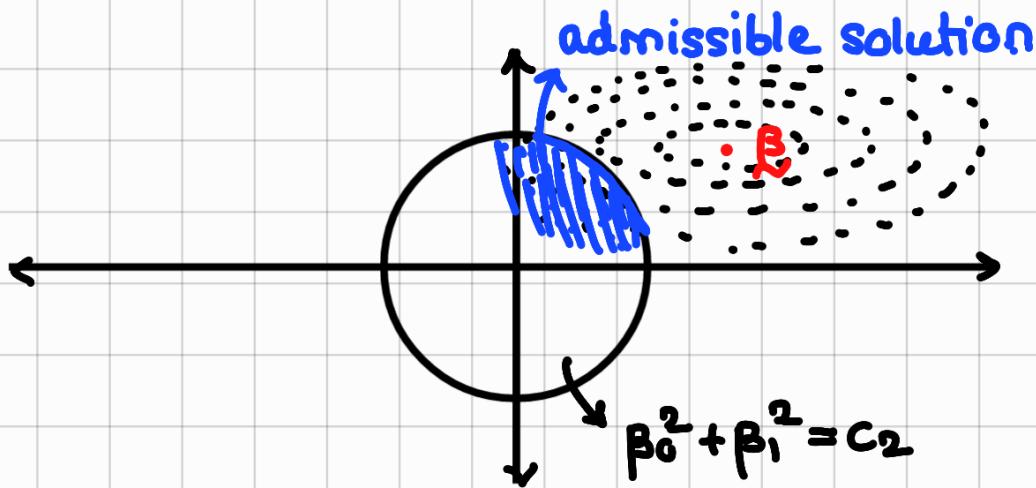
$$\text{MSE } (\hat{\beta}_R) = \text{tr}[D(\hat{\beta}_R)] + \hat{\beta}^T W^2 \hat{\beta} \lambda^2$$

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

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

$\lambda_i$ 's are eigenvalues of  $X^T X$ .

## Solution space for RIDGE Geometrical Intuition



Similar steps for LASSO (differentiation is a bit complex)

### Variable Selection

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

$(p+r=k+1)$

Data we received :  $(Y, x_p)$

We need  $r$  more features to fit data into true model

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

Reduced Model  $Y = x_p \hat{\beta}_p + \epsilon$  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  $\sigma_p^2$  (estimate value of  $\sigma^2$  in reduced model) an unbiased estimator of  $\sigma^2$ ?

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

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

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

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

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

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

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,

$$\Sigma^{-1} = \begin{pmatrix} \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} & \Sigma_{12} \Sigma_{22}^{-1} \\ \Sigma_{21} \Sigma_{22}^{-1} & \Sigma_{22} \end{pmatrix}^{-1}, \text{ and}$$
$$(\mathbf{A} + \mathbf{B}\mathbf{C}\mathbf{D})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{C}^{-1} + \mathbf{D}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{D}\mathbf{A}^{-1}$$

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 (analysis②  
of dataset)

Justification of  
maths + dataset  
description ③  
using things  
taught in  
class)

+ PDF of ④  
executed code  
(Python or R)

(10 mins)

30 marks  $\rightarrow$  video recording ④  
of project (Teams)

Upload all the stuff in a  
drive folder

} 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:  $y_i = \beta^T \tilde{x}_i + \epsilon_i$ ,  $i=1,2,\dots,n$

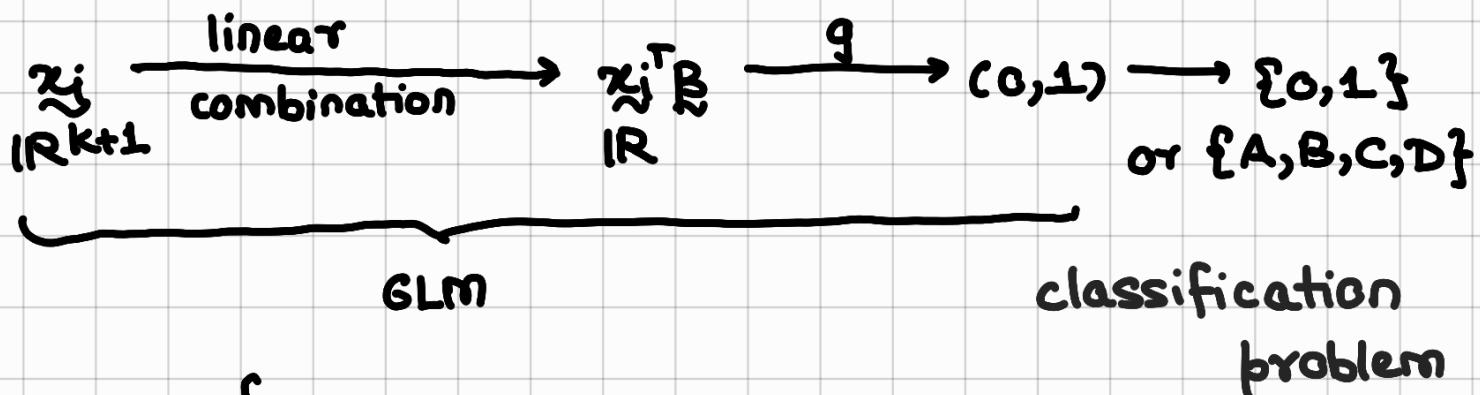
$$E(\epsilon_i) = 0 \quad V(\epsilon_i) = \sigma^2$$

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

$$E(y_i | \tilde{x}_i) = \beta^T \tilde{x}_i = I(\beta^T \tilde{x}_i) \quad (I(x) = x.)$$

For categorical variables as output, we need

$$E(y_i | \tilde{x}_i) = g(\beta^T \tilde{x}_i)$$



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

any cdf of a cont. random variable

## Logistic Regression

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

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

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

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

$$E(y_i | \mathbf{x}_i) = \pi_i = \frac{e^{\mathbf{x}_i^\top \beta}}{1 + e^{\mathbf{x}_i^\top \beta}}$$

Likelihood fn. of ( $\beta$ )

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

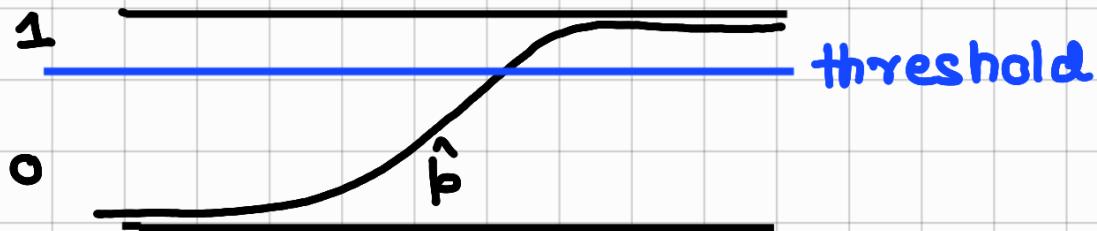
$$\log(L(\beta | y)) = \sum_{i=1}^n \left[ y_i \log\left(\frac{\pi_i}{1-\pi_i}\right) \right] + \sum_{i=1}^n \log(1-\pi_i)$$

$$L(\beta) = \sum_{i=1}^n [y_i(\mathbf{x}_i^\top \beta)] - \sum_{i=1}^n \log(e^{1+e^{\mathbf{x}_i^\top \beta}})$$

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

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

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



## Error Analysis of MLR

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

$$y_i = \mathbf{x}_i^\top \beta + \epsilon_i \quad \text{data}$$

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

$$H = P_X$$

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

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

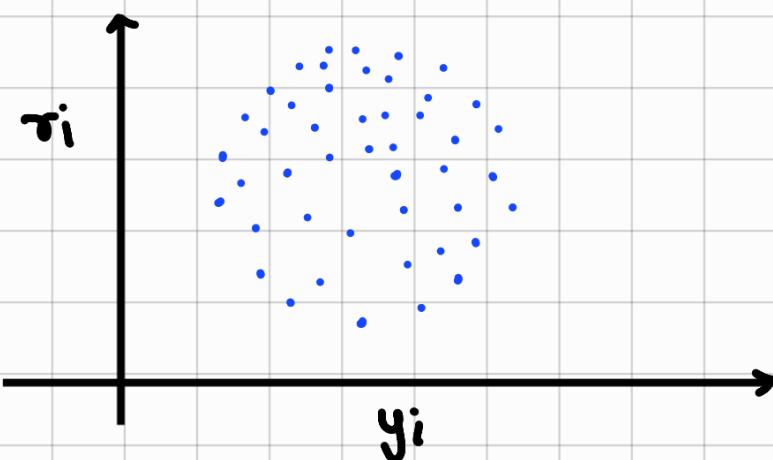
$$\text{cov}(e_i, e_j) = \begin{cases} \sigma^2(1-h_{ii}) & i=j \\ -\sigma^2 h_{ij} & i \neq j \end{cases}$$

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

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

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

If all the analysis is correct, the pattern resembles :



**Homework:** Show that

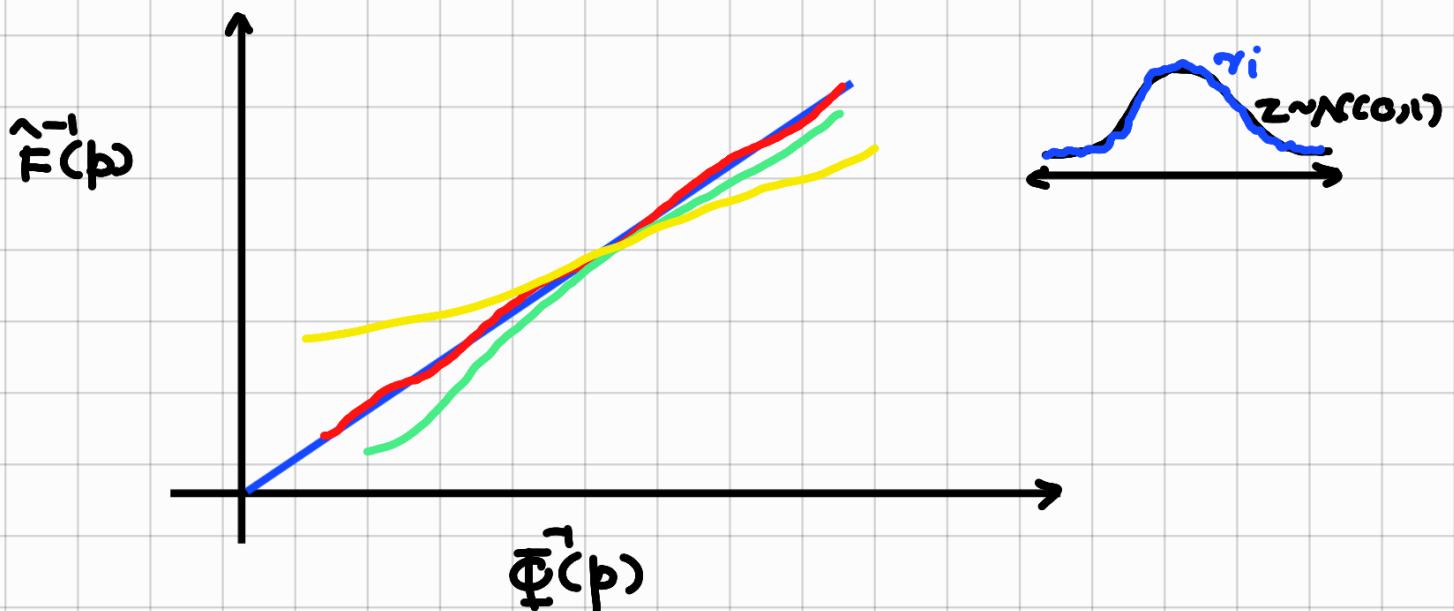
$$\tau_i = \frac{e_i}{\left( MSE_{\text{error}} \left[ 1 - \left( \frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{xx}} \right) \right] \right)^{1/2}} \text{ for SLR.}$$

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

To validate this,

- Q-Q plot
- $\chi^2$ -test
- KS test

## Q-Q Plot



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

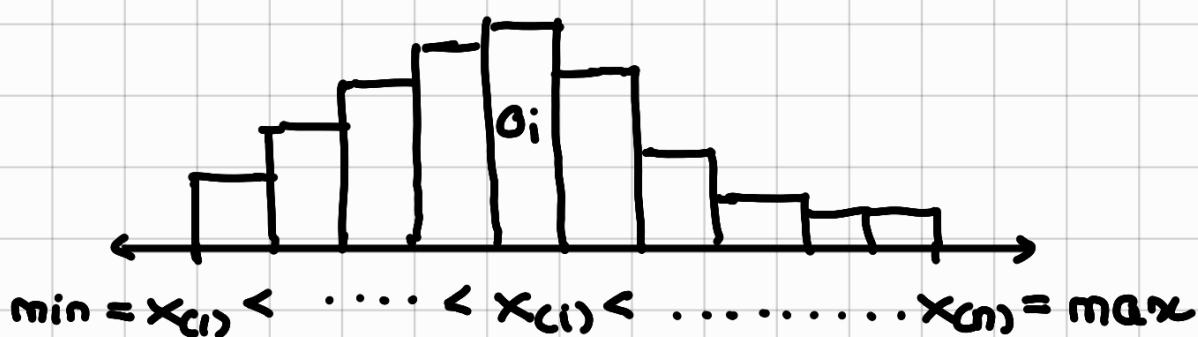
$\Phi^{-1}(p)$  cdf of  $N(0,1)$   
 $\hat{F}$  is the empirical cdf of plot

## $\chi^2$ -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
- $E_i = \text{expected no. of observations within } i\text{th interval}$ 
  - $= n \times p_i$
  - $= (\text{total obs}) \times \text{prob. of } i\text{th interval}$

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

$df = k-1$  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.

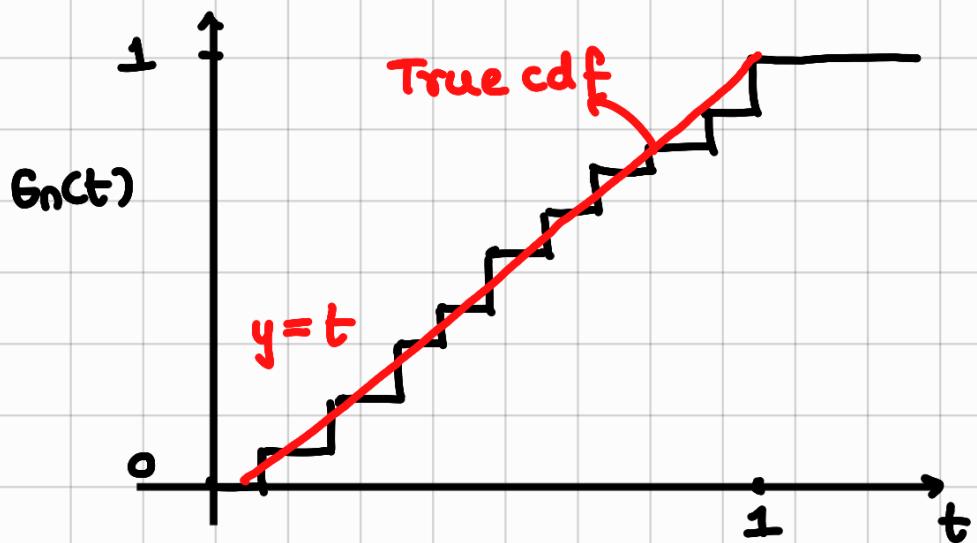
→ test only applicable to continuous r.v.

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

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

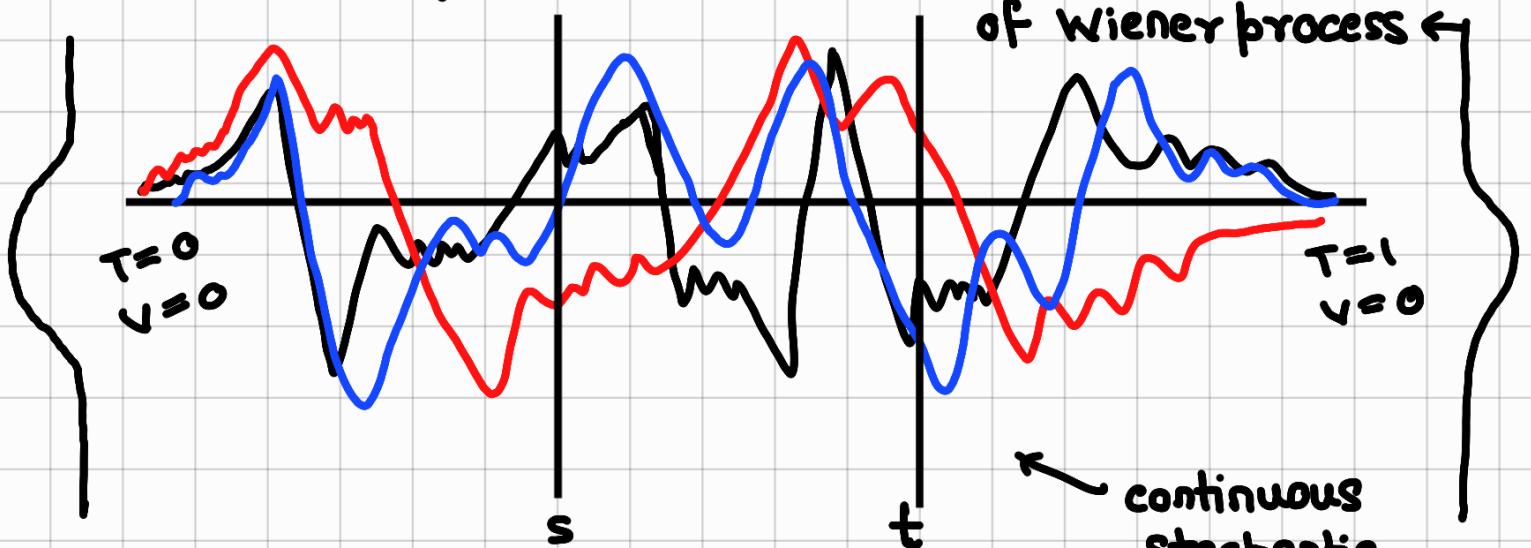
$D$  is said to follow Kolmogorov distribution as  $n \uparrow \infty$ .

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

Kolmogorov Distribution :

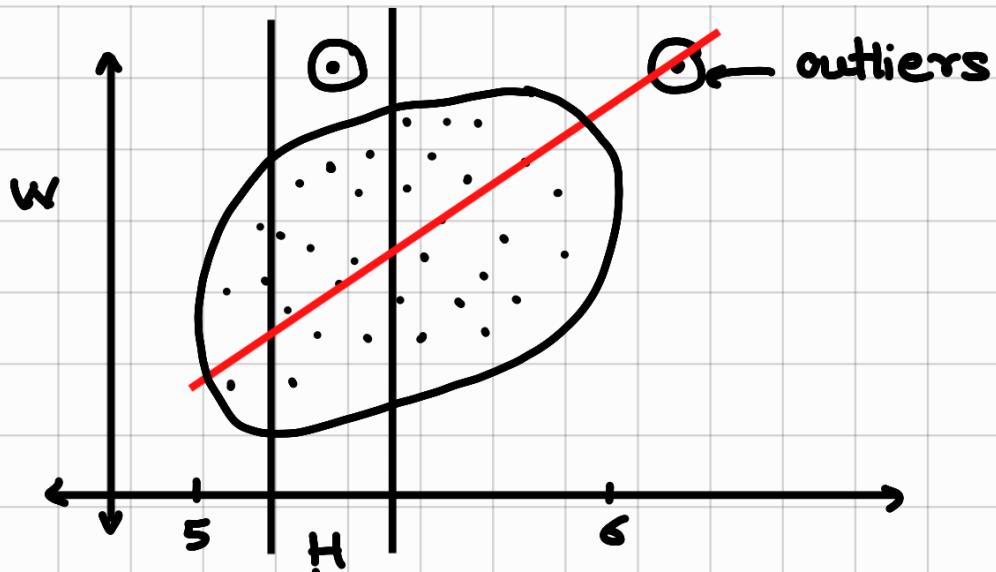
Brownian Bridge :

conditional distribution  
of Wiener process ↪



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

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



$$e_i = y_i - \hat{y}_i \\ = y_i - \tilde{x}_i^T \hat{\beta}$$

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

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

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

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

Hence, fitting the whole  $n$  observations serves the purpose.

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

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

$\hat{\sigma}^2 = \text{MSE}_{\text{Error}}$  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

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

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

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  $Y = [x]$

$$\begin{aligned} P(Y = r) &= P([x] = r) \\ &= P(r \leq x < r+1) \\ &= P(x < r+1) - P(x < r) \\ &= (1 - e^{-\lambda(r+1)}) - (1 - e^{-\lambda r}) \\ &= (e^{-\lambda})^r (1 - e^{-\lambda}) \\ &\approx q^r p \quad \leftarrow \text{geometric } (p = 1 - e^{-\lambda}) \\ &\qquad\qquad\qquad \rightarrow \text{discrete} \end{aligned}$$

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

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

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

(WN stands for white noise)

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

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

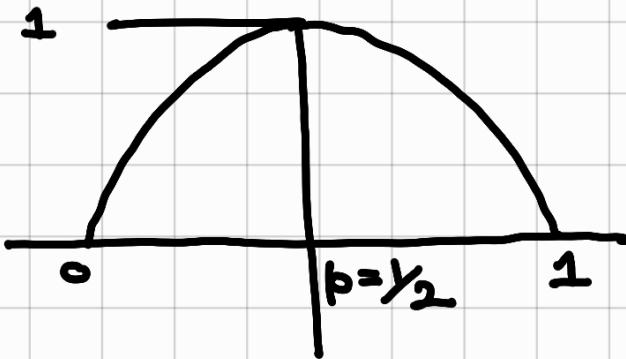
$$x_t = \begin{cases} +1 & \text{with prob. } p = \frac{1}{2} \\ -1 & \text{with prob. } 1-p = \frac{1}{2} \end{cases}$$

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

$$Var(x_t) = 4p(1-p)$$

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

$$x_t = 2z_t - 1$$



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



$$x_0 = 0$$

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

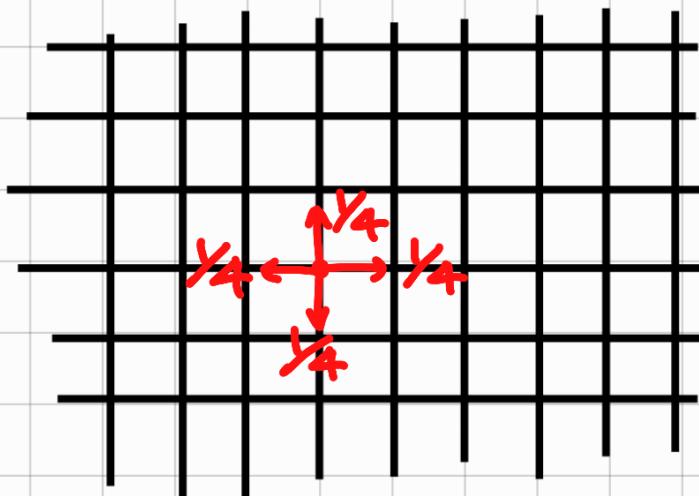
where  $w_t$  is a binary process with prob. ( $p$ ).

$\binom{2n}{n} p^n (1-p)^n$ : 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

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

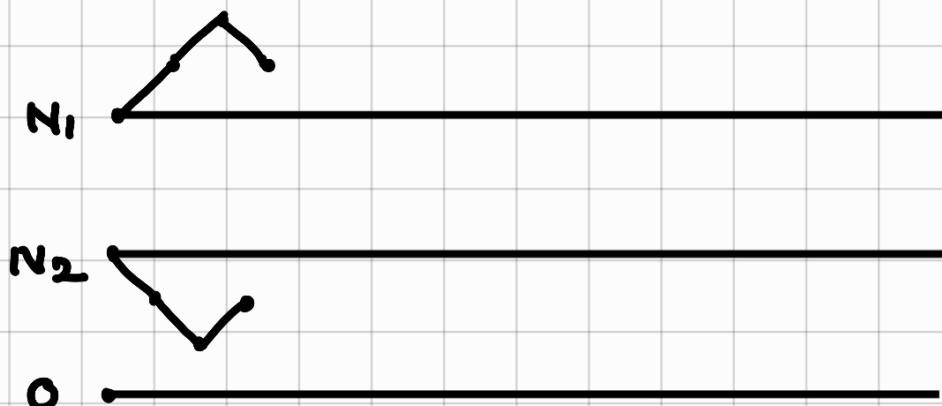
$$V(W_t) = \sigma_w^2$$

$$X_t = X_0 + \sum_{i=1}^t W_i$$

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

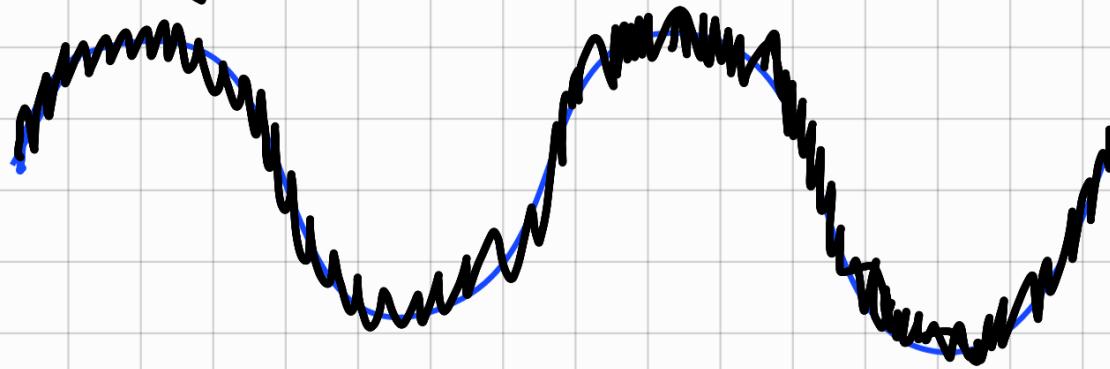
$$V(X_t) = \sigma^2 t$$

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)

MAC1 process

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

$$X_t = \alpha W_t + \beta W_{t-1}$$

## 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)

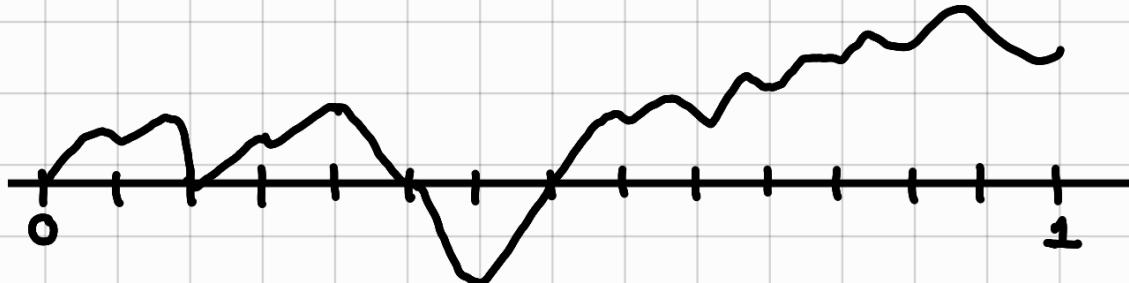
For definition, refer to the slides.

Generating data from Brownian motion (#)

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

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

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  $y_n(t) = \frac{1}{\sqrt{n}} \sum_{k=1}^{\lfloor nt \rfloor} w_k$  converges to Wiener process  $x_t$  on  $[0,1]$  for large  $n$ .

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

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

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

Brownian Bridge :

$$s < t \quad \begin{array}{c} \leftarrow \rightarrow \\ 0 \quad s \quad t \quad 1 \end{array}$$

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

$$\begin{array}{ccccccccc} \leftarrow & 0 & t_1 & t_2 & t_3 & \dots & \dots & t_{k-1} & t_k & 1 \end{array} \rightarrow$$

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



Description of Brownian motion :

$$t = 0, X_0 = 0 \quad \begin{array}{c} \overbrace{\hspace{1cm}}^{t=0} \quad \overbrace{\hspace{1cm}}^{t=t_1} \quad \overbrace{\hspace{1cm}}^{t=t_2} \end{array} \quad t_2 - t_1$$

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

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

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

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  $[0,1]$

$$(t=0, B_0(t)=0) \text{ and } (t=1, B_0(t)=0) \quad \begin{cases} B_0(t) = B(t) - tB(1) \\ 0 \leq t \leq 1 \end{cases}$$

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

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

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

## Brownian Bridge ( $B_0(t)$ )

→ KS Test for goodness of fit

→ Change point detection in temporal data

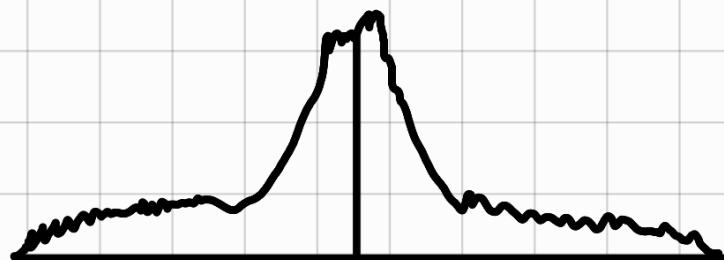
$$S_t := \frac{1}{\sqrt{n}} \sum_{i=1}^{[tn]} \frac{(x_i - \bar{x})}{\hat{\sigma}} ; 0 \leq t \leq 1 \quad \text{CUSUM statistic}$$

↪ estimated process variance

$$S_t \Rightarrow B_0(t)$$

$$|S(t)|$$

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



→ these statistics are such that software cannot directly estimate.

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

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

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

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

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

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

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

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

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

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

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



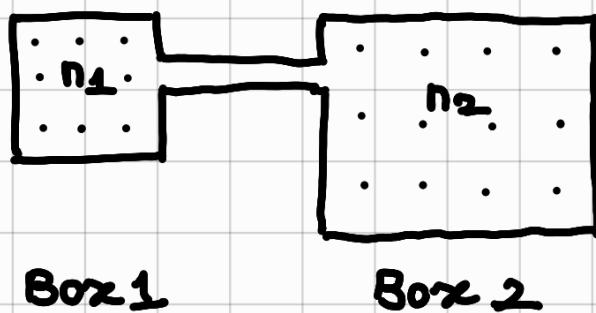
$$\text{cov} = \min(s, t) - st \\ = s(1-t)$$

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

## Stationary Timeseries

- strong stationarity
- weak stationarity

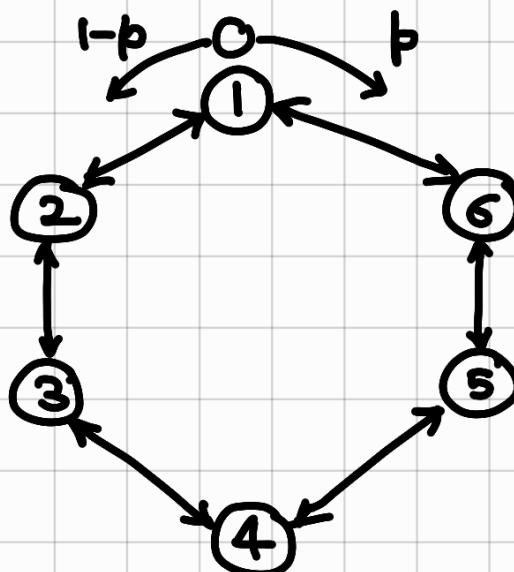
e.g. 1  $x_t \sim \text{bin}(n_1 + n_2, p_2)$  as  $t \rightarrow \infty$



$$x_0 = n_1$$

$x_t = \# \text{ balls on box 1 at time } t$

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.

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

'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  $\{x_t\}$   $\overset{\text{WN or iid}}{\sim}$   $E(x_t) = 0, \text{Var}(x_t) = \sigma^2$

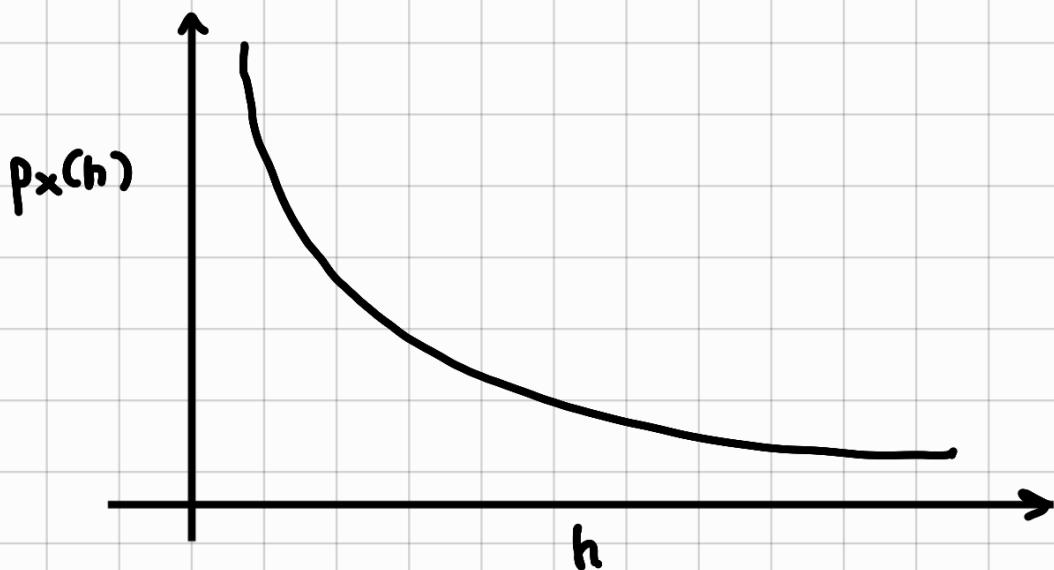
$$\gamma_x(h) = \begin{cases} \sigma^2 & h=0 \\ 0 & h \neq 0 \end{cases}$$

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

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)

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

$\phi = 0$   
Zt is weakly stationary

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

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

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

$$V(z_t) = E(z_t^2) = E(\varphi^2 z_{t-1}^2 + w_t^2 + 2\varphi z_{t-1} w_t)$$

$$= \varphi^2 E(z_{t-1}^2) + E(w_t^2) + 0$$

$$= \varphi^2 E(z_{t-1}^2) + \sigma^2$$

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

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

$$E(z_t) = 0$$

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

→ this is similar to BWN case.

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

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

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

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

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

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

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

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

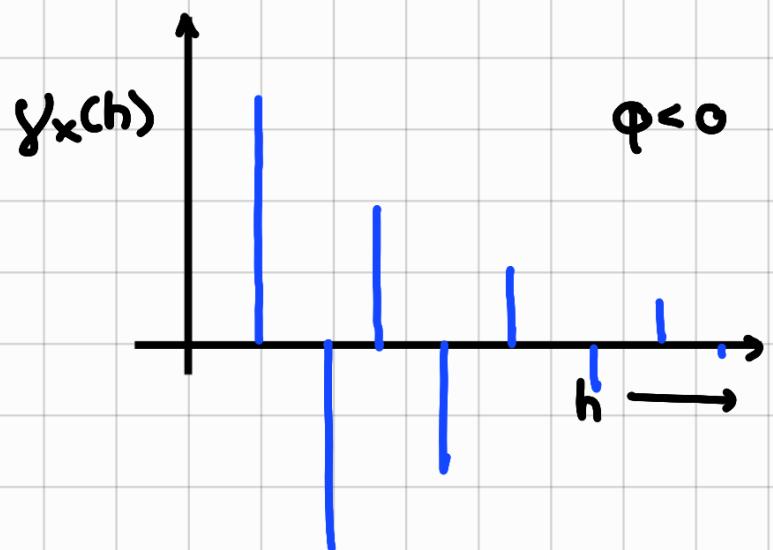
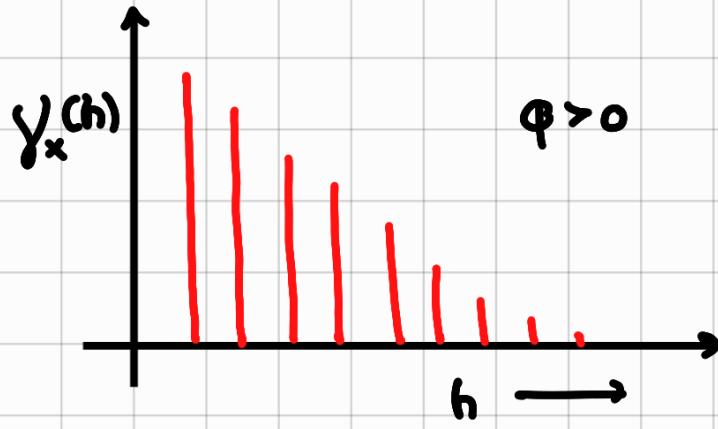
....

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

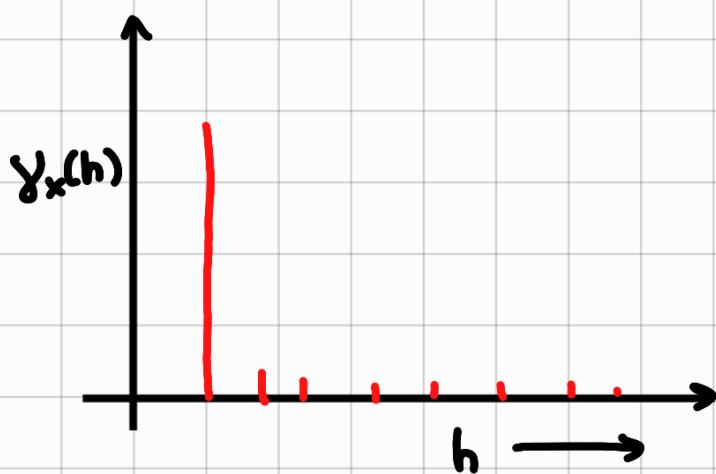
$$= \varphi^h \left( \frac{\sigma^2}{1-\varphi^2} \right) = \gamma_z(-h)$$

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 ( $\nabla$ )

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

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

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

3. Seasonal Operator

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

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

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

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

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

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

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

	$\nabla$	$\nabla^2$	$\nabla^3$
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 :

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

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

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

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

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

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

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

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

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

Note. Absolutely summable series is always summable

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

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

Ez 2.  $x_t \sim \text{MAC}(1)$

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

Ez 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

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

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

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

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

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

Autoregressive 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

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

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

$$\Rightarrow (x_t) - (\sum_{j=1}^k \varphi_j b_j) x_t = z_t.$$

$$\Rightarrow (I - \sum_{j=1}^k \varphi_j b_j) x_t = z_t.$$

Let  $\bar{\Phi}_B(u) = (I - \sum_{j=1}^k \varphi_j b_j)$ , then

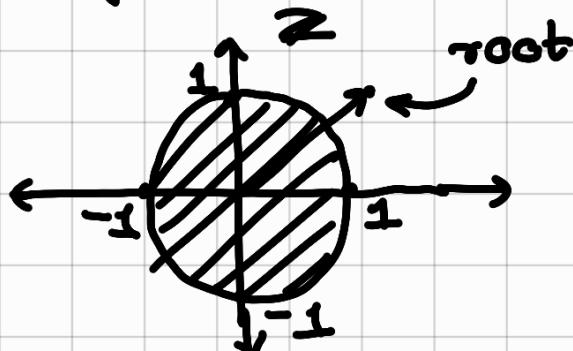
$$\bar{\Phi}_B(B) x_t = z_t$$

$$\Rightarrow x_t = \left(\frac{1}{\bar{\Phi}_B(B)}\right) z_t. \quad \rightarrow \textcircled{*}$$

$\frac{1}{\bar{\Phi}_B(B)} = \frac{1}{1 - \sum_{j=1}^k \varphi_j b_j}$ . This should satisfy

unit root condition for  $\textcircled{*}$  to be feasible, that is, for eqn.  $1 - g u \equiv 0$ , root  $u = \sqrt{g} > 1$ .  
 (for our case,  $g u \equiv \sum_{j=1}^k \varphi_j b_j$ )

Also,  $|g| < 1$ ,  $g \neq 0$ .



AR(1) process can be thought as MA( $\infty$ ) process.

$$AR(1) = \lim_{q \rightarrow \infty} MA(q).$$

We can represent AR(1) process as MA( $\infty$ ) process

$$x_t = \varphi x_{t-1} + z_t \quad |\varphi| < 1 \quad \varphi \neq 0$$

$$\rightarrow (I - \varphi B) x_t = z_t$$

$$\rightarrow x_t = \left(\frac{1}{1 - \varphi B}\right) z_t$$

$$\Rightarrow x_t = \left(\sum_{j=0}^{\infty} (\varphi B)^j\right) z_t$$

$$\Rightarrow x_t = \sum_{j=0}^{\infty} \varphi^j z_{t-j}$$

$$x_t = (\sum_{j=0}^k \varphi^j z_{t-j}) + \varphi^{k+1} x_{t-(k+1)}$$

$$E(x_t - \sum_{j=0}^k \varphi^j z_{t-j})^2 = E\{\varphi^{k+1} x_{t-(k+1)}\}^2$$

$$\lim_{k \rightarrow \infty} E(x_t - \sum_{j=0}^k \varphi^j z_{t-j}) = \lim_{k \rightarrow \infty} E\{\varphi^{k+1} x_{t-(k+1)}\}^2$$

$$= \lim_{k \rightarrow \infty} (\varphi^{2k+2}) E(x_{t-(k+1)}^2)$$

( bounded )

$$\left( \lim_{n \rightarrow \infty} E(|x_n - y|^2) \rightarrow 0 \Rightarrow x_n \xrightarrow{d} y. \right)$$

Hence,  $x_t = \sum_{j=0}^{\infty} \varphi^j z_{t-j}$

$$\lim_{n \uparrow \infty} (x_t - \sum_{j=0}^k \varphi^j z_{t-j})^2 = \lim_{k \uparrow \infty} \{ \varphi^{2k} \underbrace{E(x_{t-k}^2)}_{\text{bounded}} \}$$

↓  
AR(1)  
↓  
MA(k)

$$y_k = \sum_{j=0}^k \varphi^j z_{t-j} \text{ such that } E|x_t - y_k|^2 \rightarrow 0$$

$$\Rightarrow y_k \xrightarrow{d} x_t \text{ as } k \uparrow \infty.$$

Note.  $x_1, x_2, \dots, x_n$  iid  $E(x) = \mu$   $V(x) = \sigma^2$

$$T_k = \sqrt{k}(\bar{x} - \mu)/\sigma$$

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

$$\text{Then } E(T_k - Z)^2 \rightarrow 0 \Rightarrow T_k \xrightarrow{d} Z$$

### 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

$$y_t = \sum_{j=-\infty}^{\infty} \psi_j x_{t-j}, \text{ we have}$$

$$\rightarrow E(y_t) = 0$$

$$\rightarrow \gamma_{y(h)} = \sum_{j=-\infty}^{\infty} \sum_{k=-\infty}^{\infty} \psi_j \psi_k \gamma_x(h+k-j)$$

**Proof.**

$$E|Y_t| = E\left|\sum_{j=-\infty}^{\infty} \psi_j x_{t-j}\right| \leq \sum_{j=-\infty}^{\infty} |\psi_j| E|x_{t-j}| \leq M_1 \sum_{j=-\infty}^{\infty} |\psi_j| < \infty$$

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

$$(x_t - \sum_{j=1}^p \phi_j x_{t-j}) = z_t + \sum_{j=1}^q \theta_j z_{t-j} \quad \forall t \in \mathbb{Z}, z \sim WN.$$

$$\begin{aligned} &\rightarrow (I - \sum_{j=1}^p \phi_j B^j) x_t = (I + \sum_{j=1}^q \theta_j B^j) z_t \\ &\rightarrow \Phi_p(B) x_t = \Theta_q(B) z_t \text{ where} \end{aligned}$$

$$\Phi_p(z) = 1 - \sum_{j=1}^p \phi_j z^j \text{ and } \Theta_q(z) = 1 + \sum_{j=1}^q \theta_j z^j$$

$\Phi_p(z)$  and  $\Theta_q(z)$  should not have any common root.

For example,  $ARMA(1,1)$

$$\begin{aligned} x_t - \phi x_{t-1} &= z_t + \theta z_{t-1} \\ \Rightarrow x_t &= \left( \frac{I + \theta B}{I - \phi B} \right) z_t \end{aligned} \quad \begin{cases} |\phi| < 1, \phi \neq 0 \\ \theta + \phi \neq 0 \end{cases}$$

$$\begin{aligned} x_t &= \left\{ (I + \theta B) \left( \sum_{j=0}^{\infty} (\phi B)^j \right) \right\} z_t = \Psi(B) z_t \\ &= (\psi_0 + \sum_{j=1}^{\infty} \psi_j B^j) z_t \end{aligned}$$

$$\psi_0 = 1$$

$$\psi_j = \begin{cases} (\theta + \phi) \phi^{j-1} & \forall j \geq 1 \end{cases}$$

Linear process representation

$$x_t = z_t + (\theta + \phi) \sum_{j=1}^{\infty} \phi^{j-1} z_{t-j}$$

$$\begin{aligned} \gamma_x(0) &= \sigma^2 \sum_{j=0}^{\infty} \psi_j^2 = \sigma^2 [1 + (\theta + \phi)^2 \sum_{j=1}^{\infty} \phi^{j-2}] \\ &= \sigma^2 \left[ 1 + \frac{(\theta + \phi)^2}{1 - \phi^2} \right] \end{aligned}$$

$$\gamma_x(1) = \sigma^2 \sum_{j=0}^{\infty} \psi_j \psi_{j+1} = \sigma^2 [(\theta+\phi)^2 + \frac{(\theta+\phi)^2 \phi}{1-\phi^2}]$$

$$\gamma_x(h) = \phi^{h-1} \gamma_x(1). \text{ Homework.}$$

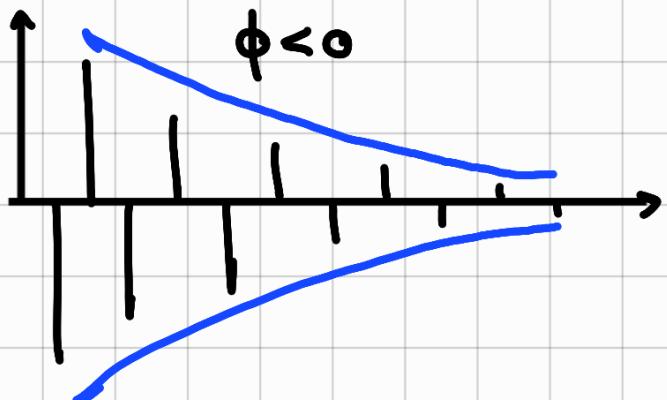
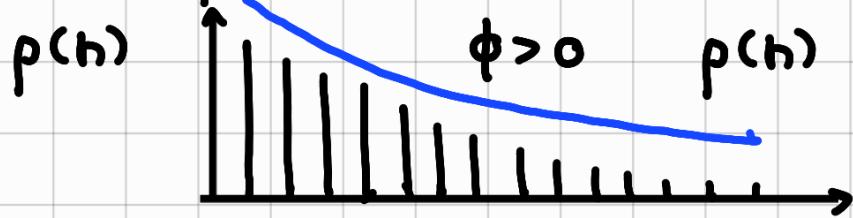
Autocorrelation coefficient of weakly stationary process.

$$p(h) = \frac{\gamma_x(h)}{\gamma_x(0)}$$

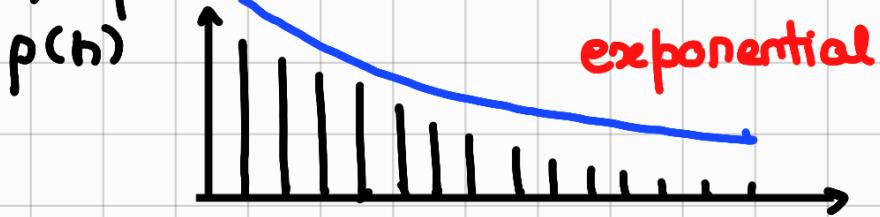
MA(q) process :



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,  $\begin{pmatrix} Y \\ Z \\ \vdots \\ X \end{pmatrix} \sim N \left( \begin{pmatrix} \mu_Y \\ \mu_Z \\ \vdots \\ \mu_X \end{pmatrix}, \begin{pmatrix} \sum_{YZ} & \sum_{YZX}^T \\ \sum_{YZX} & \sum_X \end{pmatrix} \right)$

$Y \in \mathbb{R}^1$   
 $Z \in \mathbb{R}^1$   
 $X \in \mathbb{R}^p$

$$E(Y) = \mu_Y \quad E(Z) = \mu_Z \quad E(X) = \mu_X$$

$$\Sigma_{YZ} = \begin{pmatrix} \sigma_{YY} & \sigma_{YZ} \\ \sigma_{YZ} & \sigma_{ZZ} \end{pmatrix}_{2 \times 2} \quad \Sigma_X = \sum_X = ((\sigma_{X_i X_j}))_{p \times p}$$

$$\sigma_{YX} = \text{cov}(Y, X) = \begin{pmatrix} \sigma_{YX_1} \\ \sigma_{YX_2} \\ \vdots \\ \sigma_{YX_p} \end{pmatrix}$$

$$\sigma_{ZX} = \text{cov}(Z, X) = \begin{pmatrix} \sigma_{ZX_1} \\ \sigma_{ZX_2} \\ \vdots \\ \sigma_{ZX_p} \end{pmatrix}$$

$$\Sigma_{YZX} = (\Sigma_{YX}, \Sigma_{ZX})_{p \times 2} \quad \Sigma_{YZX}^T = \begin{pmatrix} \Sigma_{YX}^T \\ \Sigma_{ZX}^T \end{pmatrix}_{2 \times p}$$

Partial correlation coefficient is the correlation b/w  $Y$  and  $Z$  after removing the effect of  $X$  and is denoted as  $\rho_{YZ \cdot X}$

$$\rho_{YZ \cdot X} = \frac{\text{cov}(Y - E(Y|X), Z - E(Z|X))}{\sqrt{\text{var}(Y - E(Y|X)) \text{var}(Z - E(Z|X))}}$$

$$\Rightarrow \rho_{YZ \cdot X} = \frac{\text{cov}(e_{Y|X}, e_{Z|X})}{\sqrt{\text{var}(e_{Y|X}) \text{var}(e_{Z|X})}}$$

$$= \frac{\sigma_{yz|x}}{\sqrt{\sigma_{yy|x} \sigma_{zz|x}}}$$

When we use it for time series we call it PACF.

conditional distribution of  $(\begin{matrix} y \\ z \end{matrix}) | x = \tilde{x}$ .

This also follows normal distribution,

$$N \left( \begin{pmatrix} \mu_y \\ \mu_z \end{pmatrix}_{2 \times 1} + \sum_{yzx}^T \sum_x^{-1} (\tilde{x} - \tilde{\mu}_x), \sum_{yz} - \sum_{yzx}^T \sum_x^{-1} \sum_{yzx} \right)_{2 \times p \quad p \times p \quad p \times 1 \quad 2 \times 2 \quad 2 \times p \quad p \times p \quad p \times 2}$$

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

$\bar{x}$ , even when  $(\frac{y}{x})$  are not jointly normally distributed  
but 2nd order moments exists.