



## Applied Statistics Notes

Applied statistics (Politecnico di Milano)



Scansiona per aprire su Studocu

**G**

**T**

**G**

**T**

**5**

**E**

**5**

**N**

**O**

**N**

**O**

**O**

**S**

**O**

**G**

**T**

**G**

**T**

**5**

**E**

**5**

**N**

**O**

**N**

**O**

**O**

**S**

**O**

**G**

**T**

**G**

**T**

**5**

**E**

**5**

## Generalities of statistical learnings

### \* Dataframes

	$x_1$	$x_2$	$x_3$	...	$x_p$	$y$
Unit 1	$x_{11}$	$x_{12}$	$x_{13}$	...	$x_{1p}$	$y_1$
Unit 2	$x_{21}$	$x_{22}$	$x_{23}$	...	$x_{2p}$	$y_2$
Unit 3	$x_{31}$	$x_{32}$	$x_{33}$	...	$x_{3p}$	$y_3$
:	:					:
Unit $n$	$x_{n1}$	$x_{n2}$	$x_{n3}$	...	$x_{np}$	$y_n$

Here  $\underline{x} = (x_1, x_2, \dots, x_p) \in \mathbb{R}^p$

$y \in \mathbb{R}$

Goal: To explain the variability of  $y$  in terms of  $\underline{x}$

Basic model:  $y = f(\underline{x}) + \epsilon$

Best solution:  $f = \min E[(y - f(\underline{x}))^2]$   
(min: expected square errors)

Answer to the optimal model

If  $g$  is a function such that  $g: \mathbb{R}^p \rightarrow \mathbb{R}$

$$E[(y - g(\underline{x}))^2] = E[y - E[y(\underline{x})]^2] + E[(E[y(\underline{x})])^2] - E[y^2]$$

Best  $f \Rightarrow f(\underline{x}) = E[y|\underline{x}]$   
 ↗ conditional mean of  $y$   
 given  $\underline{x}$   
 (also called regression function)

To find the best  $f$ , we try to approximate  $E[y|\underline{x}]$   
 by means of a member of a family.

$$\{f_\theta \mid f_\theta: \mathbb{R}^p \rightarrow \mathbb{R} \text{ known when } \theta \in \Theta \subseteq \mathbb{R}^m \text{ is known}\}$$

$\theta$  is a multi-dimensional point in space

Problem

To find  $\hat{\theta}$  such that

$$\hat{\theta} = \min_{\theta} E[(y - f_\theta(\underline{x}))^2]$$

Once it is found then

$\Rightarrow \hat{f}_\theta$  should be taken as approx of  $E[y|\underline{x}]$

### \* Nomenclature in $y = f(\underline{x}) + \epsilon$

- $\underline{x} \in x_1, \dots, x_p \rightarrow$  features/predictors/regressors /inputs
- $y \in \mathbb{R} \rightarrow$  response, target, output
- $\epsilon \rightarrow$  errors, residuals
- using data to find best  $f \rightarrow$  model fitting
- $\hat{f} \rightarrow$  fitted model with best fit & parameters  $\hat{\theta}$

### \* Linear Models

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

$\underbrace{\epsilon}_{E[y|\underline{x}]}$  residual error term

$\underline{x} = (x_1, \dots, x_p) \in \mathbb{R}^p$

### \* Classical Linear Model

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i \quad i=1, \dots, n$$

here  $\epsilon_i$  follows a normal distribution

$$\epsilon_i \sim N(0, \sigma^2) \quad i=1, \dots, n \quad (\text{same variability})$$

The above model can be mentioned in the form of a design matrix of shape  $n \times p$  such that

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

$$\text{where } \bar{Y} = (y_1, y_2, \dots, y_n)'$$

$$\bar{\beta} = (\beta_0, \beta_1, \dots, \beta_p)'$$

$$\bar{\epsilon} = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)'$$

**Remark**

Here  $\beta_0, \beta_1$  in true notation

recover  $\star$  by setting  $x_{i1} = 1$

Data is in the format  
 $(n \times p)$  design matrix

$x_{11} \dots x_{1p}$	$y_1$
$x_{21} \dots x_{2p}$	$y_2$
$\vdots$	$\vdots$
$x_{n1} \dots x_{np}$	$y_n$

### Estimating the parameters

#### • Method 1: Ordinary Least Squares

Find  $\hat{\beta}$  that makes the values closer to  $\bar{Y} = \bar{X} \bar{\beta} + \bar{\epsilon}$

$$\hat{\beta} = \min_{\beta \in \mathbb{R}^p} \|\bar{Y} - \bar{X} \bar{\beta}\|^2 = \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n (y_i - x_i^\top \beta)^2$$

where  $x_i^\top \beta = \beta_0 x_{i1} + \dots + \beta_p x_{ip}$

### Geometry of OLS

#### ① $\bar{X} \bar{\beta}$

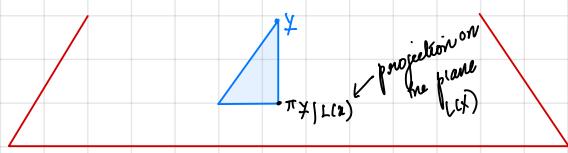
$$\bar{X} = [x^1, \dots, x^p]$$

$$x^{(j)} = \begin{bmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{nj} \end{bmatrix} \text{ where } j = 1, \dots, p$$

$$\therefore \bar{X} \bar{\beta} = \beta_0 x^{(1)} + \beta_1 x^{(2)} + \dots + \beta_p x^{(p)}$$

If  $L(\underline{x})$  is the linear space generated by the design matrix

$$L(\underline{x}) = \text{Span}(x^{(1)}, \dots, x^{(p)})$$



$\pi_{\gamma|L(x)}$  is an element of the space  $L(x)$

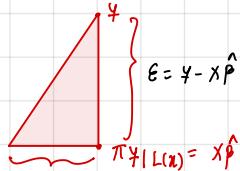
since  $\pi_{\gamma|L(x)} \rightarrow L(x)$

$$\pi_{\gamma|L(x)} = \hat{\beta}^T x + \dots + \hat{\beta}^T p = x \hat{\beta}$$

so

$\hat{\beta}$  is a vector of coefficient of  $\pi_{\gamma|L(x)}$

In a space,  $\gamma$  can be described as the sum of 2 vectors



$$\begin{aligned} \therefore \gamma &= x \hat{\beta} + \hat{\epsilon} \\ &= \hat{\gamma} + \hat{\epsilon} \\ \text{fixed } &\quad \text{residuals} \end{aligned}$$

Assume that  $x^{(1)}, \dots, x^{(n)}$  are linearly independent

$\dim(L(x)) = p$  (dimension of space of the model)

$\dim(L^\perp(x)) = n-p$  (dimension of space of the residuals)

Prop: If  $x$  is full rank means that the  $n \times p$  such that  $(n \times p)$

$$\hat{\beta} = (x'x)^{-1}x'\gamma$$

if collinear  $x'x$  will be = 0

$$\hat{\gamma} = x(x'x)^{-1}x'\gamma = H\gamma$$

hat matrix

such that  $H = x(x'x)^{-1}x'$

$$\hat{\epsilon} = \gamma - x \hat{\beta} = \gamma - \hat{\gamma} = \gamma - H\gamma = (I - H)\gamma$$

while estimating the variability for  $\epsilon$

$$\text{the variability } \hat{\sigma}^2 = \frac{1}{n-p} \sum_{i=1}^n \hat{\epsilon}_i^2$$

because  $\hat{\epsilon}$  lies in a plane  $(n \times p)$

Properties of  $\hat{\beta}$  &  $\hat{\sigma}^2$

① They are both unbiased such that

$$E[\hat{\beta}] = \beta$$

$$E[\hat{\sigma}^2] = \sigma^2$$

Minimum likelihood Estimator (MLE) for  $\beta$  &  $\sigma^2$

$$\gamma = x\beta + \epsilon$$

$$\epsilon \sim N(0, \sigma^2 I)$$

Under this assumption

$$\gamma \sim N(x\beta, \sigma^2 I) \quad (\text{gaussian})$$

This document is available free of charge on

Then the density of the gaussian distribution

$$f_{\beta, \sigma^2}(y_1, \dots, y_n) = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - x_i \beta)^2}$$

Likelihood

$$L(\beta, \sigma^2 | y_1, \dots, y_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} (y_i - x_i \beta)^2}$$

MLE for  $\beta$  &  $\sigma^2$

$$(\hat{\beta}, \hat{\sigma}^2) = \max_{\beta, \sigma^2} L(\beta, \sigma^2 | y)$$

$$= \max_{\beta, \sigma^2} \log L = \max_{\beta, \sigma^2} \ell(\beta, \sigma^2 | y)$$

$$= \max \left[ -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - x_i \beta)^2 \right]$$

Proposition

If  $x$  is full rank

$$\hat{\beta}_{MLE} = (x'x)^{-1}x'\gamma \Rightarrow E[\hat{\beta}_{MLE}] = \beta$$

$$\hat{\sigma}^2_{MLE} = \frac{1}{n-p} \sum \hat{\epsilon}_i^2 \Rightarrow E[\hat{\sigma}^2_{MLE}] = \frac{n-p}{n} \sigma^2$$

Restricted MLE for  $\sigma^2$

Notice:  $\gamma \sim N_p(x\beta, \sigma^2 I)$

$$\hat{\beta} = (x'x)^{-1}x'\gamma \sim N_p(\beta, \sigma^2 (x'x)^{-1})$$

Recall: If we have a variable  $w$

$$w \sim N(\mu, \Sigma)$$

& we take a linear transformation  $A$  ( $q \times n$ )

$$Aw \sim N_q(A\mu, A\Sigma A')$$

(if you linearly transform gaussian, you get gaussian)

if  $\hat{\beta} = (x'x)^{-1}x'\gamma \sim N_p(\beta, \sigma^2 (x'x)^{-1})$   
then

$$\hat{\epsilon} = (I - H)\gamma \sim N_n(0, \sigma^2 (I - H))$$

$\therefore$  RMSE mean

$$\hat{\sigma}^2 = \max_{\sigma^2} L(\sigma^2 | \hat{\epsilon})$$

$$\Rightarrow \hat{\sigma}^2 = \frac{1}{n-p} \sum_{i=1}^n \hat{\epsilon}_i^2$$

Inferences

$$\begin{aligned} \text{Var}(\hat{\beta}) &= \text{Cov}(\hat{\beta}) = \sigma^2 (x'x)^{-1} \\ \text{Var}(\hat{\beta}) &= \frac{\hat{\sigma}^2}{n-p} (x'x)^{-1} \end{aligned}$$

\* Methods to run test on linear combinations of models

If  $L$  is a  $q \times p$  matrix

Inference:  $L\hat{\beta}$

$$F = (L\hat{\beta} - L\beta)' (L \text{Var}(\hat{\beta}) L')^{-1} (L\hat{\beta} - L\beta)$$

Testing :  $H_0: \beta_B = \beta_0$   
vs  
 $H_1: \beta_B \neq \beta_0$

### Example

$H_0: p_i = 0$  vs  $H_1: p_i \neq 0$

If  $L = [0 - 1 - 0]$   
 $q = 1$

then  
 $F = \frac{(\hat{\beta}_i - \beta_0)^2}{\text{Var}(\hat{\beta}_i)}$

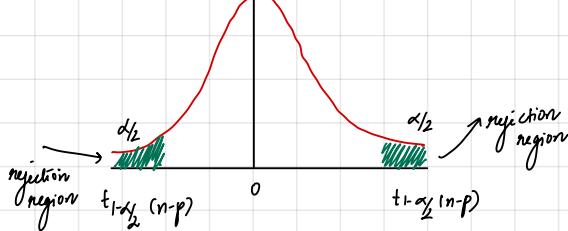
$H_0: p_i = \beta_0$  vs  $H_1: p_i \neq \beta_0$

$F$  should be rejected if  $\frac{(\hat{\beta}_i - \beta_0)^2}{\text{Var}(\hat{\beta}_i)} > F_{1-\alpha}(q, n-p)$

which is equivalent to

$$\frac{|\hat{\beta}_i - \beta_0|}{\sqrt{\text{Var}(\hat{\beta}_i)}} > \sqrt{F_{1-\alpha}(1, n-p)} \quad \text{since } q=1$$

$\Rightarrow \frac{|\hat{\beta}_i - \beta_0|}{\sqrt{\text{Var}(\hat{\beta}_i)}} \sim t(n-p) \quad (\text{we know})$



$$\frac{|\hat{\beta}_i - \beta_0|}{\sqrt{\text{Var}(\hat{\beta}_i)}} \sim t(n-p)$$

$$\Rightarrow [\hat{\beta}_i \pm t_{1-\alpha/2}(n-p) \sqrt{\text{Var}(\hat{\beta}_i)}] = C_{1-\alpha}(\hat{\beta}_i)$$

### Model Selection

If there are 2 models for the same data  
 $M_0, M_1$

Example:

$$M_0: Y = \beta_0 + \beta_1 X_1 + \epsilon$$

$$M_1: Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon$$

Here we need to identify if we can eliminate the coefficients  $\beta_2$  &  $\beta_3$  such that  $M_1$  becomes  $M_0$

Null hypothesis  $H_0: \beta_2 = \beta_3 = 0$

Alternative hypo  $H_1: \beta_2 = \beta_3 \neq 0$

- If the models are nested, we always use the F-test
- If the models are not nested?

Let  $L(\beta, \sigma^2 | \text{data}) \rightarrow \text{likelihood}$

For example: If  $\epsilon_1, \dots, \epsilon_n$  (iid)  $\sim N(0, \sigma^2)$   
 then

$$L(\beta, \sigma^2 | \text{data}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left( \frac{\epsilon_i - \beta_0 - \beta_1 x_i}{\sigma^2} \right)^2}$$

then  $\log(\text{likelihood})$  will be

$$\log L(\hat{\beta}, \hat{\sigma}^2 | \text{data})$$

Assume that

$\hat{\beta}_0, \hat{\sigma}_0^2$  are estimates of  $\beta_0, \sigma^2$  under model  $M_1$   
 $\hat{\beta}_0, \hat{\sigma}_0^2$  are estimates of  $\beta_0, \sigma^2$  for model  $M_0$

We will move from Model  $M_0$  to  $M_1$

if  $L(\hat{\beta}_0, \hat{\sigma}_0^2) > L(\hat{\beta}_1, \hat{\sigma}_1^2)$

i.e. if  $L(\hat{\beta}_0, \hat{\sigma}_0^2) - L(\hat{\beta}_1, \hat{\sigma}_1^2)$  is large

i.e. if  $\frac{L(\hat{\beta}_0, \hat{\sigma}_0^2)}{L(\hat{\beta}_1, \hat{\sigma}_1^2)}$  is large

### Akaike Information criterion (AIC)

It is an index used for model selection  
 $AIC(\beta, \sigma^2) = 2p - 2\log(L(\beta, \sigma^2))$

We should choose the model such that

we choose the model  $M_1$  if

$$AIC(\hat{\beta}_0, \hat{\sigma}_0^2) < AIC(\hat{\beta}_1, \hat{\sigma}_1^2)$$

### Bayesian Information criterion (BIC)

$$BIC(\beta, \sigma^2) = (\log n)p - 2\log(L(\beta, \sigma^2))$$

### Diagnostics

- Residual Analysis

- Influential cases (follows the same concept as Cook's distance)

$$D_i = \frac{(\hat{\beta} - \hat{\beta}_{(-i)})' (\text{Var}(\hat{\beta}))^{-1} (\hat{\beta} - \hat{\beta}_{(-i)})}{p} \quad i = 1, \dots, n$$

\$ after removing

1 item from dataset

Here  $D_i \in$  Cook's distance

Influential cases are those for which  $D_i$  is big

$D_i$  can be re-written as

$$D_i = \frac{1}{\sigma_{\text{REML}}^2} \left( \frac{\hat{\epsilon}_i}{\sqrt{1-h_{ii}}} \right)^2 \frac{h_{ii}}{1-h_{ii}}$$

$$H = X(X'X)^{-1}X'$$

Here  $h_{ii}$  is the  $i^{\text{th}}$  element on the diagonal of  $H$

In this  $\text{var}(\hat{\epsilon}_i) \propto 1/h_{ii}$

## Likelihood displacement

$$LD_i = 2 \left[ L(\hat{\beta}, \hat{\sigma}^2) - L(\hat{\beta}(-i), \hat{\sigma}^2(-i)) \right]$$

difference in likelihood before  
and after removing an influential  
element

$i$  is influential if  $LD_i$  is very large

## • Example : Recall

$$Y_i = X_i' \beta + \epsilon_i$$

$$\text{Var}(Y_i) = \text{Var}(\epsilon_i) = \sigma^2 \lambda_i$$

Suppose that  $y_i$  is a mean

$$y_i = \frac{1}{n_i} \sum_{j=1}^{n_i} z_{ij} \text{ if } z_{ij} \text{ iid with } \text{Var}(z_{ij}) = \sigma^2$$

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

$$x_i = \frac{1}{n_i}$$

## • Example :

$$y_i = \frac{1}{n_i} \sum_{j=1}^{n_i} z_{ij}$$

$$z_{ij} \text{ iid with } \text{Var}(z_{ij}) = \sigma^2$$

in this case

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

$$\therefore \lambda_i = n_i$$

In group 1

$$\lambda_i = \lambda_i(y_i) = \text{Var Fixed}(y_i)$$

considering  $\Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$

$$Y = X\beta + E$$

$$E \sim N(0, \sigma^2 I)$$

$$\Lambda = \Lambda \Lambda'$$

Consider  $\Lambda^{-1} = \begin{bmatrix} 1/\lambda_1 & & \\ & \ddots & \\ & & 1/\lambda_n \end{bmatrix}$

$$\begin{aligned} \Lambda^{-1} Y &= \Lambda^{-1} X\beta + \Lambda^{-1} E \\ \tilde{Y} &= \tilde{X}\tilde{\beta} + \tilde{E} \end{aligned}$$

Now

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

$$\text{where } \tilde{E} \sim N_n(0, \sigma^2 \Lambda^{-1} R \Lambda^{-1})$$

$\Lambda^{-1} \Lambda^{-1}$  can be written as

$$\Lambda^{-1} \Lambda \Lambda^{-1} \Lambda = I$$

$$\therefore \tilde{E} \sim N_n(0, \sigma^2 I)$$

using OLS or ML

$$\begin{aligned} \hat{\beta} &= (\tilde{X}' \tilde{X})^{-1} \tilde{X}' \tilde{Y} = (X' \Lambda^{-1} \Lambda' X)^{-1} X' \Lambda' \tilde{Y} \\ &= (X' R^{-1} X)^{-1} X' R^{-1} Y \end{aligned}$$

covariance matrix enters the picture

## Linear Models with heterogeneous variances

Model for true phenomenon

$$y = X\beta + E \sim N(0, \sigma^2 I)$$

$$E[y|x]$$

Model for observations

$$y_i = x_i' \beta + \epsilon_i \quad i=1 \dots n$$

$$\epsilon_i \sim N(0, \sigma_i^2)$$

parameters:  $\beta, \sigma_1^2 \dots \sigma_n^2$   
 $n \times p$  parameters

## Rephrasing the model into vector form

- $Y \in \mathbb{R}^n$
- $X$   $n \times p$  design matrix
- $\beta \in \mathbb{R}^p$
- $E \in \mathbb{R}^n \quad E \sim N(0, \sigma^2 I)$
- $R = \Lambda \Lambda'$  where  $\Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$   
Lambda (diagonal matrix)  
ie  $\sigma_i^2 = \sigma^2 \lambda_i^2$   
scaling parameter
- $\lambda_i = \lambda_i(\delta, m_i, v_i) \rightarrow$  variance function
- $\delta \in \mathbb{R}^k$  vector of parameters common to all cases
- $m_i = E[Y_i | X_i] \quad i=1 \dots n$   
 $= x_i' \beta$
- $v_i \in \mathbb{R}^k$  vector of known covariates

## Classification of variance functions

Variance functions can be classified into 4 groups  
(pinheiro & Bates 2000 book)

### Group 1 : Known variance function

$\lambda_i = \lambda_i(v_i)$  Variance function depends only on  
known covariates

This document is available free of charge on

Scaricato da Vishwesh Mente (vishweshmente@gmail.com) -  $\tilde{\beta} = \Lambda^{-1}(Y - X\beta)$

$$\sigma^2 = \frac{1}{n-p} \sum_{i=1}^n e_i^2 = \frac{1}{n-p} (\mathbf{y} - \hat{\mathbf{x}}\beta)' \Lambda^{-1} \Lambda^{-1} (\mathbf{y} - \hat{\mathbf{x}}\beta)$$

$$= \frac{1}{n-p} (\mathbf{y} - \hat{\mathbf{x}}\beta)' \Lambda^{-1} (\mathbf{y} - \hat{\mathbf{x}}\beta)$$

### \* Weighted Least Squares (WLS)

$$\begin{aligned}\hat{\beta} &= \min_{\beta} \sum_{i=1}^n (\mathbf{y}_i - \hat{\mathbf{x}}_i\beta)^2 \\ &= \min_{\beta} (\mathbf{y} - \hat{\mathbf{x}}\beta)' (\mathbf{y} - \hat{\mathbf{x}}\beta) \\ &= \min_{\beta} (\mathbf{y} - \hat{\mathbf{x}}\beta)' \Lambda^{-1} \Lambda^{-1} (\mathbf{y} - \hat{\mathbf{x}}\beta) \\ \hat{\beta} &= \min_{\beta} \sum_{i=1}^n \frac{(\mathbf{y}_i - \hat{\mathbf{x}}_i\beta)^2}{\lambda_i^2}\end{aligned}$$

This is called the weighted least squares (WLS)

$$\hat{\beta} = (\mathbf{X}' \mathbf{R}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{R}^{-1} \mathbf{y}$$

A

$$\therefore \hat{\beta} = A \mathbf{y}$$

$$\begin{aligned} \text{so } \text{Var}(\hat{\beta}) &= A \cdot \text{Var}(\mathbf{y}) A' \\ &\quad \text{||}_R \\ &= \sigma^2 (\mathbf{X}' \mathbf{R}^{-1} \mathbf{X})^{-1}\end{aligned}$$

$$\hat{\beta} \sim N_p(\hat{\beta}, \sigma^2 (\mathbf{X}' \mathbf{R}^{-1} \mathbf{X})^{-1})$$

### Inference:

$$H_0: L\beta = c_0 \text{ vs } H_1: L\beta \neq c_0$$

$$L \rightarrow q \times p$$

then

Test statistic

$$F = \frac{(L\hat{\beta} - c_0)' (L \text{Var}(\hat{\beta}) L')^{-1} (L\hat{\beta} - c_0)}{q}$$

so if  $H_0$  is true  $\Rightarrow F \sim F(q, n-p)$

$$\text{8 confidence interval for } \beta_i = [\hat{\beta}_i \pm t_{1-\alpha/2}(n-p) \sqrt{\text{Var}(\hat{\beta}_i)}]$$

$$\hat{\beta} = (\mathbf{X}' \mathbf{R}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{R}^{-1} \mathbf{y}$$

### \* Group 2: $\langle S \rangle$ group

$$\textcircled{1} \quad \lambda_i = \lambda_i(\underline{\epsilon}, v_i)$$

$\delta$   $\in$  unknown parameters

$v_i \in$  known covariates

Here the units are clustered in  $S \geq 1$  groups

Example:

$$\textcircled{i} \quad \lambda_i(\underline{\epsilon}, v_i) = |v_i| \quad \text{is unknown}$$

$$= \text{var power}(\underline{\epsilon}; v_i, s_i)$$

$$\begin{aligned} \textcircled{2} \quad \lambda_i(\underline{\epsilon}; v_i) &= \exp(v_i \underline{\epsilon}_{s_i}) \quad \text{unknown} \\ &= \text{var exp}(\underline{\epsilon}; v_i, s_i)\end{aligned}$$

$$\begin{aligned} \textcircled{3} \quad \lambda^*(\underline{\epsilon}; v_i) &= \delta_1 s_i + |v_i| \underline{\epsilon}_{s_i} \\ &= \text{var const. power}(\underline{\epsilon}; v_i, s_i)\end{aligned}$$

$$\begin{aligned} \textcircled{4} \quad \lambda_i(\underline{\epsilon}; v_i) &= \delta s_i \\ \text{unknown: } \delta s_1, \delta s_2, \dots, \delta s_n \\ \text{usually } \delta s_i = 1 \text{ for identifiability}\end{aligned}$$

### \* Estimation of the parameters

Maximum likelihood

$$L(\beta, \underline{\epsilon}, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2\lambda_i^2}} \exp\left[-\frac{1}{2\sigma^2} \frac{(\mathbf{y}_i - \hat{\mathbf{x}}_i\beta)^2}{\lambda_i^2}\right]$$

$$\ell(\beta, \underline{\epsilon}, \sigma^2) = \log L(\beta, \underline{\epsilon}, \sigma^2)$$

$$(\hat{\beta}_{ML}, \hat{\delta}_{ML}, \hat{\sigma}_{ML}^2) = \max_{(\beta, \underline{\epsilon}, \sigma^2)} \ell(\beta, \underline{\epsilon}, \sigma^2)$$

### \* Optimizing profile likelihood

For any fixed  $\delta \Rightarrow \hat{\beta} = \hat{\beta}(\delta)$  obtained by using the diamond formula

$$\lambda^*(\underline{\epsilon}, \sigma^2) = \lambda(\hat{\beta}(\delta), \underline{\epsilon}, \sigma^2)$$

objective would be to maximize

$$\max_{\delta, \sigma^2} \ell^*(\underline{\epsilon}, \sigma^2) = (\hat{\delta}, \hat{\sigma}^2)$$

There are chances that can get a biased estimator so instead of writing the likelihood for observations, we should find the likelihood for the residuals ( $\epsilon^*$ )

This gives us true:

$$\Rightarrow \hat{\delta}_{REML}, \hat{\sigma}_{REML}^2$$

$$\Rightarrow \hat{\beta} = \hat{\beta}(\hat{\delta}_{REML})$$

### Inference

$$\text{Var}(\hat{\beta}) = \sigma^2 (\mathbf{X}' \mathbf{R}^{-1} \mathbf{X})^{-1}$$

$$\Lambda \Lambda^{-1} \quad \Lambda = \begin{bmatrix} \lambda(s_1, v_1) \\ \vdots \\ \lambda(s_n, v_n) \end{bmatrix}$$

conduct testing using F statistic estimating  $\text{Var}(\hat{\beta})$  with  $\hat{\sigma}^2 (\mathbf{X}' \hat{\mathbf{R}}^{-1} \mathbf{X})^{-1}$

$$H_0: L\beta = c_0 \text{ vs } H_1: L\beta \neq c_0$$

distribution of test statistic will be

$$F \sim F(q, n-p)$$

but if  $n$  is large  $\chi^2$  distribution

$$F \sim \frac{1}{q} \chi^2(q)$$

## Special Important case:

$$\lambda_i = \lambda_i(S_i, v_i) = \delta_{S_i}$$

Consider the LR test based :

$$2 \left[ \ell_i(\hat{\beta}_i^k, \hat{v}_i^k, \hat{\sigma}^2) - \ell_i(\hat{\beta}_0^k, \hat{\sigma}^2_0) \right] > \chi^2_{\text{df}}(c\alpha)$$

If this is satisfied then reject the  $H_0$

## Group 3: $\langle \mu \rangle$

$$\lambda_i = \lambda_i(\mu_i, v_i)$$

$$\mu_i = \underline{x}_i' \hat{\beta}$$

Notable example:

$$\lambda_i = \mu_i$$

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

$$\text{cv}(y_i) = \frac{\text{std}(y_i)}{\mu_i} = \sigma$$

Examples :

$$\textcircled{1} \quad \lambda_i = |\mu_i| \begin{cases} \delta_{S_i} \\ \text{these are known} \end{cases} \quad S_i \in \{1, \dots, S\} \quad \begin{matrix} \text{Group membership} \\ \text{these is known} \end{matrix}$$

$$= \text{var power}(\mu_i, v_i, S_i)$$

$$\textcircled{2} \quad \lambda_i = \exp(\mu_i \delta_{S_i})$$

$$\textcircled{3} \quad \lambda_i = (\delta_{S_i} + |\mu_i| \delta_{S_i}) \begin{cases} \delta_{S_i} \\ \text{these are known} \end{cases}$$

## Group 4: $\langle \mu, \varepsilon \rangle$

$$\lambda_i = \lambda_i(\varepsilon, \mu_i, v_i)$$

example would be same as group 3 ( $\langle \mu \rangle$ )

but with unknown  $\varepsilon$

Note: In the  $\langle \mu \rangle$  group and  $\langle \mu, \varepsilon \rangle$  group  
 $\hat{\beta}$  (unknown) enters the mean & the variance

1. Try to maximize  $L(\beta, \varepsilon, \sigma^2)$

2. Iterative procedure for estimating  $\beta, \varepsilon, \sigma^2$

Step 0: Initial guess for  $\hat{\beta}$ :  $\hat{\beta}_0$   
(for instance use OLS)

For  $K \geq 1$  : repeat "until convergence"

Step 1: Get  $(\hat{\delta}^{(K)}, \hat{\sigma}^2(K)) = \max_{\delta, \sigma^2} \ell(\hat{\beta}^{(K-1)}, \delta, \sigma^2)$

i.e. assume that  $\mu_i = \underline{x}_i' \hat{\beta}^{(K-1)}$

Step 2: Estimate  $\hat{\beta}^{(K)}$  using  $(\Delta)$  and  $\hat{\sigma}$  by plugging  $\hat{\delta}^{(K)}, \hat{\sigma}^2(K)$

Stop when  $\hat{\beta}^{(K)}$  and  $\hat{\beta}^{(K-1)}$  are not very different w.r.t.  
to some metric.

for instance:

$$\|\hat{\beta}^{(K)} - \hat{\beta}^{(K-1)}\| \leq \epsilon \quad (\text{Euclidean distance})$$

$$\sqrt{(\hat{\beta}^{(K)} - \hat{\beta}^{(K-1)})' \text{var}(\hat{\beta}^{(K-1)}) (\hat{\beta}^{(K)} - \hat{\beta}^{(K-1)})} \leq \epsilon$$

(Mahalanobis distance)

After this the question becomes what are the distributions of  $\hat{\beta}$  and  $F$ .

## \* Linear Models with fixed effects & correlated errors

$$\begin{aligned} R \rightarrow Y &= X' \beta + \varepsilon \\ R^P \rightarrow X, \quad \beta \in R^P, \quad \varepsilon \sim N(0, \sigma^2 I) \end{aligned} \quad \left. \begin{matrix} \text{Model for phenomenon} \\ \text{units } i \in \{1, 2, \dots, N\} \\ \text{each unit } i \text{ is observed } n_i \text{ times} \\ (n = n_1 + n_2 + \dots + n_N) \end{matrix} \right.$$

For a Unit  $i$ :

$$R^{n_i} \rightarrow y_i = x_i' \beta + \varepsilon_i$$

$$x_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{in_i} \end{bmatrix} \quad x_{ij} \in R^P \quad (n_i \times P \text{ matrix})$$

$$\varepsilon_i \sim N_{n_i}(0, \sigma^2 R_i)$$

$R_i$  need not be diagonal

Moreover: different units, say  $i$  and  $i' \in \{1, \dots, N\}$   
are independent

$R \rightarrow y_{ij}$  observation for unit  $i \in \{1, \dots, N\}$   
at instance  $j \in \{1, \dots, n_i\}$

- $E[Y_{ij}|X_{ij}] = X_{ij}' \beta$
- $\text{var}[Y_{ij}] = \text{var}[\varepsilon_{ij}] = \sigma^2 R_i(j, j)$
- $\text{cov}[Y_{ij}, Y_{ij'}] = \sigma^2 R_i(j, j')$  such that  $j, j' \in \{1, \dots, n_i\}$   
 $\Rightarrow \text{cov}(\varepsilon_{ij}, \varepsilon_{ij'}) = 0$

Recall:

$$\text{corr}(z, w) = \frac{\text{cov}(z, w)}{\sqrt{\text{var}(z) \cdot \text{var}(w)}}$$

$$\text{cov}(z, w) = \sqrt{\text{var}(z)} \text{corr}(z, w) \sqrt{\text{var}(w)}$$

$$\begin{aligned} \therefore \sigma^2 R_i(j, j') &= \text{cov}(Y_{ij}, Y_{ij'}) = \text{cov}(\varepsilon_{ij}, \varepsilon_{ij'}) \\ &= \sqrt{\text{var}(Y_{ij})} \text{corr}(Y_{ij}, Y_{ij'}) \sqrt{\text{var}(Y_{ij})} \\ &= \lambda_{ij} G(j, j') \lambda_{ij'} \end{aligned}$$

For  $\lambda_{ij} \rightarrow$  variance function already introduced  $\langle \delta \rangle, \langle u \rangle, \langle \delta, u \rangle$

- $C_{ij} \rightarrow$  correlation structures

Let

$$\Lambda_i = \begin{bmatrix} \lambda_{ii} & 0 \\ 0 & \lambda_{in_i} \end{bmatrix} = \text{diag}[\lambda_{ii}, \dots, \lambda_{in_i}] \quad (\text{mixn}_i)$$

$$C_i = \begin{bmatrix} 1 & C_i(j_i, j'_i) \\ C_i(j'_i, j_i) & 1 \end{bmatrix} \Rightarrow C_i(j_i, j'_i) \text{ because of symmetry}$$

then

$$R^{ni} \rightarrow \varepsilon_i \sim N_{nr} (0, \sigma^2 R_i)$$

$$R_i = \Lambda_i C_i \Lambda_i'$$

&

$$\text{Var}(Y_{ij}) = \sigma^2 \tilde{\lambda}_{ij}^2$$

### \* Correlation structures

Most of the times we consider structure of this form:

$$\text{corr}(\varepsilon_{ij}, \varepsilon_{ij'}) = C_i(j_i, j'_i) = h(a(t_{ij}, t_{ij'}), f)$$

where

$a$  is a metric

$t_{ij}$  is a vector in (time, space)

identifying the instance where  $Y_{ij}$  has been observed

Assumptions on  $h$ :

- continuous (regularity)
- $h(0, f) = 1$

Examples

- $\text{corr}(\varepsilon_{ij}, \varepsilon_{ij'}) = \rho_i, | \rho_i | \leq 1$   
it is corr-comp-symm in R-studio
- $\text{corr}(\varepsilon_{ij}, \varepsilon_{ij'}) = \rho_i \delta_{jj'}$   
 $= \text{corr AR}(i)$

These are very commonly used and are called serial autocorrelation functions

- $\text{corr}(\varepsilon_{ij}, \varepsilon_{ij'}) = e^{-\frac{d(t_{ij}, t_{ij'})}{\beta_i}}, \beta_i > 0$   
 $t_{ij}$ : location in space  
where  $Y_{ij}$  has been observed

- $\text{corr}(\varepsilon_{ij}, \varepsilon_{ij'}) = e^{-\frac{d^2(t_{ij}, t_{ij'})}{\beta_i^2}}$   
 $= \text{corr-Gauss}$
- $\text{corr}(\varepsilon_{ij}, \varepsilon_{ij'}) = \left[ 1 - \frac{a(t_{ij}, t_{ij'})}{\beta_i} \right]$  (this is an indicator)  
 $= \text{corr-Lin}$

$$\text{H}[\delta \rho_i] = \begin{cases} 0 & \text{if } d > \rho_i \\ 1 & \text{if } d \leq \rho_i \end{cases}$$

Indicated symbol

$$\cdot \text{corr}(\varepsilon_{ij}, \varepsilon_{ij'}) = \left[ 1 - \frac{3}{2} \frac{d(t_{ij}, t_{ij'})}{\beta_i} + \frac{1}{2} \frac{d^3(t_{ij}, t_{ij'})}{\beta_i^3} \right] + [\delta \rho_i]$$

$= \text{corr-sphr.}$

### \* Excuses on stationarity and ARMA models

$\{x_t\}_{t \geq T}$  stochastic process

$\{x_t\}_{t \geq T}$  is said to be strongly stationary

if:  $\forall n$  and  $(t_1, \dots, t_n)$  and  $T > 0$

$$F(x_{t_1}, \dots, x_{t_n}) = F(x_{t_1+T}, \dots, x_{t_n+T})$$

$$\Rightarrow F(x_{t_1}) = F(x_{t_2})$$

distributions of  $x_{t_1}$  &  $x_{t_2}$  are the same

If we have the stationary then

$$\text{① } E[x_t] = \mu$$

$$\text{② } \text{cov}(x_t, x_{t+h}) = h(1, t-h)$$

If we have the above 2 conditions satisfied

then  $\{x_t\}$  is said to be weakly stationary

Note:

$$\text{Var}(x_t) = h(0)$$

Consider this process  $\{x_t\}$  weakly stationary

$$\text{① } E[x_t] = 0 \quad t = -3, -1, 0, 1, 2, \dots$$

$$\text{② } x_t = \phi x_{t-1} + \varepsilon_t \quad \varepsilon_t \text{ are independent of } x_t; \quad \varepsilon_t \sim N(0, \sigma^2)$$

$$\text{Var}(x_T) = \text{Var}(\phi x_{T-1} + \varepsilon_T)$$

$$= \phi^2 \text{Var}(x_{T-1}) + \sigma^2$$

$$\text{Var}(x_t) = \phi^2 \text{Var}(x_{t-1}) + \sigma^2$$

$$\therefore \text{Var}(x_t) = \frac{\sigma^2}{1-\phi^2} \quad |\phi| < 1$$

$$\text{corr}(x_{t+1}, x_t) = \frac{\phi \sigma^2}{1-\phi^2}$$

$$\text{cov}(x_{t+1}, x_t) = \text{cov}(\phi x_t + \varepsilon_t, x_t) \\ = \phi \frac{\sigma^2}{1-\phi^2} + 0$$

$$\text{Similarly } \text{corr}(x_{t+2}, x_t) = \phi^2$$

$$\text{corr}(x_{t+k}, x_t) = \phi^k \quad |\phi| < 1$$

In an ARMA model

$$x_t = \sum_{i=1}^p \phi_i x_{t-i} + \sum_{i=1}^q \theta_i \varepsilon_{t-i} + \varepsilon_t$$

$\underbrace{\quad}_{\text{AR}(p)}$        $\underbrace{\quad}_{\text{MA}(q)}$

### \* Estimating $\beta$

$$y = X\beta + \varepsilon$$

If  $R$  is known:

$$R^{-\frac{1}{2}} y = R^{-\frac{1}{2}} X \beta + R^{-\frac{1}{2}} \varepsilon$$

such that  $R^{-\frac{1}{2}} \varepsilon \sim N_n (0, \sigma^2 I)$

$$\hat{\beta} = (X' A^{-1} X)^{-1} X' R^{-\frac{1}{2}} y$$

Take variance functions in  $\langle \delta \rangle$  group  
 $\mu = \hat{x}'\beta$

$$\hat{\beta} = \hat{f}_1(\delta, \beta)$$

To estimate  $\delta$  and  $\beta$ , estimate the profile likelihood

$$l^* \text{lik}(\sigma^2, \delta, \beta) = l \text{lik}(\hat{\beta}(\delta, \beta), \sigma^2, \delta, \beta)$$

$$(\hat{\sigma}^2, \hat{\delta}, \hat{\beta}) = \max_{\sigma^2, \delta, \beta} l^* \text{lik}(\hat{\beta}(\delta, \beta), \sigma^2, \delta, \beta)$$

$$\Rightarrow \hat{\beta} = \hat{f}(\hat{\delta}, \beta) = (X' R^{-1} X)^{-1} X' R^{-1} Y$$

- For unbiasedness of  $\hat{\sigma}^2$ , consider REML:

$$\begin{aligned} l^* \text{lik}(\sigma^2, \delta, \beta) & \text{ are } l \text{lik} \text{ for} \\ \hat{\beta} &= Y - X \hat{\beta} \\ \Rightarrow \hat{\sigma}^2 &= \frac{1}{n-p} (Y - X \hat{\beta})' R^{-1} (Y - X \hat{\beta}) \end{aligned}$$

Profile again

$$\begin{aligned} l^* \text{lik}(\delta, \beta) &= l \text{lik}(\hat{\beta}(\delta, \beta), \sigma^2, \delta, \beta) \\ (\hat{\delta}, \hat{\beta}) &= \max \text{lik}(\delta, \beta) \\ \therefore \hat{\beta} &= \hat{f}(\hat{\delta}, \beta) \\ \hat{\sigma}^2 &= \hat{\sigma}^2(\hat{\delta}, \beta) \end{aligned}$$

- For  $\langle u \rangle$  and  $\langle u, \delta \rangle$  group

$\rightarrow$  use iterative procedure introduced in previous chapters

The problem with these procedures is that the  $R$  matrix is a covariance

$$\hat{R} = \begin{bmatrix} R_1 & & \\ & \ddots & \\ & & R_N \end{bmatrix} \text{ and we should make sure that } \hat{R} \text{ is a covariance matrix}$$

Covariance  $\equiv$  of a vector  $w$ :

For a vector to be covariant

- $Z$   $n \times n$  should be symmetric

$w \in R^n$

Linear combinations of components of  $w$   
 $\text{Var}(a_1 w_1 + a_2 w_2 + \dots + a_n w_n) \quad \text{if } a \in R$   
 $= a' Z a > 0$

## Linear Mixed Models

Model for the phenomenon

$$R \rightarrow Y = X'\beta + Z'b + \varepsilon$$

random vector of coefficients

$X \in R^p$  vector of known regressors

$Z \in R^{q \times p}$  " " " "

$\beta \in R^p$  fixed but unknown

$b \sim N_q(0, \sigma^2 D)$  here  $D$  is a  $q \times q$  matrix

$\varepsilon \sim N_p(0, \sigma^2 I)$

Model for the data

Suppose that there is 1 layer of clustering  
with clusters  $\{1, 2, \dots, N\}$

So in cluster  $i$ , there contains  $n_i$  units

$$R^{n_i} \rightarrow Y_i = X_i \beta + Z_i b_i + \varepsilon_i$$

Random effects

where  $X_i : n_i \times p$  design matrix  
 $Z_i : n_i \times q$  design matrix  
 $b_i \sim N_q(0, \sigma^2 D)$   
 $\beta \in R^p$  fixed coeff (unknown)  
 $\varepsilon_i \sim N_{n_i}(0, \sigma^2 R_i)$

## Layers for hierarchy

1 layer with  $N$  groups  $\{1, 2, \dots, N\}$

2 layers: in group  $i \in \{1, 2, \dots, N\}$

$\rightarrow n_i$  groups

$\Rightarrow$  for  $i \in \{1, \dots, N\}$   $j \in \{1, \dots, n_i\}$

there are  $n_{ij}$  observations

$$R^{n_{ij}} \rightarrow Y_{ij} = X_{ij} \beta + Z_{1,ij} b_i + Z_{2,ij} b_{ij} + \varepsilon_{ij}$$

$(n_{ij} \times p)$

here

$$\begin{aligned} X_{ij} &: n_{ij} \times p \text{ design matrix} \\ Z_{1,ij} &: n_{ij} \times q_1 \text{ design matrix} \\ Z_{2,ij} &: n_{ij} \times q_2 \text{ design matrix} \end{aligned}$$

These are all known

$$\begin{aligned} b_i &\sim N_q(0, \sigma^2 D_1) \\ b_{ij} &\sim N_{q_1}(0, \sigma^2 D_2) \\ \varepsilon_{ij} &\sim N_{n_{ij}}(0, \sigma^2 R_{ij}) \end{aligned}$$

These are all independent of each other

## \* Overall model with 1 layer

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^{n_1+n_2+\dots+n_N = n}$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} \text{ } n \times p$$

$$\mathbf{z} = \begin{bmatrix} z_1 & 0 \\ z_2 & \ddots \\ \vdots & z_N \end{bmatrix} \text{ } n \times N_q$$

$$\mathbf{b} = N_{N_q} [0, \sigma^2 \begin{bmatrix} D & 0 \\ 0 & D \end{bmatrix}]$$

$$\mathbf{v} \sim N_n (0, \sigma^2 \begin{bmatrix} R_1 & & \\ & \ddots & \\ & & R_N \end{bmatrix})$$

Model :

$$\mathbf{y} = \mathbf{x}\beta + \mathbf{z}\mathbf{b} + \mathbf{v}$$

$\mathbf{v} \perp \mathbf{b}$  (they are independent)

## \* conditional distribution of $\mathbf{y}$ given $\mathbf{b}$

for  $i \in \{1, 2, \dots, n\}$

$$y_i | b_i \sim N_{n_i} (x_i \beta + z_i b_i, \sigma^2 R_i)$$

$y_i | b_i$  and  $y_j | b_i$  are independent

## \* Marginal distribution of $\mathbf{y}$

$$y_i | b_i \sim N_{n_i} (x_i \beta + z_i b_i, \sigma^2 R_i)$$

$$b_i \sim N_q (0, \sigma^2 D)$$

$$\Rightarrow y_i \sim N_{n_i} (x_i \beta, \sigma^2 z_i' D z_i + \sigma^2 R_i)$$

$$z_i' D z_i + R_i = v_i$$

substitute here

$$\therefore y_i \sim N_{n_i} (x_i \beta + \sigma^2 v_i)$$

$$\text{where } v_i = \begin{bmatrix} v_1 \\ \vdots \\ v_N \end{bmatrix}$$

## \* Structuring $D$ & $R_i$

$$\text{when } q \text{ is small } D = \begin{bmatrix} d_{11} & d_{12} & d_{1q} \\ d_{21} & \ddots & \vdots \\ \vdots & \vdots & d_{qq} \end{bmatrix}$$

$$d_{ij} \text{ free : } \frac{q(q+1)}{2} \text{ free parameters}$$

$$\text{when } q \text{ is large : } D = \begin{bmatrix} q_{11} & 0 \\ 0 & q_{NN} \end{bmatrix}$$

$$R_i = \Lambda_i C_i \Lambda_i'$$

where

$$\Lambda_i = \begin{bmatrix} \lambda_{ii} & 0 \\ 0 & \lambda_{inc} \end{bmatrix}$$

where  $\lambda_{ini}$  = variance function of  $\langle \mathbf{v} \rangle$  group

$C_i$  :  $n_i \times n_i$  correlation of spatial, spatio-temporal structure

## \* Estimating $\beta$

If  $v$  was known (ie.  $D$  &  $R_i$ ,  $R_N$  is known)

$$\hat{\beta} = (X' V^{-1} X)^{-1} X' V^{-1} Y$$

use  $\hat{\beta}$  to profile log likelihood and then maximize w.r.t. the other parameters.

$$\hat{\sigma}^2 = \frac{1}{n-p} \sum_{i=1}^n (y_i - x_i \hat{\beta}) v_i^{-1} (y_i - x_i \hat{\beta})$$

## \* Posterior distribution of $b_i$ given $y_i$

$$y_i | b_i \sim N_{n_i} (x_i \beta + z_i b_i, \sigma^2 R_i)$$

$$b_i \sim N_q (0, \sigma^2 D)$$

$$y_i \sim N_{n_i} (x_i \beta, \sigma^2 z_i' D z_i + \sigma^2 R_i)$$

$\Rightarrow b_i | y_i$  is gaussian in  $R^q$

$$\begin{aligned} E[b_i | y_i] &= \text{mean}(b_i | y_i) \\ &= D z_i' v_i^{-1} (y_i - x_i \beta) \end{aligned}$$

Predictors of  $b_i$ :

$$\hat{b}_i = D z_i' v_i^{-1} (y_i - x_i \beta)$$

( $b_i$  would be random & different for different observations)

$\text{Var}(b_i)$  see formula (B.2) in G&B

Note:  $\forall \mathbf{c} \in R^q$

$$\text{Var}(\mathbf{c}' b_i) \leq \sigma^2 \mathbf{c}' D \mathbf{c} = \text{var}(\mathbf{c}' b_i)$$

## Spatial Statistics

### Table's first law of Geography

Everything is related with everything else, but near things are more related than distant things.

#### Process

$$\{z_s, SED\}$$

D: spatial domain

$z_s$ : observation at location SED

Spatial locations:  $s_1, \dots, s_n \in D$

Data:  $z_{s_1}, \dots, z_{s_n}$

n data objects: real numbers, labels, vectors, functions.

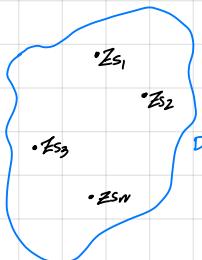
### Types of Spatial Data

#### Geostatistical Data

In this  $D \subseteq \mathbb{R}^d$   $d=1, 2, 3, \dots$   
(D is a subset of  $\mathbb{R}^d$ )

and D is a rectangle

$z_s$  is defined for every SED

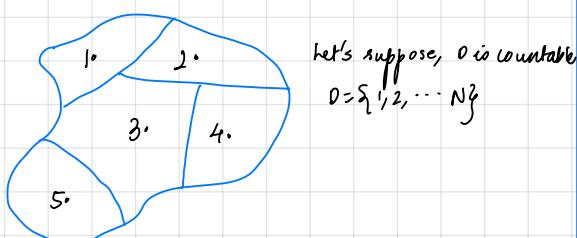
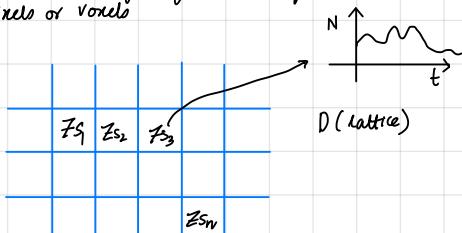


Goal: • To estimate the spatial dependence of variables  
(how variables are related due to their positions in D)

- Structured spatial covariances → estimate parameters
- Prediction at an unsampled location (so  $\notin D$ )

#### Lattice Data

D is a countable (often finite) set of sites (administrative areas)  
on pixels or voxels



Neighbouring structure:

Given  $i \in D$

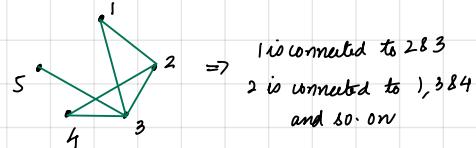
$$N_i := \{j \in D \text{ in the neighbourhood of } i\}$$

This document is available free of charge on

After the neighbouring structure is specified, D becomes a graph

Ex:  $i$  &  $j$  are the same neighbour if they share a border

Nodes of the graphs are the sites



Formally:

$W$  is a  $n \times n$  matrix  
where  $n$  is no. of sites  
elements of lattice

$$W = [w_{ij}]_{n \times n}$$

$$w_{ij} = \begin{cases} \neq 0 & \text{if } j \text{ is a neighbour of } i \\ 0 & \text{otherwise} \end{cases}$$

	1	2	3	4	5
1	0	1	1	0	0
2	1	0	1	1	0
3			0		
4				0	
5					0

### Standardizing the weight matrix

standardizing by rows

$$\tilde{w}_{ij} = \frac{w_{ij}}{\sum_{j=1}^n w_{ij}} \Rightarrow \sum_{j=1}^n \tilde{w}_{ij} = 1$$

$$\sum_{i=1}^n \sum_{j=1}^n \tilde{w}_{ij} = n$$

Suppose  $d: D \times D \rightarrow [0, \infty)$  is a metric

$d(i, j)$  is the distance b/w  $i, j$  in D

Example 1: Fix  $\delta > 0$

$$w_{ij} = \begin{cases} 1 & \text{if } d(i, j) < \delta \\ 0 & \text{otherwise} \end{cases}$$

Example 2: K-nearest neighbour

$$w_{ij} = \begin{cases} 1 & \text{if } j \text{ is among the } k \text{-nearest sites to } i \\ 0 & \text{otherwise} \end{cases}$$

Example 3:  $w_{ij} = f(d(i, j), \theta)$

$$\cdot w_{ij} = \frac{1}{d(i, j)^\alpha} \quad \alpha > 0$$

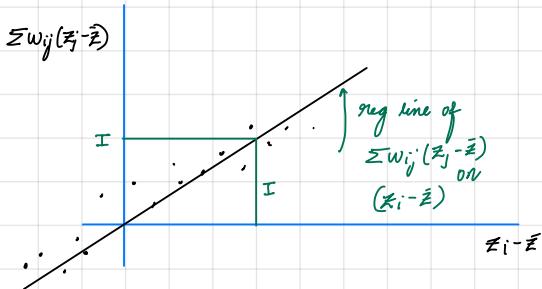
$$\cdot w_{ij} = \exp[-\beta d(i, j)], \beta > 0$$

## \* Spatial Auto-correlation

### Moran's Index (I)

$$(a) I = \frac{\sum \sum w_{ij} (z_i - \bar{z})(z_j - \bar{z})}{\sum (z_i - \bar{z})^2} \quad \text{if weights are standardized}$$

$$(b) I = \frac{\sum \sum w_{ij} (z_i - \bar{z})(z_j - \bar{z})}{\sum \sum w_{ij}} \cdot \frac{1}{\frac{\sum (z_i - \bar{z})^2}{n}} \quad \text{if weight are not standardized}$$



### \* checking for auto-correlation

$H_0$ : no spatial auto corr.

$H_1$ : spatial auto corr.

we use the permutation test

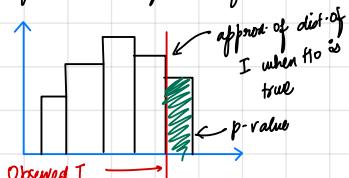
for  $k=1, \dots, M$  (large)

- permute  $(z_1, \dots, z_n)$  (sites/locations)

we get  $\pi^{(1)}, \dots, \pi^{(n)}$  where  $\pi$  is a permutation of  $\{1, \dots, n\}$

- compute  $I: I^K$

- compute the histogram of  $I^{(1)}, \dots, I^{(K)}$



LISA

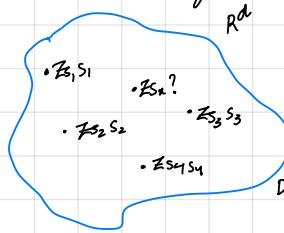
$$I_i := \frac{(z_i - \bar{z}) \sum w_{ij} (z_j - \bar{z})}{\sum (z_i - \bar{z})^2}$$

## \* Geostatistical Data

We observe a variable

$R \in \mathbb{R}^d$  where  $S \subseteq D \subseteq \mathbb{R}^d$ ,  $d=1, 2, 3, \dots$

$D$  should contain a rectangle in  $\mathbb{R}^d$



If data: fix  $s_1, \dots, s_n$   
 $z_{s_1}, \dots, z_{s_n}$

Used to predict the value of  $z_s$  in a specific location

### General Model

# SED

we assume that

$$z_s = m_s + \delta(s) \quad \begin{matrix} \text{random} \\ \text{part} \\ \text{deterministic} \end{matrix}$$

where

$m_s \rightarrow$  drift / trend of the process (captures large scale variation)

$\delta(s) \rightarrow$  residual (captures the local scale variation)

### \* Identifying process stationarity

$Z = \{z_s, SED\}$  is a random field

The Law of  $z$  is known if

$\forall n \geq 1$  and  $\forall s_1, \dots, s_n \in D$

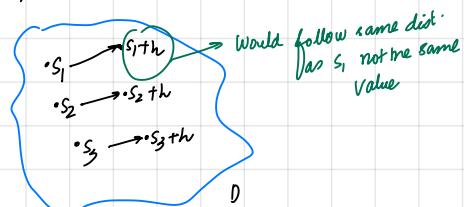
$$F_{s_1, \dots, s_n}(z_1, \dots, z_n) = P[z_{s_1} \leq z_1, \dots, z_{s_n} \leq z_n] \quad \forall z_1, \dots, z_n \in \mathbb{R}$$

These distributions should follow some consistency conditions

(Kolmogorov condition)

Definition: (Strong Stationarity)

$\{z_s, SED\}$  is said to be strongly stationary if for every  $n \geq 1$  and  $s_1, \dots, s_n$  the distribution of  $z_{s_1}, \dots, z_{s_n}$  is the same as the distribution of  $z_{s_1+h}, \dots, z_{s_n+h} \quad \forall h \in \mathbb{R}^d$



Definition: (weak stationarity)

$\{Z_s, SED\}$  is said to be weakly stationary if

- (1)  $E[Z_s] = m \quad \text{+ SED}$  Cov function of  $h$
  - (2)  $\text{cov}(Z_s, Z_{s+h}) = c(h) \quad \text{+ SED, + herd}$  This happens if  $E[Z_s^2] < \infty \quad \text{+ SED}$
- 

Proposition

Strong stationarity implies weak stationarity but not vice versa

The function

$c: R^d \rightarrow R$  such that  $c(h) = \text{cov}(Z_s, Z_{s+h})$

is called a covariogram.

Properties of  $c$ :

(1)  $\forall n \geq 1$ , and any  $s_1, \dots, s_n \in D$  define

$$c_{ij} = \text{cov}(Z_{s_i}, Z_{s_j}) = c(s_i - s_j)$$

$C = [c_{ij}]$   $n \times n$  matrix

$$\begin{aligned} \text{If } (\lambda_1, \dots, \lambda_n) \in R^n &= \lambda \\ &\leq \text{var}\left(\sum_{i=1}^n \lambda_i Z_{s_i}\right) = \lambda C \lambda \\ &= \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j c_{ij} \\ &= \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j c(s_i - s_j) \end{aligned}$$

$\therefore c$  is a positive definite function

(2)  $c(h) = c(-h)$

$$\Rightarrow \text{cov}(Z_s, Z_{s+h}) = \text{cov}(Z_{s-h}, Z_s)$$

(3)  $|\text{cov}(Z_s, Z_{s+h})| \quad \text{+ } h \in R^d$

$$= \frac{|\text{cov}(Z_s, Z_{s+h})|}{\sqrt{\text{var}(Z_s) \cdot \text{var}(Z_{s+h})}} \quad \text{var}(Z_s) = \text{cov}(Z_s, Z_s) = c(0) \text{ since } h=0$$

$$= \frac{|c(h)|}{\sqrt{c(0) \cdot c(0)}} = \frac{|c(h)|}{c(0)}$$

$$\therefore |c(h)| \leq c(0)$$

Definition: Intrinsic stationarity

$\{Z_s, SED\}$  is said to be intrinsically stationary if

(1)  $E[Z_s] = m \quad \text{+ SED}$

(2)  $\text{var}(Z_s - Z_{s+h}) = 2 \cdot \delta(h) \quad \text{+ SED, + herd}$   
where

2  $\gamma: R^d \rightarrow R$  is called variogram

$\gamma: R^d \rightarrow R$  is called semivariogram

Proposition: If a field is weakly stationary, it is also intrinsically stationary

Let  $\{Z_s, SED\}$  is weakly stationary

$$\begin{aligned} \text{var}(Z_s - Z_{s+h}) &= \text{var}(Z_s) + \text{var}(Z_{s+h}) - 2 \cdot \text{cov}(Z_s, Z_{s+h}) \\ &= c(0) + c(0) - 2 \cdot c(h) \\ &= 2c(0) - 2c(h) \\ &= 2 \cdot \delta(h) \end{aligned}$$

$\therefore$  for weakly stationary process

$$\delta(h) = c(0) - c(h) \quad \text{+ herd}$$

Properties of  $\gamma$

(1)  $\forall n \geq 1, s_1, \dots, s_n \in D$

and for  $\lambda_1, \dots, \lambda_n$  such that  $\sum \lambda_i = 0$

$$\sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j \gamma(s_i - s_j) \leq 0 \quad (\text{negative definiteness})$$

(2)  $\gamma(h) = \gamma(-h) \quad \text{+ herd}$

(3)  $\gamma(h) \geq 0 \quad \text{+ herd}$

(4)  $\gamma(0) = \text{var}(Z_s - Z_s) = 0$

(5)  $\lim_{\|h\| \rightarrow \infty} \frac{\gamma(h)}{\|h\|^2} = 0 \quad (\text{sub quadratic})$

$\|h\| \rightarrow \text{distance}$

Definition

If  $\gamma(h)$  is indeed a function of  $\|h\|$   
then  $\{Z_s, SED\}$  intrinsically stationary is  
also isotropic.

i.e.

$$\text{var}(Z_s - Z_{s+h}) = 2 \gamma(\|h\|) \quad \text{+ SED, herd}$$

WKT

①  $\gamma(0) = 0$  since  $\gamma(0) = \text{Var}(Z_S - Z_S) = \text{Var}(0)$

$\lim_{\|h\| \rightarrow 0} \gamma(h) = \tau^2 \geq 0$  if  $\tau^2 > 0 \Rightarrow \gamma(h)$  is not continuous in 0

$\tau^2$  is called nugget  
(measures instantaneous variability)

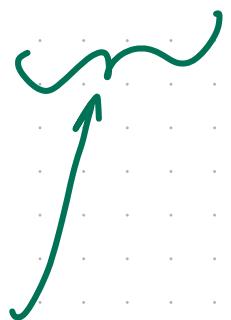
②  $\lim_{\|h\| \rightarrow \infty} \gamma(h) = \gamma^2 + \tau^2$  if it is finite

$$= \text{Var}(Z_S - Z_{S+h}) = 2 \left[ c(0) - c(ch) \right]$$

when  $h \rightarrow \infty, c(h) \rightarrow 0$

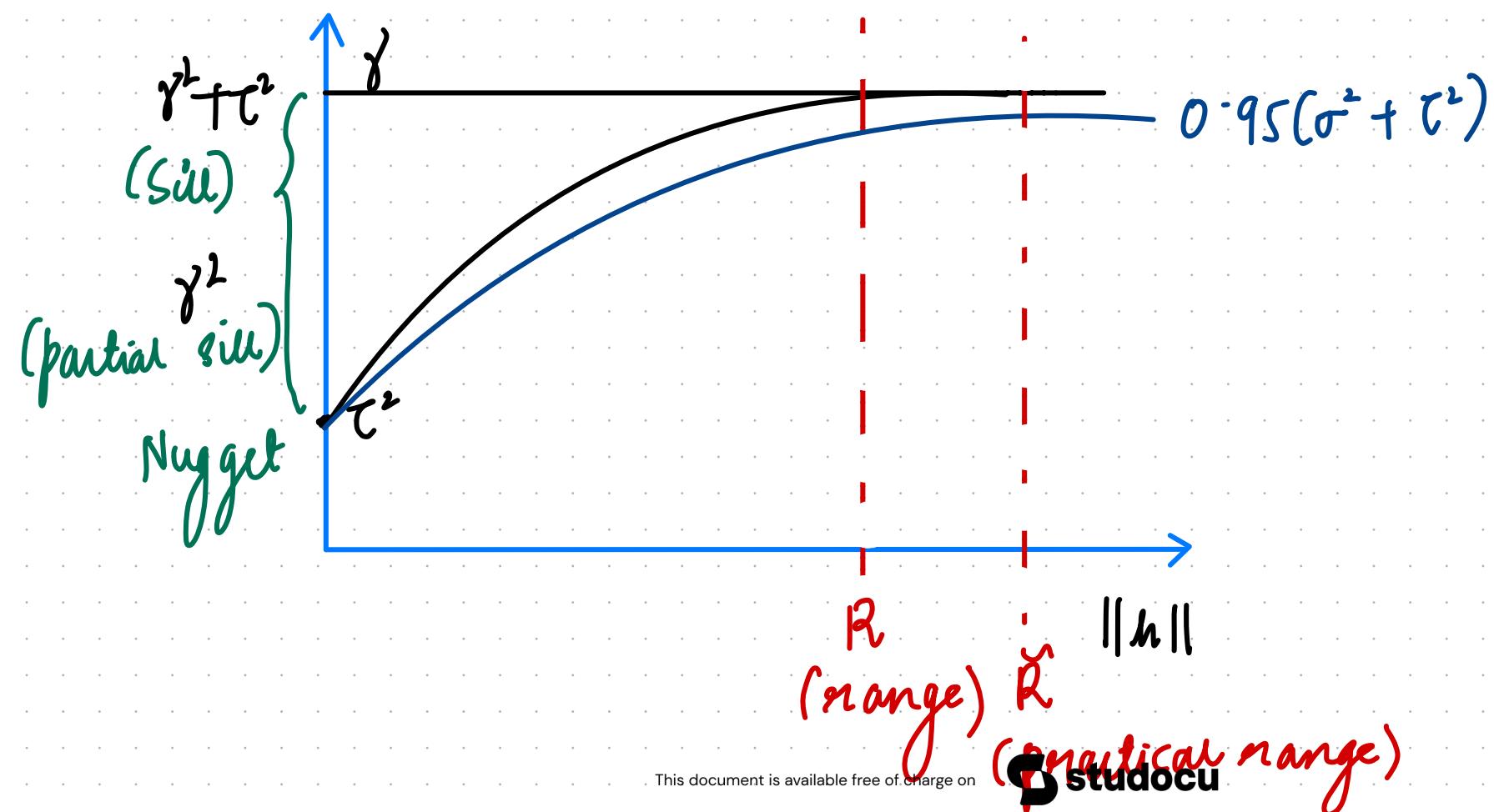
$$\therefore = 2 \times c(0)$$

If  $\gamma^2 + \tau^2 = \lim_{\|h\| \rightarrow \infty} \delta(h)$  is finite then



$\{z_s, s \in D\}$  is weakly stationary

is called Sill



# Parametric Models for semi-variogram

- Pure Nugget

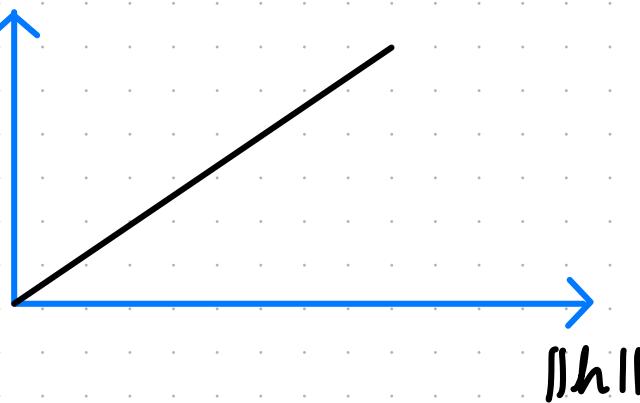
$$\gamma(h) = \begin{cases} 0 & h=0 \\ T^2 & h \neq 0 \end{cases}$$



- Linear Model (Isotropic)

$$\gamma(h) = a^2 \|h\|$$

(no sill, no nugget)



- Exponential Model

$$\gamma(h) = \sigma^2 \left( 1 - \exp \left[ -\frac{\|h\|}{a^2} \right] \right) \quad a, \sigma \neq 0$$

No nugget, sill =  $\sigma^2$ , practical range =  $3a^2$

## Gaussian Model

$$\gamma(h) = \sigma^2 \left[ 1 - \exp \left( - \frac{\|h\|^2}{a^2} \right) \right] \quad a, \sigma \neq 0$$

No nugget, sill =  $\sigma^2$ , practical range = ?

## Spherical Model

$$\gamma(h) = \begin{cases} \sigma^2 \left\{ \frac{3}{a} \frac{\|h\|}{a^2} - \frac{1}{2} \left( \frac{\|h\|}{a^2} \right)^3 \right\} & \text{for } 0 \leq \|h\| \leq a^2 \\ \sigma^2 & \text{for } \|h\| \geq a^2 \end{cases}$$

No nugget, sill :  $\sigma^2$ , range =  $a^2$

## Proposition

If  $\gamma_1$  &  $\gamma_2$  are valid semi-variograms then

- ①  $\gamma_1 + \gamma_2$  is a valid semi-variogram.
- ②  $b^2\gamma_1$  is a valid semi-variogram.

## Estimating the semi-variogram

$$\gamma(h) = c(0) - c(h)$$

$$\gamma(h) = E[(z_s - z_{s+h})^2] \quad h \in R^d$$

$$\hat{\gamma}(h) = \frac{1}{2|N(h)|} \sum_{(s_i, s_j) \in N(h)} [z_{s_i} - z_{s_j}]^2 \quad \left. \begin{array}{l} \text{Empirical estimate of } \gamma \\ \text{denoted by } \hat{\gamma}(h) \end{array} \right\}$$

$$N(h) = \{ (s_i, s_j) : \|s_i - s_j\| = \|h\| \}$$

There are only a finite no. of  $(s_i - s_j)$  ie. only a finite number of  $h$ 's for which  $\hat{\gamma}$  is computed

To deal with the variability, we adjust the formula such that

$$\hat{\gamma}(h_i) = \frac{1}{2|N(h)|} \sum_{(s_i - s_j) \in N(h)} [z_{s_i} - z_{s_j}]^2$$

$$N(h) = \{(s_i, s_j) : \|s_i - s_j\| \in T_i\}$$

where  $h_i$  is the center of  $T_i$   
 $i = 1, \dots, K$

$K$  is no. of bins

## 2-step procedure of estimating $\hat{\gamma}$

- ① Estimate empirical semi-varogram  $\hat{\gamma}$  by binning data  
 $\Rightarrow \hat{\gamma}(h_1), \dots, \hat{\gamma}(h_k) \in k$  bins
- ② Fit a parametric model to  $(\hat{\gamma}(h_1), \dots, \hat{\gamma}(h_k))$   
ie.  $\gamma(h; \theta)$  parametric model (spherical, gauss., lin. etc)  
and then

$$\min_{\theta} \sum_{i=1}^n [\hat{\gamma}(h_i) - \gamma(h_i; \theta)]^2$$

Better approach to

$$\min_{\theta} \sum_{i=1}^k w_i [\hat{\gamma}(h_i) - \gamma(h_i; \theta)]^2$$

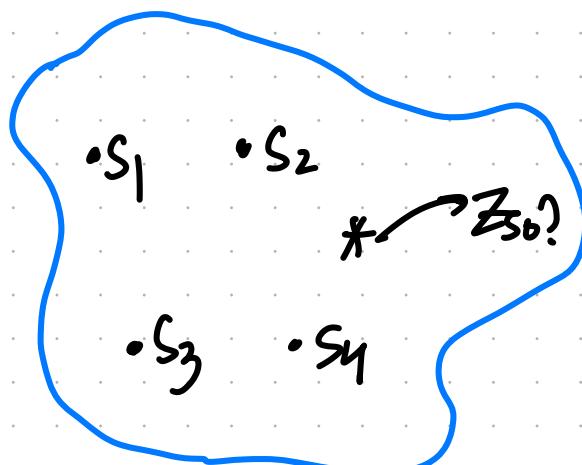
weight prop. to size of bin i

Even better approach

$$\hat{f}(h) = (\hat{y}(h_1), \dots, \hat{y}(h_k))^\top, \quad (\text{Generalised Least squares- GLS})$$

$$\min_{\theta} \left[ \frac{\hat{f}(h)}{\gamma(h; \theta)} \right] \text{cov}^{-1} \left[ \frac{\hat{f}(h)}{\gamma(h; \theta)} \right]$$

## \* Prediction (Kriging)



D

Given the position of  $Z_s$   
at locations  $Z_{s1}, \dots, Z_{sn}$   
 $\Rightarrow$  predict the  $Z_s$  at location  $Z_{s0}$

$\underline{z} = (Z_{s1}, \dots, Z_{sn})$  gives a prediction of  
 $Z_{s0}$

$f(\underline{z}) = \underline{z}_{s_0}^*$  (predicted)  $\rightarrow$  this is called Kriging

Goal is to minimize

$$E[(\underline{z}_{s_0} - \underline{z}_{s_0}^*)^2] = E[(\underline{z}_{s_0} - f(\underline{z}))^2]$$

$$\underline{z}_{s_0}^* = E[\underline{z}_{s_0} | \underline{z}_{s_1}, \dots, \underline{z}_{s_n}]$$

Properties: ①  $E[\underline{z}_{s_0}^*] = E[\underline{z}_{s_0}]$  unbiased

$$\textcircled{2} \quad \text{cov}(\underline{z}_{s_0}^*, \underline{z}_{s_0} - \underline{z}_{s_0}^*) = 0$$

③ If  $s_0 \in \{s_1, \dots, s_n\}$  then  $E[\underline{z}_{s_0} | \underline{z}_{s_1}, \dots, \underline{z}_{s_n}] = \underline{z}_{s_0}$   
(*interpolation*)

④ If  $\{\underline{z}_s, s \in D\}$  is gaussian then  $E[\underline{z}_{s_0} | \underline{z}_{s_1}, \dots, \underline{z}_{s_n}] = \lambda_0 + \sum_{i=1}^n \lambda_i \underline{z}_{s_i}$  (linear)

## Problems in Kriging

Find  $\hat{\lambda}_0, \hat{\lambda}_1, \dots, \hat{\lambda}_n$  such that  $z_{s_0}^* = \hat{\lambda}_0 + \sum_{i=1}^n \hat{\lambda}_i z_{s_i}$

satisfies :

$$1. E[z_{s_0}^*] = E[z_{s_0}]$$

$$2. (\hat{\lambda}_0^*, \dots, \hat{\lambda}_n^*) = \min_{(\lambda_0, \dots, \lambda_n)} E \left[ (z_{s_0} - (\lambda_0 + \sum_{i=1}^n \lambda_i z_{s_i}))^2 \right]$$

Only if  $z_{s_0}^*$  exist, then  $z_{s_0}^*$  is the BLUP for  $z_{s_0}$

(Best linear unbiased predictor)

$$z_{s_0}^* = \hat{\lambda}_0 + \sum_{i=1}^n \hat{\lambda}_i z_{s_i}$$

## Simple Kriging

Assume that  $E[z_s] = m_s$  is known for  $s \in \{s_0, s_1, \dots, s_n\}$   
then unbiasedness requires that

$$E[\lambda_0 + \sum \lambda_i z_{s_i}] = \lambda_0 + \sum \lambda_i m_{s_i}$$

||

$$E[z_{s_0}] = m_{s_0}$$

$$\hat{\lambda}_0 = m_{s_0} - \sum_{i=1}^n \hat{\lambda}_i m_{s_i}$$

Now we need to

$$\min_{(\lambda_1, \dots, \lambda_n)} E \left[ z_{s_0} - \left( \sum_{i=1}^n \lambda_i z_{s_i} + (m_{s_0} - \sum_{i=1}^n \lambda_i m_{s_i}) \right)^2 \right]$$

$$\therefore \min E \left[ ((z_{S_0} - m_{S_0}) - \left( \sum_{i=1}^n \lambda_i z_{S_i} - \sum_{i=1}^n \lambda_i m_{S_i} \right))^2 \right]$$

$$= \min \text{Var}(z_{S_0}) + \text{Var}(\sum \lambda_i z_{S_i}) - 2 \cdot \text{Cov}(z_{S_0}, \sum \lambda_i z_{S_i})$$

$$= \min \text{Var}(z_{S_0}) + \text{Var}(\sum \lambda_i z_{S_i}) - 2 \cdot \sum \lambda_i \text{Cov}(z_{S_0}, z_{S_i})$$

## Notations

$$\sigma_{00} = \text{Var}(z_{S_0})$$

$$\sigma_{0i} = \text{Var}(z_{S_0}, z_{S_i})$$

$$\underline{\sigma}_0 = (\sigma_0, \dots, \sigma_{0n})$$

$$\sigma_{ij} = \text{Cov}(z_{S_i}, z_{S_j}) \quad i, j = 1 \dots n$$

$$\Sigma = [\sigma_{ij}] \text{ nxn cov matrix of } (z_1, \dots, z_n)$$

$$\underline{\lambda} = (\lambda_1, \dots, \lambda_n)'$$

This document is available free of charge on

$$\therefore \min_{\lambda} \sum \lambda - 2\lambda C_0 + \sigma_{00}$$

$$\Rightarrow \sum \lambda = C_0 \Rightarrow \hat{\lambda} = \sum C_0 \text{ (if } \Sigma \text{ is inverted)}$$

$$z_{s0}^* = (m_{s0} - \sum \hat{\lambda}_i m_{si}) + \sum_{i=1}^n \hat{\lambda}_i z_{si}$$

so we need to know  $\Sigma$  and  $C_0$

## \* Ordinary Kriging

$$E[z_s] = m + SED$$

where  $m$  is unknown

Unbiasedness :

$$E[\lambda_0 + \sum \lambda_i z_{si}] = E[z_{s0}] \quad \therefore \hat{\lambda}_0 + \sum \hat{\lambda}_i m = m \\ + (\lambda_1, \dots, \hat{\lambda}_n, \hat{\lambda}_0)$$

$$\begin{aligned} \hat{\lambda}_0^* &= 0 \\ \sum_{i=1}^n \hat{\lambda}_i^* &= 1 \end{aligned} \quad \left. \right\} \text{constants}$$

$$= \begin{cases} \min_{(\lambda_1, \dots, \lambda_n)} E \left[ (z_{S0} - \sum \lambda_i z_{Si})^2 \right] \\ \sum \lambda_i = 1 \end{cases}$$

$$= \begin{cases} \min_{\lambda} \lambda' \Sigma \lambda - 2 \lambda C_0 + \sigma_{00} \\ \lambda' = 1 \end{cases}$$

$$= \min_{\lambda, \xi} \lambda' \Sigma \lambda - 2 \lambda C_0 + \sigma_{00} + \xi (\lambda' - 1)$$

$$= \begin{bmatrix} \Sigma & 1 \\ 1' & 0 \end{bmatrix} \begin{bmatrix} \lambda \\ \xi \end{bmatrix} = \begin{bmatrix} C_0 \\ 0 \end{bmatrix} = \hat{\lambda}$$

## Universal Kniging

$$E[Z_S] = a_0 + a_1 f_1(s) + \cdots + a_L f_L(s)$$

where

$a_0, a_1, \dots, a_L$  are coefficients - unknown

$f_1(s), \dots, f_L(s)$  are covariates - known & SED

Model  $Z_S = m_S + \delta(S)$  and  $E[\delta(S)] = 0$

Unbiasedness :  $E[\hat{\lambda}_0 + \sum \hat{\lambda}_i z_{S_i}] = E[Z_{S_0}]$

$$\hat{\lambda}_0 + \sum_{i=1}^n \hat{\lambda}_i (a_0 + a_1 f_1(s_i) + \cdots + a_L f_L(s_i)) = a_0 + a_1 f_1(s_0) + \cdots + a_L f_L(s_0)$$

$$\Rightarrow \hat{\lambda}_0 = 0 \text{ then } \sum \hat{\lambda}_i f_i(s_i) = f_i(s_0)$$

$$\sum \hat{\lambda}_i = 1 \text{ then } \sum \hat{\lambda}_i f_j(s_i) = f_j(s_0) \quad j = 1, \dots, L$$

Hence

$$\left\{ \begin{array}{l} \min \underline{\lambda}' \sum \underline{\lambda} - 2 \underline{\lambda}' c_0 + \sigma_{00} \\ \text{s.t.} \end{array} \right.$$

$$\sum \lambda_i = 1$$

$$\sum \lambda_i f_j(s_i) = f_j(s_0) \quad j=1, \dots, L$$

$$\Rightarrow \begin{bmatrix} \sum & F \\ F' & 0 \end{bmatrix} \begin{bmatrix} \underline{\lambda} \\ \underline{\xi} \end{bmatrix} = \begin{bmatrix} c_0 \\ f_0 \end{bmatrix}$$

where  $F = \begin{bmatrix} 1 & f_1(s_1) & \dots & f_L(s_1) \\ 1 & f_1(s_2) & \dots & f_L(s_2) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & f_1(s_n) & \dots & f_L(s_n) \end{bmatrix}$

(design matrix)

and where

$$f_0 = \begin{pmatrix} f_1(s_0) \\ \vdots \\ f_L(s_0) \end{pmatrix}$$

$$\therefore z_{s_0}^* = \sum \hat{\lambda}_i z_{s_i}$$

To solve simple, ordinary & universal kriging equations

we need to know  $\Sigma \delta c_0$

Assumptions:

for simple kriging:  $\{z_s - m_s, SED\}$  is weakly stationary

for ordinary kriging:  $\{z_s, SED\}$  is weakly stationary

for universal kriging:  $\{\delta(s), SED\}$  is weakly stationary

with these assumptions

• Use data  $(z_1, \dots, z_n)$  + covariates to estimate  $\hat{\delta}$

$$\delta(h) = C_0 - C(h) \quad \text{where } C_0 = \text{Var}(z_s)$$

$$C(h) = \text{cov}(z_s, z_{s+h})$$

$$\therefore \hat{C}(h) = \hat{C}(0) - \hat{\delta}(h) \quad \text{here } \hat{C}(0) \text{ is the sill}$$

$$\text{cov}(Z_{S_i}, Z_{S_j}) = \hat{\sigma}_{ij} = \hat{C}(0) - \hat{\delta}(s_i - s_j)$$

In universal kriging, for estimating  $\delta$ , you need to know  $\delta$

$$Z_s = a_0 + a_1 f_1(s) + \dots + a_L f_L(s) + \delta(s)$$

$$\Xi = (Z_{S_1}, \dots, Z_{S_n})' \quad \underline{a} = (a_0, a_1, \dots, a_L)$$

$$Z = F\underline{a} + \underline{\delta} \quad \text{linear model } E(\underline{\delta}) = 0$$

$$\underline{\delta} = (\delta(s_1), \dots, \delta(s_n)) \quad \text{cov}(\underline{\delta}) = \Sigma$$

- Estimate  $\underline{\alpha}$  with  $\hat{\underline{\alpha}} \longrightarrow \hat{\underline{\alpha}} = (F' \Sigma F)^{-1} F' \Sigma^{-1} \underline{z}$
- Estimate  $\delta = \underline{z} - F \hat{\underline{\alpha}}$
- Estimate  $\gamma$

## Iterative Process

$$\text{Step 0: } \hat{\underline{\alpha}}^{(0)} = (F' F)^{-1} F' \underline{z} \quad (\text{OLS})$$

$$\text{Step i: } \underline{\delta} = \underline{z} - F \hat{\underline{\alpha}}^{(k-1)}$$

$$\Rightarrow \text{Estimate } \gamma^{(i)} = \hat{\Sigma}^{(i)} = \hat{\underline{\alpha}}^{(i)} = (F' \hat{\Sigma}^{(i)} F)^{-1} F' \hat{\Sigma}^{(i)} \underline{z}$$

## Variance of Kriging predictors

- $\text{Var}(\hat{Z}_{s_0}(s_k)) = E[(\hat{Z}_{s_0} - \hat{Z}_{s_0}^*(s_k))^2]$   
 $= C(0) - \sum_{i=1}^n \lambda_i C(s_0 - s_i) = \sigma^2 s_k$
- $\text{Var}(\hat{Z}_{s_0}^*(0_k)) = C(0) - \sum \hat{\lambda}_i C(s_0 - s_i) - \hat{\xi} = \sigma^2 0_k$
- $\text{Var}(\hat{Z}_{s_0}^*(v_k)) = C(0) - \sum \hat{\lambda}_i C(s_0 - s_i) - \sum \hat{\xi}_j f_j(s_0) = \sigma^2 v_k$

## Markov Processes (Chains)

We have a process  $\{X_t, t \in [0, \infty)\}$  with  $X_t \in \mathbb{R}$  is called Markov if, given  $X_t$ , the distribution of  $X_s$  for all  $s > t$  does not depend of  $X_u$  for  $u < t$ .

Equivalently,

given  $X_t$ ,

$X_u$  with  $u < t$   
 $X_s$  with  $s > t$

} are independent

is called state space

We will consider discrete time process so  $t = n = 0, 1, 2, 3, \dots$

moreover  $X_n \in S$

discrete, finite, countable  $\{X_n\} \rightarrow$  markov chains

## Markov property

$$P[X_{n+1} = j \mid X_0 = i_0, X_1 = i_1, \dots, X_n = i_{n+1}] \neq i_0, i_1, \dots, i_n \in S$$

This probability will be =  $P[X_{n+1} = j \mid X_n = i_n]$

Notation :  $\forall i, j \in S$

↳ only dependent on previous element

$$P[X_{n+1} = j \mid X_n = i] = P_{ij} \text{ (transition probability)}$$

If  $P_{ij} \neq i, j \in S$  and  $\forall n$  are stationary, then

$P_{ij}^{(n, n+1)} = P_{ij}$  and the chain is said to be homogeneous.

• For  $i \in S$ , let

$P_i = P[X_0 = i]$  for initial periods ie.  $p_i \geq 0$  and  $\sum P_i = 1$

Notice that, because of the markov property,

$\# i_0, i_1, \dots \text{ in } t^S$

$$P[X_0 = i_0, X_1 = i_1, X_2 = i_2, \dots, X_n = i_n] = P_i P_{i_0} i_1, P_{i_1} i_2 \dots P_{i_{n-1}} i_n$$

Introducing a matrix (transition matrix)

$P = [P_{ij}]$  of order  $|S| \times |S|$

$$\therefore \sum_{j \in S} P_{ij} = 1 \quad \forall P_{ij} \geq 0$$

## Central to analysis of Markov chains

$$P_{ij}^{(n)} = P[X_n=j \mid X_0=i] = P[X_{m+n}=j \mid X_m=i]$$

$$P^{(n)} = [P_{ij}^{(n)}] \quad (n\text{-step transition matrix})$$

Set  $P_{ij}^{(0)} = \begin{cases} 0 & \neq i \neq j \\ 1 & \neq i = j \end{cases}$

**Proposition:** For all  $i, j \in S$ :

$$P_{ij}^{(n)} = \sum_{k \in S} P_{ik} P_{kj}^{(n-1)} \quad \text{ie. } P^{(n)} = P \cdot P^{(n-1)}$$

$$P^{(n)} = P \cdot P^{(n-1)} = \underbrace{P \cdot P \cdots P}_{n\text{-times}} = P^n$$

$$P^{(n)} = P^n$$

(This means that for the probability at  $n^{\text{th}}$  observations ( $P^{(n)}$ ) we need to multiply  $P$ ,  $n$  times ie.  $P^n$  (which =  $P^{(n)}$ ))

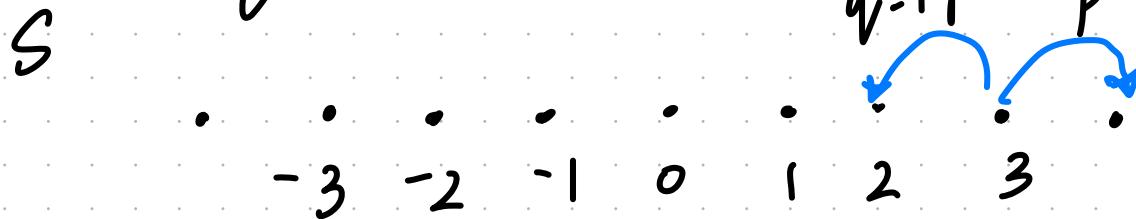
If  $0 \leq k \leq n$

$$\text{then } P^n = P^k \cdot P^{n-k} = P^{(k)} \cdot P^{(n-k)}$$

Chapman-Kolmogorov's  
equations

## Examples (1)

Considering a 1-D random walk (1)



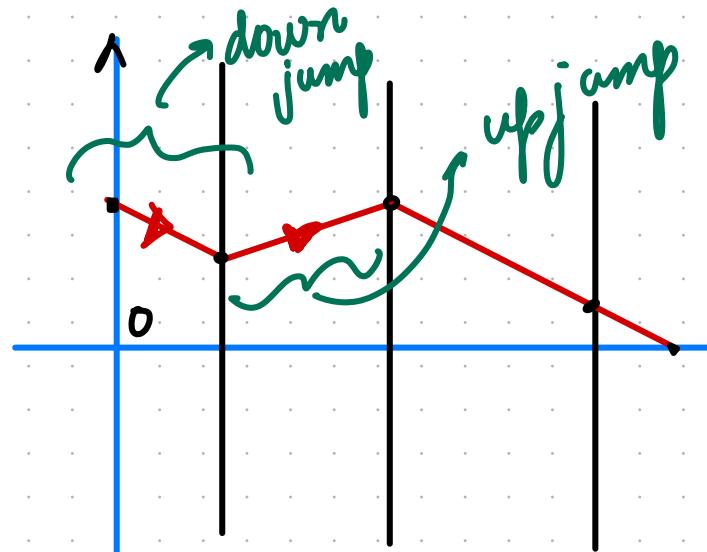
then the transition matrix

$$P = i \begin{bmatrix} 0 & 0 & q & 0 & p & 0 \\ q & 0 & p & 0 & p & - \\ 0 & p & - & 0 & 0 & q \end{bmatrix}$$

$$\therefore P_{ij}^{(n)} = P[X_n=j | X_0=i]$$

Let  $j^i z^i$

and  $u$ : sum of up jumps  
 $d$ : sum of down jumps



$$X_{n+1} = X_n + Y_n$$

where  $Y_1, Y_2 \dots Y_n$  = iid

$$Y_i = \begin{cases} +1 & p \\ -1 & 1-p \end{cases}$$

$$\therefore u = \sum_{K=1}^n 1[Y_K = +1] \text{ (jump-up)}$$

$$d = \sum_{K=1}^n 1[Y_K = -1] \text{ (jump-down)}$$

$\therefore n+j-i$  must be even

$$\left. \begin{aligned} \text{here } u+d &= n \\ u-d &= j-i \end{aligned} \right\} \text{ adding them}$$

$$2u = n+j-i$$

$$u = \frac{1}{2}(n+j-i)$$

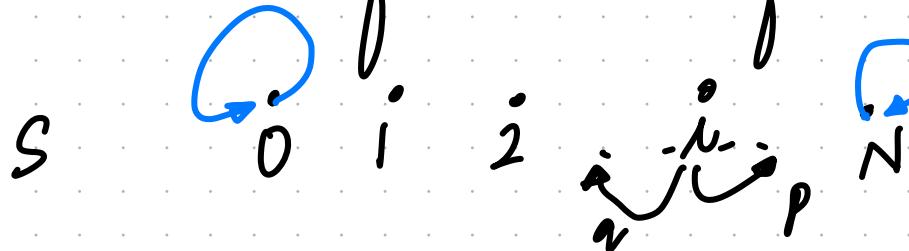
$$P_{ij}^{(n)} = \begin{cases} 0 & \text{if } n+j-i \text{ is odd} \\ \binom{n}{\frac{1}{2}(n+j-i)} p^{\frac{1}{2}(n+j-i)} \times (1-p)^{n-\frac{1}{2}(n+j-i)} & \text{if } n \text{ is even} \end{cases}$$

$$\text{Eq: } P_{00}^{(2n)} = \binom{2n}{n} p^n (1-p)^n$$

*probability of coming back from where you started*

## ② Gambler's ruin

Set  $S$  is a finite no. of integers



If you reach  $0$ , you stay at  $0$

If you reach  $N$ , you stay at  $N$

only if you are at  $i > 0$  and  $i < N$ , you can jump up or down

where  $i$  is current gambler fortune

$N-i$  is fortune of adversary

## Transition matrix

$$P = \begin{bmatrix} 0 & 1 & 2 & & N \\ 1 & 0 & - & & 0 \\ 0 & q & 0 & p & - \\ 0 & 0 & q & 0 & p & - \\ \vdots & & & & & \\ 0 & - & 0 & q & 0 & p \\ 0 & - & - & & & 1 \end{bmatrix}$$

$P_{ij} \rightarrow$  compute  $P^{(n)}$  with R

if  $T = \inf \{ n : X_n = 0 \text{ or } X_n = N \}$

$$P[X_T = 0 | X_0 = i] \quad i = 0, 1, \dots, N$$

probability of gambler going broke

$$\lambda_0 = 1 \quad (\text{probability of going broke with 0 euros})$$

$$\lambda_1 = q + p \cdot \lambda_2$$

$$\lambda_2 = q \cdot \lambda_1 + p \cdot \lambda_3$$

$$\therefore \lambda_{N-1} = q \cdot \lambda_{N-2} + p \cdot 0$$

(p of going broke with N euros)

$$\boxed{\lambda_N = 0}$$

Exercise : When  $N=3$ , prove that  $\lambda_1 = \frac{q}{1-pq}$  and  $\lambda_2 = \frac{q^2}{1-pq}$

$$\text{if } p=q=\frac{1}{2}, \lambda_1 = \frac{2}{3}, \lambda_2 = \frac{1}{3}$$

For general  $N$

$$\lambda_i = \begin{cases} \frac{N-i}{N} & \text{if } p=q=\frac{1}{2} \\ \frac{\left(\frac{q}{p}\right)^i - \left(\frac{q}{p}\right)^N}{1 - \left(\frac{q}{p}\right)^N} & \text{if } p \neq q \end{cases}$$

Here if  $N \rightarrow \infty$

$$\lim_{N \rightarrow \infty} \lambda_i = 1 \text{ if } p=q$$

$$\lim_{N \rightarrow \infty} \lambda_i = 1 \text{ if } q > p$$

$$\lim_{N \rightarrow \infty} \lambda_i = \left(\frac{q}{p}\right)^i \quad p > q$$

### ③ Inventory problem

Consider a warehouse of  $C$  stock units

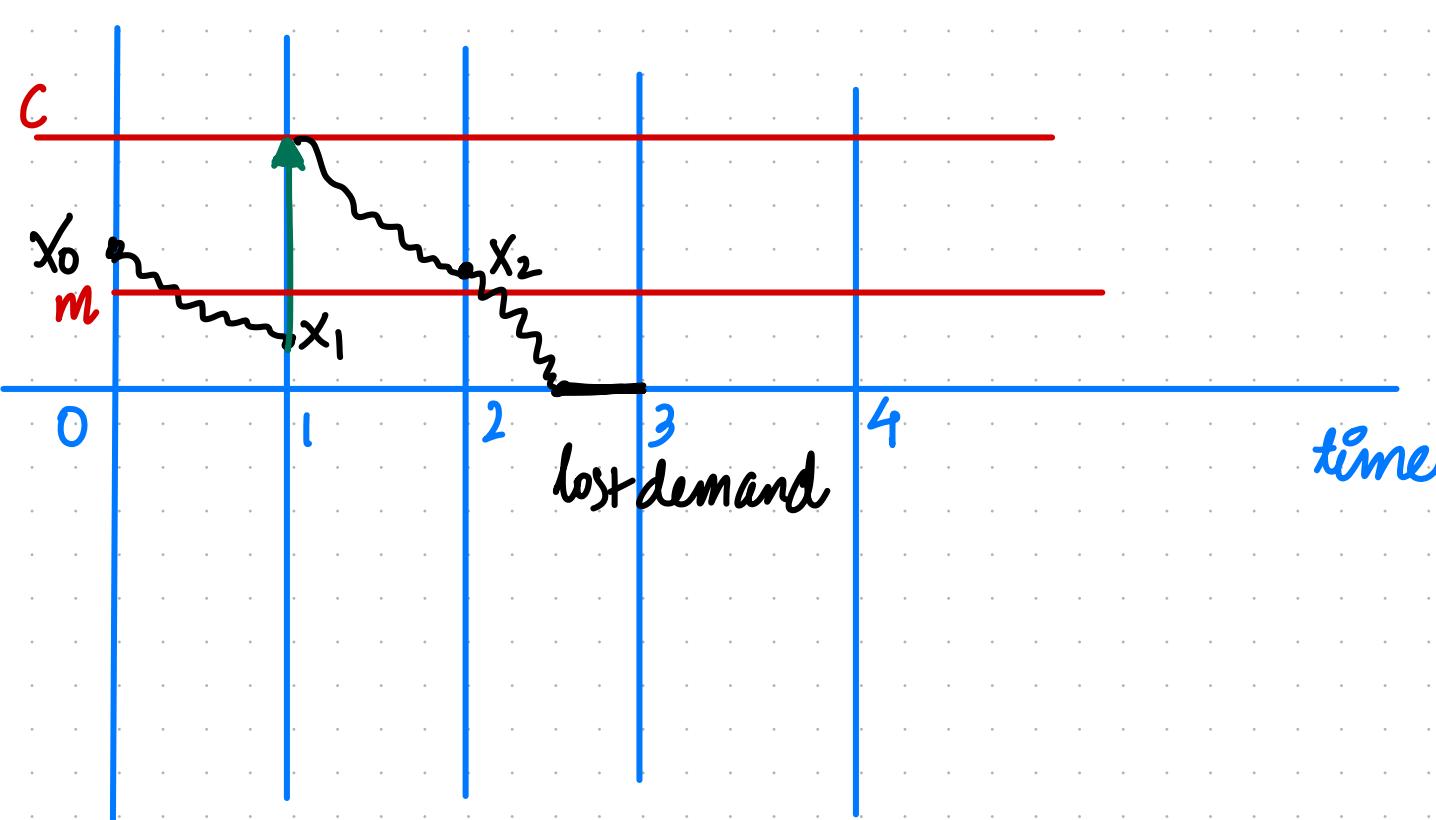
Every time period  $n$ :  $D_n \in \{0, 1, 2, \dots\}$  demand  
met, if possible

$X_n$ : residual stock in warehouse at end of period  $n$ .

Assume  $D_1, D_2, \dots, D_n$  iid

Restock strategy : let  $0 \leq m \leq C-1$

- if  $X_n > m$  don't do anything
- if  $X_n \leq m$  restock to full capacity  $C$



$$x_{n+1} = \begin{cases} C - D_{n+1} & \text{if } 0 \leq x_n \leq m \\ (x_n - D_{n+1}) & \text{if } m < x_n \leq C \end{cases}$$

is not dependent upon  $C$  and instead on previous obs.  $\therefore$  it is a markov chain

## Relevant Questions

$$\textcircled{1} \lim_{n \rightarrow \infty} P_{C_0}^{(n)}$$

$$\textcircled{2} \lim_{n \rightarrow \infty} \sum_j P_{Cj}^{(n)}$$

(prob. of finding warehouse  
empty at EOD)

$$\textcircled{3} \text{ long run avg. demand that is not met } \left( \lim_{n \rightarrow \infty} \sum_{i=0}^c u_i P_{Ci}^{(n)} \right)$$

$D_1, D_2, \dots, D_n$  iid

$$X_{n+i} \in \{0, 1, 2, \dots, c\}$$

$$M_i = \begin{cases} |D - c| & 0 \leq i \leq m \\ |D - i| & m < i \leq c \end{cases} \quad \text{where } u_i \in \text{Expected unmet demand at end of period } n+1$$

## Long run behaviour of Markov chains

Consider a markov chain  $\{X_n\}$  on  $S = \{1, 2, \dots, N\}$

Assume  $K \geq 1$  such that  $P^{(K)}$  has all non-negative entries

i.e.  $\forall i, j \in S$   
 $p_{ij}^{(K)} > 0$

Such a markov chain is called regular.

**Proposition:** If  $\{X_n\}$  is regular, for all  $i, j \in S$

$$\lim_{n \rightarrow \infty} p_{ij}^{(n)} = \pi_j > 0$$

does not depend on  $i$

Moreover

$$\sum_{j \in S} \pi_j = 1$$

$\Pi = (\pi_1, \dots, \pi_N)'$  is a probability distribution on  $S$  and is the solution of this eigen equation.

$$\Pi' = \pi' P$$

i.e.  $\Pi'$  is a left eigen vector of  $P$  corresponding to eigen vector  $1$ .

## Interpretation

- $\pi_j$  is the long term probability of finding the process  $\{X_n\}$  in state  $j$ .
- $\pi_j$  is the long run mean fraction of time that the chain spends in state  $j$ .

$$\text{Consider : } E\left[\frac{1}{m} \sum_{n=0}^m \pi [x_n=j]\right] = \frac{1}{m} \sum_{n=0}^m p_{ij}^{(n)}$$

WKT when  $m \rightarrow \infty$ ;  $p_{ij}^{(n)} \rightarrow \pi_j$

$$\therefore E\left[\frac{1}{m} \sum_{n=0}^m \pi [x_n=j]\right] = \pi_j$$

when  $m \rightarrow \infty$

## Example : Inventory problem

$$X_{n+1} = \begin{cases} (c - D_n) & \text{if } 0 \leq X_n \leq m \\ (X_n - D_n) & \text{if } m < X_n \leq c \end{cases}$$

Now if  $D \sim D_1, D_2, \dots, D_n$  iid is such that

$P[D=k] > 0$  for  $k = 0, 1, 2, \dots, M$  with  $M \geq c$

$$S = \{0, 1, \dots, c\}$$

let  $\pi = (\pi_0, \dots, \pi_c)'$  such that  $\pi' = \pi' P$

then

- $\pi_0$  (long run freq. of residual stock) is 0
- $\sum_{j \in S} j \pi_j$  (long run average no. of items in stock)

- Enforced un-met demand

For  $i \in S$

$$u_i^o = \begin{cases} E[D - C] & \text{if } 0 \leq i \leq m \\ E[D - i] & \text{if } m < i \leq C \end{cases}$$

- Long run average of un-met demand

$$u(m) = \sum_{i \in S} u_i^o \pi_i \quad \text{here } u(m) \downarrow \text{as } m \uparrow$$

Suppose if  $a$  is the cost of restocking and  $b$  is the profit for each item sold, then

Long run average of restocking

$$g(m) = \sum_{j=0}^m \pi_j \quad \text{if } m(1) \text{ then } g(m)(1)$$

$$\text{cost: } a \sum_{j=0}^m \pi_j + b \cdot \sum_{i \in S} u_i \pi_i$$

Here the optimization problem would be to find  $m^*$  such that  $m^* = \min \text{cost}(m)$

Exercise: Find the optimal 'm' when  $c=3$  (possible values for  $m$  are 2, 1, 0) ;  $D \in \{0, 1, 2, 3, \dots\}$  and  $P[D=k] = \frac{1}{2^{k+1}}$   $k=0, 1, \dots$

Solution :

$$m = \begin{cases} 2 & \text{if } a \leq \frac{1}{4} \\ 1 & \text{if } \frac{1}{4} < a \leq 1 \\ 0 & \text{if } a \geq 1 \end{cases}$$

Definitions Let  $i, j \in S$

$j$  is reachable from  $i$ , if  $n \geq 1$  such that  $P_{ij}^{(n)} > 0$

If  $j$  is reachable from  $i$   
 $i$  is reachable from  $j$

$\left. \begin{array}{c} j \\ i \end{array} \right\}$   $i$  &  $j$  communicate  
 $(i \leftrightarrow j)$

## Proposition

$\leftrightarrow$  is an equivalence relation for states in  $S$

- ie - reflexive
- transitive
- symmetric

Hence since it is an equivalence relation

$$S = C_0 \cup C_1 \cup C_2 \dots \cup C_R \cup \dots$$

$C_i$  equiv classes induced by  $\leftrightarrow$

ie  $c_i \neq \emptyset$ , and  $c_i \cap c_j = \emptyset \forall i, j$

• Starting in one class  $c_i$ , the process can move to another class  $c_j$  but it cannot come back

In case of Gambler's ruin, classes would be  $\{0\}, \{1, \dots, N-1\}, \{N\}$

Definition: If  $S = C$  (one equivalence class)

then Markov chain is said to be reducible

Definition: Period of a state  $i \in S$

Period is a max. common divisor of  $n$  such that  $P_{ii}^{(n)} > 0$

Example: In a simple random walk

$$\begin{pmatrix} i \xrightarrow{\cdot} i+1 \\ i+1 \xrightarrow{\cdot} i \end{pmatrix}$$

$i-1 \quad i \quad i+1$  the only way to come back to the  
original pos  $\xrightarrow{\cdot}$  2 steps  $\therefore$  Period = 2

Proposition: All the states in the same class have the same period.

Definition: If  $i \in S$  is said to be recurrent if

$$P[\text{returning to } i \mid X_0 = i] = 1$$

(belongs to a single state)

If  $P[\text{returning to } i \mid X_0 = i] < 1$ , it is called transient

(In case of transient state, after some runs, the probability of returning back becomes zero)

Proposition: If  $i \leftrightarrow j$  then either  $i \& j$  are recurrent or  $i \& j$  are transient.

Example: Gambler's ruin

$$\{0\} \{1, 2, \dots, N-1\} \{N\}$$

0 is recurrent

$\{1, 2, \dots, N-1\}$  is transient

N is recurrent

(because eventually you either reach 0 or N)

Random walk

If  $p = q = \frac{1}{2}$  then all states are recurrent

$p \neq q$  then states are all transient

let

$T_i = \inf\{n : X_n = i \mid X_0 = i\}$  time of first return

$$E[T_i \mid X_0 = i] = \mu_i^*$$

If  $i$  is transient,  $\mu_i^* = \infty$

If  $i$  is recurrent,

$\mu_i^* = \infty$  (null recurrent)

$\mu_i^* < \infty$  (positive recurrent)

\* All states in the same class are only one b one of this:

- ① Transient
- ② positive recurrent
- ③ Negative recurrent

Theorem: If a Markov chain is reducible and the states are positive recurrent then  $\pi = (\pi_0, \pi_1, \dots)$  distribution on  $S$  which is an unique solution of  $\pi' = \pi'P$  and  $\pi_i = \frac{1}{\mu_i}$

Here  $\pi_i$  is called the stationary distribution of chain.

Theorem: If a Markov chain is irreducible and aperiodic then

$$\pi_j^* = \lim_{n \rightarrow \infty} p_{ij}^{(n)} = \frac{1}{\mu_i}$$

## \* Markov Decision Processes

We are in a state space  $S$

$\{X_n\}$  is a process such that  $X_n \in S$  for  $n = 1, 2, 3, \dots$

and  $A$  is an action set such that  $A = \{a_1, a_2, \dots, a_m\}$

$q : S \times A \rightarrow$  prob. distribution on  $S$

i.e.

$$P[X_{n+1} = j | X_n = i, \sigma_n = a] = q(s, a)(j) \quad [\text{Law of motion}]$$

$\sigma_{\text{strategy}} = (\sigma_0, \sigma_1, \sigma_2, \dots)$  where  $\sigma_n \rightarrow h^n$  (history)

$$\sigma_n : (s_0, a_0), (s_1, a_1), \dots, (s_{n-1}, a_{n-1}), s_n$$

$u : (S \times A) (S \times A) \dots \rightarrow \mathbb{R}$   
(utility)

i.e. if  $h = ((s_0, a_0), (s_1, a_1), \dots)$   
then  $u(h)$  is the gain.

Goal: Find  $\sigma^*$  such that  $\max_{\sigma} E[u(h)]$

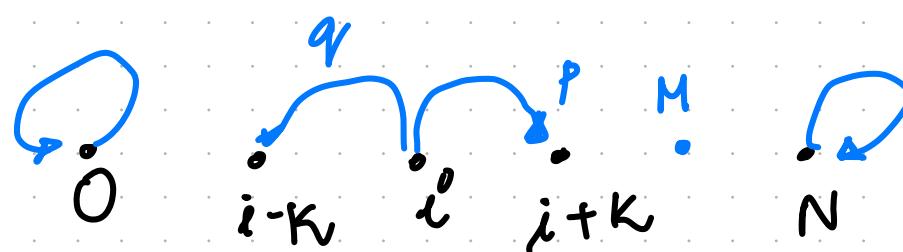
Example:

$$u(h) = \sum_{K=1}^{\infty} \beta^{K-1} u(s_K)$$

$$u : S \rightarrow \mathbb{R}$$

Example: Extension of Gambler's ruin

→ Red & Black



Bids: (action)  $K = X$  euro  $\leq i$  if  $X_n$  is in  $t$

$$u = \begin{cases} 1 & \text{if } h \text{ is s.t } X_n \text{ reached } M \leq N \\ 0 & \text{otherwise} \end{cases}$$

Optimal strategy:

⇒ Bold play: Play everything but do not overshoot the target  
(maximize probability of reaching M)

## \* Hidden Markov Models (HMM)

A process made of 2 random variables  $\{(X_n, Y_n)\}_{n=0,1,2,\dots}$

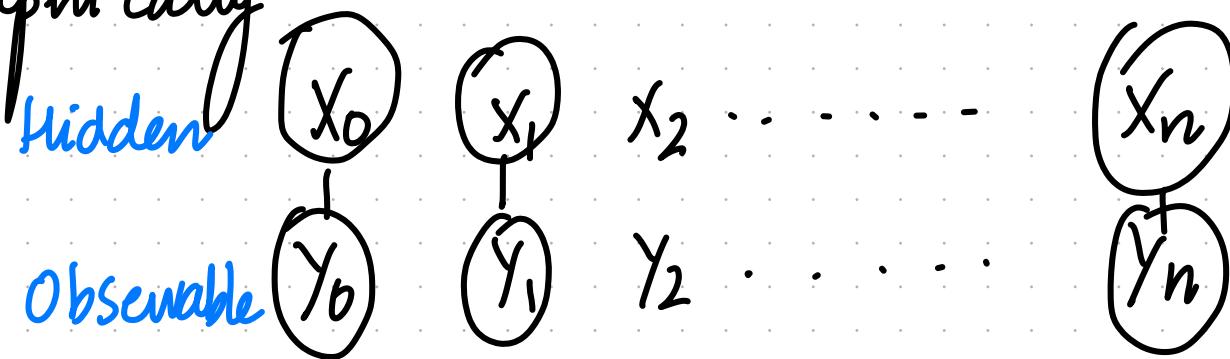
$\{X_n\}$  is a MC on S - state space finite

Conditionally on  $\{x_n\}$

$y_0, y_1, \dots, y_n$  are independent (random variables with values in  $\mathbb{R}$ )

and  $P(y_n \in A | x_0, x_1, \dots, x_{n-1}, \dots) = P(y_n \in A | x_n)$

Graphically



Ingredients :

→ state space of chain  $\{x_n\}$ ,  $S = \{s_1, s_2, \dots, s_N\}$   
→ finite

→ observables : finite  $\mathcal{O} = \{v_1, \dots, v_m\}$  infinite  $R: \mathbb{Z}, \mathbb{N}$

- $P$ : transition matrix on  $S$
- Initial distribution :  $\lambda$  is a dist. on  $S$   
 $\lambda(i) = P[x_0=i] \quad i \in S$
- Emission Matrix: if  $\mathcal{O}$  is finite :  $b_j(v_k) = P[x_n=v_k | x_n=j]$   
if  $\mathcal{O}$  is  $R$  :  $b_j$  is the dist. of  $y_n | x_n=j$   
Eg:  $y_n | x_n=j \sim N(\mu_j, \sigma_j^2)$  is gaussian  
we will have a matrix  $B$  ( $N$  rows  $\times M$  col.) if  $\mathcal{O}$  is finite

$$B = \begin{bmatrix} b_1(v_1) & b_2(v_2) & \dots & b_n(v_n) \\ \vdots & \vdots & & \vdots \\ b_N(v_1) & b_N(v_2) & \dots & b_N(v_N) \end{bmatrix}$$

• if  $\sigma$  is R &  $Y_n | X_n=j \sim N(\mu_j, \sigma_j^2)$  then

$$B = \begin{bmatrix} N(\mu_1, \sigma_1^2) \\ N(\mu_2, \sigma_2^2) \\ \vdots \\ N(\mu_N, \sigma_N^2) \end{bmatrix}$$

### Example:

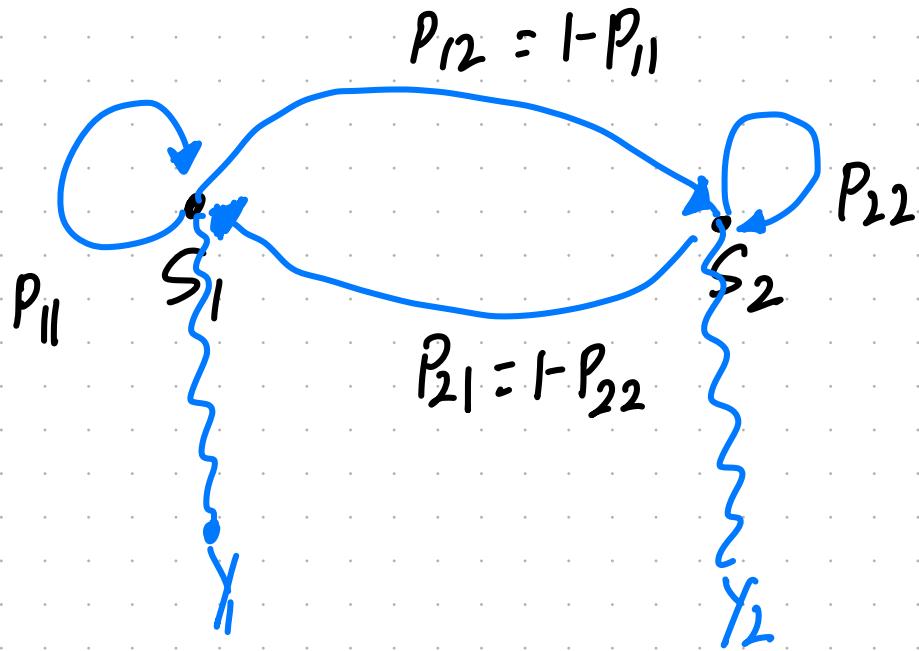
we have a signal data: 100 1110010 ...

Here  $X_n$  iid  $\sim$  Bernoulli( $p$ )

this is called HMM-0 (order 0)

If  $X = \{s_1 \quad \pi_1\}$  then  $Y|s_1 \sim \text{Bernoulli}(p_1)$  this is HMM-1  
 $Y|s_2 \sim \text{Bernoulli}(p_2)$

## HMM-2

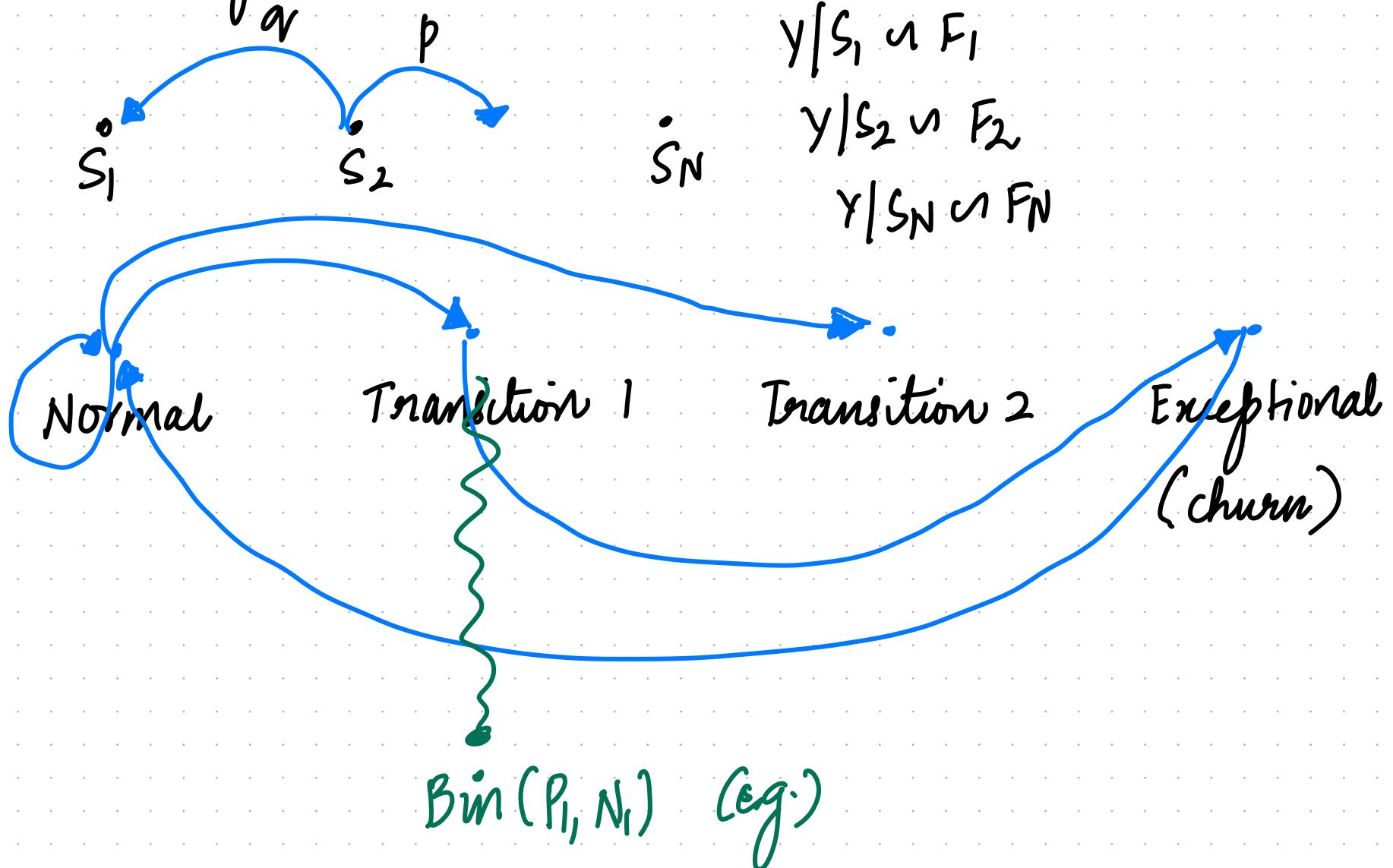


$Y | S_1 \sim Ben(p_1)$

$Y | S_2 \sim Ben(p_2)$

## A general example

We have a state space  $S$



## Problems

- Given a model  $m = (\lambda, P, B)$

compute the probability

$$P(y_0 = y_0, \dots, y_T = y_T) \text{ for all } T=0, 1, 2, \dots$$

Baum's Algorithm

- Estimate optimally the history of process  $\{X_n\}$  upto time  $T$ ,

$T \geq 1$  ie.  $X_0 = x_0, \dots, X_T = x_T$  (smoothing problem)

Viterbi's Algorithm

Given a model,  $m = (\lambda, P, B)$

(Incomplete)

- Estimate  $m$  (computationally heavy)  
(Baum-welch Algorithm), estimates of parameters in  $m$  through ML.

## \* State space Models and Kalman Filters

## Linear Models - Interaction terms

Consider the linear model

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2$$

$$m_1 \leftarrow \text{lm}(y \cup x_1 + x_2) \quad [\hat{y} = b_0 + b_1 x_1 + b_2 x_2]$$

$$m_2 \leftarrow \text{lm}(y \cup x_1 * x_2) \quad [\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_1 x_2]$$

## AIC & BIC

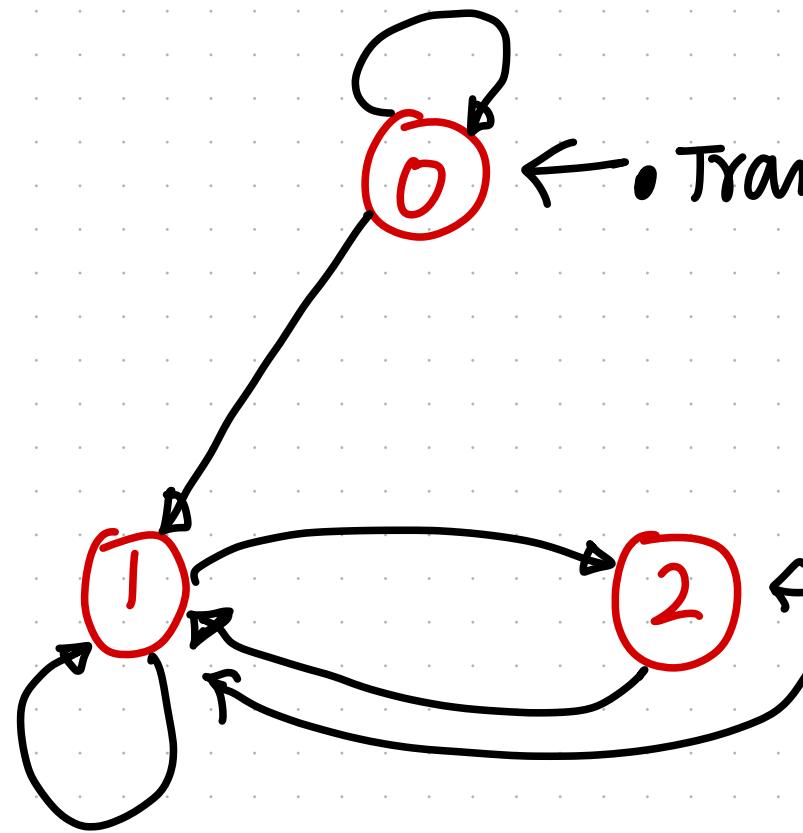
$$AIC = 2 \times \text{no. of parameters} - 2 \times L\text{-Ratio}$$

$$BIC = \text{no. of parameters} \times \log(\text{sample size}) - 2 \times L\text{-ratio}$$

$[0, \infty)$

If sample size  $> 1$ ;  $BIC > AIC$

# Markov chains & Markov Models



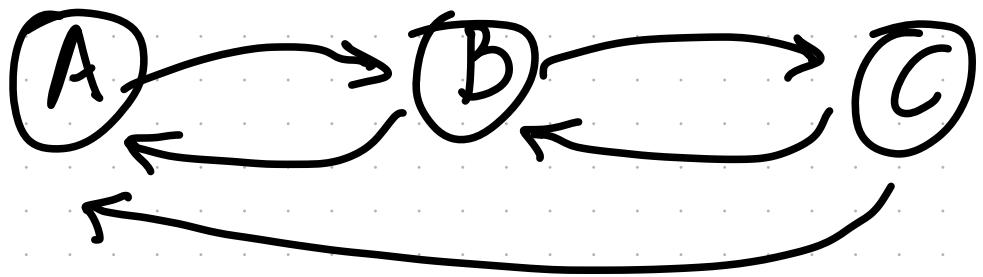
• Transient state: [If left from the state then no chance of coming back]

• Recurrent state

: [There is a probability of coming back in a random walk]

- Reducibility : If some states are unreachable from others
- Irreducibility : If all states can be left & returned back
- Absorbing: One entered, cannot be left

• If there is 1 self loop, it makes the whole chain aperiodic. (period=1)



Possible paths:

$$A \rightarrow B \rightarrow A - 2$$

$$A \rightarrow B \rightarrow C \rightarrow A - \textcircled{3}$$

$$A \rightarrow B \rightarrow C \rightarrow B \rightarrow A - \textcircled{4}$$

If period=1 ; (aperiodic)

$$\text{Period} = \text{GCD}(2, 3, 4) = 1$$