

Bio stat 250c

Class Notes



WEEK 1

Biostat 250C Lec 1

Spring 2021 TA: Mr. Rogers

Introduction to Graphical Modelling by David Edwards
(springer)

Intro. to multivariate statistical analysis by T.W. Anderson

Univariate $\leftarrow X \sim F(\cdot)$. $P(X \in A) \text{ or } P(X \leq x)$

$$\underbrace{X_1, \dots, X_n}_{\text{Realization from } F(\cdot)} \stackrel{iid}{\sim} F(\cdot)$$

$$\textcircled{Ex} \quad X \sim N(\mu, \sigma^2) \quad X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

You can not model using iid realizations

$$X = \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix} \sim F(\cdot) \quad (\text{multivariate})$$

$$P(X_1 \in A_1, X_2 \in A_2, \dots, X_n \in A_n) \leftarrow \begin{matrix} \text{How} \\ \text{construct it?} \end{matrix}$$

Elements of X are dependant! Model dependency.

$$n=2: \quad P(X, Y) = P(X) \cdot P(Y|X)$$

Notations:

$$[X, Y] = [X][Y|X]$$

$$[X] = \int [X, Y]; \quad [Y] = \int [X, Y]$$

$$[X][Y|X] = [X, Y] = [Y][X|Y]$$

$$\Rightarrow [X|Y] = \frac{[X][Y|X]}{[Y]}$$

$$\Rightarrow \underbrace{[X|Y]}_{\text{Bayes' theorem}} \propto [X][Y|X] = [X, Y]$$

Bayes' theorem

Q: Suppose I give you

$$[X|Y] \quad \& \quad [Y|X]$$

Let us suppose $[X, Y]$ exists.

$$\int \int [X, Y] = 1$$

- Is it possible to find $[X, Y]$ from $[X|Y]$ & $[Y|X]$?

$$\frac{[Y|X]}{[X|Y]} = \frac{[Y]}{[X]} \Rightarrow \int \frac{[Y|X]}{[X|Y]} = \frac{1}{[X]}$$

$$\Rightarrow [X] = \left(\int \frac{[Y|X]}{[X|Y]} \right)^{-1}$$

$$\Rightarrow [X, Y] = [X][Y|X]$$

$$= [Y|X] \left(\int \frac{[Y|X]}{[X|Y]} \right)^{-1}$$

What is the key?

$$\frac{[X]}{[Y]} = \frac{[X|Y]}{[Y|X]} \Rightarrow [X] = \frac{[Y][X|Y]}{[Y|X]}$$

$$\text{or } [X] \propto \frac{[X|Y]}{[Y|X]}.$$

I can evaluate $[X]$.

DO NOT NEED INTEGRATION
EXPLICITLY!

Conditional Independence

$X \perp Y$ or $X \amalg Y$ iff

$$P(X, Y) = P(X) \cdot P(Y) \text{ or}$$

$$[X, Y] = [X] \cdot [Y]$$

$$\Leftrightarrow \underbrace{\frac{[X, Y]}{[Y]}}_{\text{def}} = [X]$$

$$\Leftrightarrow [x|Y] = [x]$$

$$\Leftrightarrow [Y|x] = [Y]$$

Intro. a 3rd variable Z

$$X \amalg Y | Z \stackrel{\text{def}}{\Leftrightarrow} [x, Y | Z] = [x|Z][Y|Z]$$

or $P(X, Y | Z) = P(X|Z) \cdot P(Y|Z)$

what about $[x|Y, Z]$?

$$[x|Y, Z] = [x|Z]$$

Because

$$[x, Y | Z] = [Y | Z] [x | Y, Z]$$

$$\xrightarrow{\text{C.I.}} = [Y | Z] [x | Z]$$

Decomposition property of Cond. Indep.

$$X \amalg (Y, W) | Z \Rightarrow$$

$$X \amalg Y | Z \wedge X \amalg W | Z$$

$$\text{Proof: } [X, Y, W | Z] =$$

$$\underset{\text{def}}{=} [x|Z][Y|x, Z][w|x, Y, Z]$$

$$= [x|Z][Y, w|Z]$$

Integrate/marginalize out W from both sides.

$$\text{LHS: } \int [x, Y, w | Z] = [x, Y | Z]$$

$$\text{RHS: } \int [x|Z][Y, w | Z] = [x|Z][Y | Z]$$

Done!

HW: Given $X \amalg (Y, w) | Z$, prove $X \amalg w | Z$ using factorization & integration.

HW: $X \amalg (Y, w) | Z$, prove $X \amalg Y | w, Z$

HW: Contraction.

$$(X \amalg w | Z, Y) \wedge (X \amalg Y | Z) \Rightarrow$$

$$X \amalg (Y, w) | Z$$

Important HW rule:

Any HW given in a week is due on or before the beginning of class of the Wed. of the following week.

Comment: If $[x, Y]$ exists
 $[x|Y]$ & $[Y|x]$ defined uniquely

But what if $[x, Y]$ does NOT \exists ?

$$p(x, y) \propto e^{-\frac{1}{2}(x-y)^2}$$

$$\iint p(x, y) dx dy = +\infty.$$

$$[x|Y] \sim N(y, 1), [Y|x] \sim N(x, 1)$$

WEEK 5 Lec 9 04-26

Determinant version of SWM

Recall

$$\begin{aligned} \textcircled{1} & [I \ 0] [A \ B] = [A \ B] \\ \textcircled{2} & [X \ I] [C \ D] = [C \ D] \end{aligned}$$

$$\textcircled{2} \leftarrow \textcircled{2} + X\textcircled{1}$$

$$[C:D] = [C:D] + X[A:B]$$

Row(i) of $[C:D] \leftarrow$

$$\text{Row}(i) \text{ of } [C:D] + \underbrace{\text{row}(i) \text{ of } X}_{\text{row}(i) \text{ of } X[A:B]}$$

- $C = AB$ (Digression)

$$\text{row}(i) \text{ of } C = (\text{row } i \text{ of } A) \cdot B$$

Matrix multiplication:

$$\begin{aligned} C &= [C_{*1} : C_{*2} : \dots : C_{*p}] \\ &= \begin{bmatrix} C_{1*}^T \\ C_{2*}^T \\ \vdots \\ C_{m*}^T \end{bmatrix} \xrightarrow{\text{row 2 of } C} \text{column } p \text{ of } C \end{aligned}$$

$$C^T = [c_{1*} \ c_{2*} \ \dots \ c_{m*}] = \begin{bmatrix} c_{*1}^T \\ \vdots \\ c_{*p}^T \end{bmatrix}$$

$$C = AB \Rightarrow$$

$$C_{ij} = (\text{row } i \text{ of } A) \cdot (\text{col } j \text{ of } B)$$

$$= a_{i*}^T b_{*j}$$

$$\begin{aligned} C_{ij}^T &= a_{i*}^T [b_{*1} ; b_{*2} ; \dots ; b_{*p}] \\ &= a_{i*}^T B \end{aligned}$$

$$C_{ij} = A b_{*j} \quad \text{Similarly.}$$

- $\begin{bmatrix} a_{1*}^T x_{*j} \\ \vdots \\ a_{p*}^T x_{*j} \end{bmatrix} = \begin{bmatrix} a_{1*}^T \\ \vdots \\ a_{p*}^T \end{bmatrix} x_{*j}$ NEED A PROOF!

(Another Digression)

$$x^T A = (x_1, \dots, x_m) \begin{pmatrix} a_{1*}^T \\ \vdots \\ a_{m*}^T \end{pmatrix}$$

$$= \sum_{i=1}^m x_i (\text{row } i \text{ of } A)$$

$x^T A$: linear comb. of rows of A

Ax : linear comb. of cols of A .

Back to Block Elimination.

$$\begin{bmatrix} I & 0 \\ X & I \end{bmatrix} \begin{bmatrix} A & B \\ C & D \end{bmatrix} = \begin{bmatrix} A & B \\ XA+C & XB+D \end{bmatrix}$$

$$[XA+C : XB+D] = X[A:B] + [C:D]$$

So As a transformation,

$$[C:D] \leftarrow [C:D] + X[A:B]$$

$$\text{row } (i) \text{ of } [C:D] \leftarrow$$

$$\text{row } (i) \text{ of } [C:D] + \underbrace{(\text{row } (i) \text{ of } X)[A:B]}_{\text{linear comb. of rows of } [A:B]}$$

$\begin{bmatrix} I & 0 \\ X & I \end{bmatrix}$: execute a series of type-I row operations as a block.

$\begin{bmatrix} I & O \\ X & I \end{bmatrix}$: execute a series of type-I row operations as a block.

• Type I op:

$$\text{row}(i) \leftarrow \text{row}(i) + \alpha \text{row}(j)$$

$$\text{row}(i) \leftarrow \text{row}(i) + \sum_{k=1}^r \alpha_k \text{row}(j_k)$$

This is a series of elem. row-ops.

$$\text{row}(i) \leftarrow \text{row}(i) + \alpha_1 \text{row}(j_1)$$

$$\text{row}(i) \leftarrow \text{row}(i) + \alpha_2 \text{row}(j_2)$$

& so on.

• For k in $1:r$ {

$$\text{row}(i) \leftarrow \text{row}(i) + \alpha_k \text{row}(j_k)$$

}

$$\bullet \begin{bmatrix} A & B \\ C & D \end{bmatrix} \rightarrow \begin{bmatrix} \alpha A & \alpha B \\ BC & BD \end{bmatrix}$$

Is this a seq./block of elementary ops? If so, how/why?

Yes: $\begin{bmatrix} \alpha I & O \\ O & \beta I \end{bmatrix} \begin{bmatrix} A & B \\ C & D \end{bmatrix}$

$$\begin{bmatrix} I & O \\ O & \beta I \end{bmatrix} \begin{bmatrix} \alpha I & O \\ O & I \end{bmatrix}$$

X

• Back to Determinant SWM.

$$\underbrace{\begin{bmatrix} I & O \\ -CA^T & I \end{bmatrix}}_{\text{Block elementary operation}} \begin{bmatrix} A & B \\ C & D \end{bmatrix} = \begin{bmatrix} A & B \\ O & D-CA^T B \end{bmatrix}$$

Block elementary operation

Recall elementary row operation of the form

$$\text{row}(i) \leftarrow \text{row}(i) + \alpha \text{row}(j)$$

\Rightarrow

$$\det \begin{pmatrix} A & B \\ C & D \end{pmatrix} = \det \begin{pmatrix} A & B \\ O & D-CA^T B \end{pmatrix}$$

(*) Most mature def. (algebraic)

$$\det(A) = \sum_{\pi}^{\text{sign}(\pi)} (-1)^{\pi} a_{1\pi_1} a_{2\pi_2} \dots a_{n\pi_n}$$

Ex $A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$

π	term	sign
(1,2)	$a_{11} a_{22}$	+1
(2,1)	$a_{12} a_{21}$	-1

$$a_{11} a_{22} - a_{12} a_{21} = \det |A|$$

$$\Rightarrow \det \begin{pmatrix} A & B \\ C & D \end{pmatrix} = \det(A) \det(D-CA^T B)$$

Let $B=0$:

$$\underbrace{\begin{bmatrix} I & -BD^{-1} \\ O & I \end{bmatrix}}_{\text{block elem. ops. mat}} \begin{bmatrix} A & B \\ C & D \end{bmatrix} = \begin{bmatrix} A-BD^{-1}C & O \\ C & D \end{bmatrix}$$

block elem. ops. mat

$$\Rightarrow \det \begin{pmatrix} A & B \\ C & D \end{pmatrix} = \det(0) \det(A-BD^{-1}C)$$

• $|A|(D-CA^T B) = |D||A-BD^{-1}C|$

HW WEEK 5

① Let A be 3×3 . Fill out a table of all permutations of $(1, 2, 3)$

π	Term	Sign
-------	------	------

Find $\det(A)$.

② A is $p \times p$, D is $n \times n$

$|A| \neq 0$, $|D| \neq 0$. Prove:

$$\det \begin{pmatrix} A & B \\ 0 & D \end{pmatrix} = \det A \det D.$$

③ what is the SWM form for

$|U + VWX|$, where U is $n \times n$, W is $p \times p$,
 V is $n \times p$, X is $p \times n$. & $|U| \neq 0$, $|W| \neq 0$.

_____ X _____

Application of SWM.

Linear mixed models.

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

$$\epsilon \sim \mathcal{N}(0, D), \mu \sim \mathcal{N}(0, D)$$

Unknown: β , μ & maybe para in D & R

Posterior: $\mathcal{N}(y | X\beta + Zu, D) \times$

$$\mathcal{N}(\mu | 0, R) \times$$

$$\mathcal{N}(\beta | 0, V_\beta), V_\beta \rightarrow 0.$$

$$\mathcal{N}(\beta | 0, V_\beta) \rightarrow C. \text{ (improper)}$$

So posterior (or likelihood)

$$\mathcal{N}(u | 0, R) \times \mathcal{N}(y | X\beta + Zu, D)$$

Dim of unknown $\{u \& \beta\}$.

Integrate out the "u".

Stats: marginalize.

Machine Learning: Collapsing.

$$\int \mathcal{N}(u | 0, R) \times \mathcal{N}(y | X\beta + Zu, D) du$$

WEEK 5 Lec10 04-28

Sequential Bayesian Learning

Idea: Single data point $y_1 \sim p(y|\theta)$

$$\theta \sim p(\theta)$$

$$\text{Posterior: } p(\theta|y_1) \propto p(\theta) p(y_1|\theta)$$

Second data point: $y_2 \sim p(y|\theta)$

$$y_2 \perp y_1 | \theta$$

Update:

$$p(\theta|y_1, y_2) \propto p(\theta) p(y_1|\theta) p(y_2|\theta)$$

$$\propto \underbrace{p(\theta) p(y_1|\theta)}_{\propto p(\theta|y_1)} p(y_2|\theta)$$

$$\propto p(\theta|y_1) p(y_2|\theta)$$

$$\begin{aligned} \text{Prior on } \theta &\xrightarrow{\text{Updated}} p(\theta|y_1) \xrightarrow{\text{Updated}} p(\theta|y_1, y_2) \\ &\rightarrow \dots \xrightarrow{\text{Updated}} p(\theta|y_1, y_2, \dots, y_n) \end{aligned}$$

Ex $y_i | \theta \stackrel{iid}{\sim} \mathcal{N}(\theta, \sigma^2)$, σ^2 known
 $i=1, \dots, n$, n : sample size, n_0 : prior sample size
 $\theta \sim \mathcal{N}(\mu, \frac{\sigma^2}{n_0})$

n = Sample size

n_0 = "prior sample size"

(Larger $n_0 \Rightarrow$ more precise priors)

$$\begin{aligned} p(\theta|y_1) &\propto \mathcal{N}(\theta|\mu, \frac{\sigma^2}{n_0}) \times \mathcal{N}(y_1|\theta, \sigma^2) \\ &\propto e^{-\frac{n_0}{2\sigma^2}(\theta-\mu)^2} \times e^{-\frac{1}{2\sigma^2}(y_1-\theta)^2} \\ &\propto e^{-\frac{1}{2\sigma^2}\{n_0(\theta-\mu)^2 + (y_1-\theta)^2\}} \end{aligned}$$

$$\begin{aligned} n_0(\theta-\mu)^2 + (y_1-\theta)^2 &= (n_0+1)\theta^2 - 2(\mu n_0 + y_1)\theta + y_1^2 \\ &= (n_0+1)\left\{\theta - \frac{n_0\mu + y_1}{n_0+1}\right\}^2 + \text{Const.} \end{aligned}$$

$$\therefore p(\theta|y_1) = \mathcal{N}(\theta | \frac{n_0\mu + y_1}{n_0+1}, \frac{\sigma^2}{n_0+1})$$

mean in weighted average:

$$\underbrace{\frac{n_0}{n_0+1}\mu}_{\text{prior weight}} + \underbrace{\frac{1}{n_0+1}y_1}_{\text{data weight}}$$

$$\text{Precision enhanced: } \frac{\sigma^2}{n_0} \rightarrow \frac{\sigma^2}{n_0+1}$$

$$\begin{aligned} p(\theta|y_1, y_2) &\propto p(\theta|y_1) p(y_2|\theta) \\ &\propto e^{-\frac{(n_0+1)}{2\sigma^2}(\theta - \theta_*)^2} \times e^{-\frac{1}{2\sigma^2}(y_2 - \theta)^2} \end{aligned}$$

$$\theta_* = \frac{n_0\mu + y_1}{n_0+1} \Rightarrow$$

$$p(\theta|y_1, y_2) =$$

$$\mathcal{N}(\theta | \frac{(n_0+1)\theta_* + y_2}{(n_0+1)+1}, \frac{\sigma^2}{(n_0+1)+1})$$

$$(n_0+1)\theta_* + y_2 =$$

$$(n_0+1)\left\{\frac{n_0\mu + y_1}{n_0+1}\right\} + y_2 =$$

$$\Downarrow (\Delta) = \frac{n_0\mu + 2y_1 + y_2}{n_0+2}$$

Carry on:

$$p(\theta|y_1, y_2) = \mathcal{N}(\theta | \frac{n_0\mu + 2y_1 + y_2}{n_0+2}, \frac{\sigma^2}{n_0+2})$$

$$\dots \dots$$

$$p(\theta|y_1, \dots, y_{n-1}) = \mathcal{N}(\theta | \frac{n_0\mu + (n-1)y_{n-1}}{n_0+(n-1)}, \frac{\sigma^2}{n_0+n-1})$$

$$\Downarrow \mathcal{N}(\theta | \frac{n_0\mu + (n-1)y_{n-1}}{n_0+(n-1)}, \frac{\sigma^2}{n_0+n-1})$$

$$p(\theta|y_1, \dots, y_n) \propto \mathcal{N}(\theta | \theta_*, \frac{\sigma^2}{n_0+n-1}) \times \mathcal{N}(y_n | \theta, \sigma^2)$$

$$= \mathcal{N}(\theta | \frac{(n_0+n-1)\theta_* + y_n}{(n_0+n-1)+1}, \frac{\sigma^2}{n_0+n-1})$$

$$\begin{aligned}\Theta_* &= \frac{n_0\mu + (n-1)\bar{y}_{n-1}}{n_0+n-1} \Rightarrow \\ (n_0+n-1)\Theta_* &= n_0\mu + (n-1)\bar{y}_{n-1} + y_n \\ \Rightarrow p(\theta|y_1, \dots, y_n) &= \\ \mathcal{N}(\theta | \frac{n_0\mu + \sum_{i=1}^{n-1} y_i + y_n}{n_0+n}, \frac{\sigma^2}{n_0+n}) \\ &= \mathcal{N}(\theta | \frac{n_0\mu + n\bar{y}}{n_0+n}, \frac{\sigma^2}{n_0+n})\end{aligned}$$

• Much simpler to use

Sufficient statistic (2nd AB)!

$$Y_1, \dots, Y_n | \theta \stackrel{\text{iid}}{\sim} \mathcal{N}(\theta, \sigma^2)$$

$$\Rightarrow \bar{y} \sim \mathcal{N}(\theta, \frac{\sigma^2}{n})$$

$$p(\theta|y_1, \dots, y_n) = p(\theta|\bar{y})$$

$$\propto \mathcal{N}(\theta | \mu, \frac{\sigma^2}{n}) \times \mathcal{N}(\bar{y} | \theta, \frac{\sigma^2}{n})$$

$$= \mathcal{N}(\theta | \frac{n_0\mu + n\bar{y}}{n_0+n}, \frac{\sigma^2}{n_0+n})$$

X

Return to HW:

$$\left\{ \begin{array}{l} Y = X\beta + e_y; e_y \sim \mathcal{N}(0, \sigma^2 I_n) \\ \beta = \mu_\beta + e_\beta; e_\beta \sim \mathcal{N}(0, \sigma^2 V_\beta) \\ \tilde{Y} = \tilde{X}\beta + \tilde{e}; \tilde{e} \sim \mathcal{N}(0, \sigma^2 I_m) \end{array} \right.$$

• Assume σ^2 known (since we want to find $p(\tilde{Y} | \sigma^2, y)$)

$$\bullet p(\beta | \tilde{Y} | \sigma^2, y) \propto p(\beta | \sigma^2, y) \times \text{Since } \tilde{Y} \perp Y | \beta, \sigma^2 \quad p(\tilde{Y} | \beta, \sigma^2)$$

$$\tilde{Y} = \tilde{X}\beta + \tilde{e}; \tilde{e} \sim \mathcal{N}(0, \sigma^2 I_m)$$

• β 's distribution has been updated to $\beta | y \sim \mathcal{N}(\mu_m, \sigma^2 M)$

$$-M^{-1} = V_\beta^{-1} + X^T X$$

$$-m = V_\beta^{-1} \mu_\beta + X^T y$$

• Writing as a linear model

$$\beta = \mu_m + e_{\beta | y}, e_{\beta | y} \sim \mathcal{N}(0, \sigma^2 M)$$

To find $p(\tilde{Y} | \sigma^2, y)$ we substitute the LM for β (posterior) into the LM for \tilde{Y} :

$$\tilde{Y} = \tilde{X}(\mu_m + e_{\beta | y}) + \tilde{e}$$

$$= \tilde{X} \mu_m + \tilde{X} e_{\beta | y} + \tilde{e}$$

$$= \tilde{X} \mu_m + e_*$$

$$e_* \sim \mathcal{N}(0, \sigma^2 (I_m + \tilde{X} M \tilde{X}^T))$$

$$\Rightarrow \tilde{Y} | y, \sigma^2 \sim$$

$$\mathcal{N}(\tilde{X} \mu_m, \sigma^2 (I_m + \tilde{X} M \tilde{X}^T))$$

Application of SWM to Linear Mixed Models

$$Y = X\beta + ZU + e; \quad \begin{matrix} n \times 1 \\ n \times p \\ n \times r \\ n \times 1 \end{matrix} \quad \begin{matrix} n \times p \\ n \times r \\ n \times r \\ n \times 1 \end{matrix} \quad \begin{matrix} n \times 1 \end{matrix}$$

$$e \sim N(0, D)$$

$$U \sim N(0, R)$$

Dimension Reduction: Choose

a design matrix $Z_{n \times r}$ &
 $r \ll n$ Not: R is $r \times r$

$$\text{Collapse: } Y \sim N(X\beta, ZRZ^T + D)$$

$$Y|\beta \sim N(X\beta, ZRZ^T + D)$$

Compute likelihood: to solve $\hat{\beta}$ of β

$(\underbrace{ZRZ^T + D}_{n \times n})^{-1}$ is required

What if n is BIG?

- SWM to the rescue:

$$(D + ZRZ^T)^{-1} =$$

$$D^{-1} - D^{-1}Z(R^{-1} + Z^T D^{-1} Z)^{-1} Z^T D^{-1}$$

Advantage: D is diagonal / sparse

s.t. D^{-1} is cheap.

- Z is $n \times r$ with $r \ll n$. Modeler chooses σ^2 & Z . R is $r \times r$, R^{-1} is cheap.

Prediction when $\tilde{Y} \neq Y | \beta, \sigma^2$

$$\begin{bmatrix} Y \\ \tilde{Y} \end{bmatrix} \sim N\left(\begin{bmatrix} X\beta \\ \tilde{X}\beta \end{bmatrix}, \sigma^2 \begin{bmatrix} V_Y & V_{Y\tilde{Y}} \\ V_{\tilde{Y}Y} & V_{\tilde{Y}} \end{bmatrix}\right)$$

- $p(\tilde{Y}|Y, \sigma^2, \beta)$

$$= N(\tilde{X}\beta + V_{\tilde{Y}Y}V_Y^{-1}(Y - X\beta), \sigma^2(V_{\tilde{Y}} - V_{\tilde{Y}Y}V_Y^{-1}V_{Y\tilde{Y}}))$$

- $p(\beta|Y, \sigma^2) = N(\beta|m_m, \sigma^2 M)$

- $p(\sigma^2|Y) = IG(\sigma^2|a^*, b^*)$

- Stat: multivariate Bayesian stat.

- ML: probabilistic machine learning.

WEEK 6 Lec II May 03.2021

Recall: $\mathcal{N}(\beta | \mu_\beta, \sigma^2 V_\beta) \times \mathcal{N}(y | X\beta, \sigma^2 V_y)$

$$\begin{aligned} p(y | \sigma^2) &= \int p(\beta, y | \sigma^2) d\beta \\ &\quad p(\beta | \sigma^2) p(y | \beta, \sigma^2) \\ &= \int \mathcal{N}(\beta | \mu_\beta, \sigma^2 V_\beta) \times \mathcal{N}(y | X\beta, \sigma^2 V_y) d\beta \end{aligned}$$

Easy to avoid integration:

$$[y | \beta, \sigma^2] \rightarrow y = X\beta + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2 V_y)$$

$$[\beta | \sigma^2] \rightarrow \beta = \mu_\beta + \epsilon_\beta; \quad \epsilon_\beta \sim \mathcal{N}(0, \sigma^2 V_\beta)$$

$$\begin{aligned} \Rightarrow [y | \sigma^2] &\rightarrow y = X(\mu_\beta + \epsilon_\beta) + \epsilon \\ &= X\mu_\beta + \underbrace{X\epsilon_\beta + \epsilon}_{\text{[75-year-old prof.] } \mathcal{N}(0, \sigma^2(XV_\beta X^T + V_y))} \end{aligned}$$

Similar for predictions:

$$p(\tilde{y} | y, \sigma^2) = \int p(\tilde{y} | \beta, \sigma^2, y) \cdot p(\beta | \sigma^2, y) d\beta$$

$$[\tilde{y} | \beta, \sigma^2, y] \rightarrow \tilde{y} = \tilde{X}\beta + \tilde{\epsilon}; \quad \tilde{\epsilon} \sim \mathcal{N}(0, \sigma^2 \tilde{V})$$

$$[\beta | \sigma^2, y] \rightarrow \beta = Mm + \epsilon_\beta; \quad \epsilon_\beta \sim \mathcal{N}(0, \sigma^2 M)$$

$$\begin{aligned} [\tilde{y} | \sigma^2, y] &\rightarrow \tilde{y} = \tilde{X}(Mm + \epsilon_\beta) + \tilde{\epsilon} \\ &= \tilde{X}Mm + \underbrace{\tilde{X}\epsilon_\beta + \tilde{\epsilon}}_{\sim (0, \sigma^2 (\tilde{X}M\tilde{X}^T + \tilde{V}))} \end{aligned}$$

Q Why $\epsilon \perp \epsilon_\beta$?

Key point: Look @ the integration.

Always remember:

$$\begin{aligned} [\text{Unknown} | \text{known}] &\propto [\text{unknown, known}] \\ [\beta, \sigma^2, \tilde{Y} | Y] &\propto [\beta, \sigma^2, \tilde{Y}, Y] \\ &\propto [\sigma^2 | \beta, \sigma^2] [\tilde{Y} | \beta, \sigma^2] \times \\ &\quad [Y | \beta, \sigma^2] \propto [\beta, \sigma^2 | Y] \\ &\propto [\beta, \sigma^2 | Y] [\tilde{Y} | \beta, \sigma^2] \end{aligned}$$

Prediction vs. Model-fitting

$$y_i | \beta, \sigma^2 \stackrel{\text{ind}}{\sim} \mathcal{N}(x_i^T \beta, \sigma^2)$$

$$\beta | \sigma^2 \sim \mathcal{N}(\mu_\beta, \sigma^2 V_\beta)$$

$$\sigma^2 \sim \text{IG}(a, b)$$

$$(\tilde{Y}, \tilde{X}) \rightarrow \tilde{Y} | \beta, \sigma^2 \sim \mathcal{N}(\tilde{X}^T \beta, \sigma^2)$$

$$\begin{aligned} \Rightarrow y &= X\beta + \epsilon; \quad \beta = \mu_\beta + \epsilon_\beta; \\ \tilde{Y} &= \tilde{X}\beta + \tilde{\epsilon} \end{aligned}$$

$$\tilde{Y} | y, \sigma^2 \sim \mathcal{N}(\tilde{X}Mm, \sigma^2 (I_m + \tilde{X}M\tilde{X}^T))$$

Prediction: New (but given) \tilde{X}

Predict \tilde{Y} for that new \tilde{X}

Model-fitting: \tilde{Y} but with X , not \tilde{X} !

\Rightarrow replicate

$$\tilde{Y} | y, \sigma^2 \sim \mathcal{N}(XMm, \sigma^2 (I_m + XMX^T))$$

Beneath the above derivation, we still assume

$$\tilde{y} \perp\!\!\!\perp y \mid \beta, \sigma^2$$

$P(\tilde{y} \mid \sigma^2, y) \leftarrow$ model's ability to replicate the data y !

Ex $V_{\beta}^{-1} = 0, a = -\frac{p}{2}, b = 0$

$$P(\beta, \sigma^2) \propto \frac{1}{\sigma^2} \left(\frac{1}{\sigma^2}\right)^{a+\frac{p}{2}+1}$$

Then $\begin{cases} M = (X^T X)^{-1} \\ m = X^T y \end{cases} \rightarrow Mm = \hat{\beta}_{OLS}$
 $= (X^T X)^{-1} X^T y$

model-fitting

$$\tilde{Y} \mid y, \sigma^2 \sim \mathcal{N}(X \hat{\beta}_{OLS}, \sigma^2(I_m + P_X))$$

Prediction

$$\tilde{Y} \mid y, \sigma^2 \sim \mathcal{N}(\tilde{X} \hat{\beta}_{OLS}, \sigma^2(I_m + \tilde{X}(X^T X)^{-1} \tilde{X}^T))$$

To compare models:

most model comparison matrices use:

$$P(\tilde{Y} \mid y) \rightarrow \text{pointwise}$$

We want to investigate behavior or performance in terms of

$$P(\tilde{Y}_i \mid y) \text{ for } i=1, 2, \dots, n$$

Ex $\tilde{Y} \mid y, \sigma^2 \sim \mathcal{N}(X \hat{\beta}_{OLS}, \sigma^2(I + P_X))$

$$\tilde{Y}_i \mid y, \sigma^2 \sim \mathcal{N}(X_i^T \hat{\beta}_{OLS}, \sigma^2(I + P_{X,i}))$$

$$P_X = X(X^T X)^{-1} X^T; [P_X]_{ii} = X_i^T (X^T X)^{-1} X_i$$

where

$X_i = \text{row } i \text{ of } X$ or

col i of X^T

Assume replicates are indep across i :

$$Y_i \mid \sigma^2, y \sim \mathcal{N}(X_i^T \hat{\beta}_{OLS}, \sigma^2(1 + X_i^T (X^T X)^{-1} X_i))$$

Common model comparison metric:

$$\sum_{i=1}^n \log P(\tilde{Y}_i \mid y) \xrightarrow{\text{widely}} \text{WAIC}$$

(after penalty) $\xrightarrow{\text{leads to}}$

General:

$$P(\tilde{Y} \mid y) = \int P(\tilde{Y}_i \mid \theta) p(\theta \mid y) d\theta$$

$$= \mathbb{E}_{\theta \mid y} [P(\tilde{Y}_i \mid \theta)]$$

$$= \int P(\tilde{Y}_i \mid \theta) \underbrace{P_{\text{post}}(\theta)}_{P(\theta \mid y)} d\theta$$

$$P_{\text{post}}(\tilde{Y}_i) = P(\tilde{Y}_i \mid y)$$

Suppose we have generated samples

$$\tilde{Y}_i^{(1)}, \dots, \tilde{Y}_i^{(n)} \sim P(\tilde{Y}_i \mid y)$$

$$\tilde{\mu}_i = \mathbb{E}(\tilde{Y}_i \mid y) \approx \frac{1}{n} \sum_{j=1}^n \tilde{Y}_i^{(j)}$$

Goodness of fit measure:

$$G = \sum_{i=1}^n (Y_i - \tilde{\mu}_i)^2$$

Also need to penalize models.

$$P = \sum_{i=1}^n \text{Var}(\tilde{Y}_i | Y),$$

$$\text{Var}(\tilde{Y}_i | Y) = \frac{1}{N} \sum_{j=1}^N (\tilde{Y}_i^{(j)} - \tilde{\mu}_i)^2$$

model comparison score:

$$D = \underbrace{G}_{\text{goodness of fit}} + \underbrace{P}_{\text{penalty}}$$



Marginalize σ^2

$$\tilde{Y} | Y, \sigma^2 \sim \mathcal{N}(\tilde{X}Mm, \sigma^2(I_m + \tilde{X}M\tilde{X}^T))$$

$$[\tilde{Y} | Y] = \int [\tilde{Y} | Y, \sigma^2] [\sigma^2 | Y]$$

$$= \int \mathcal{N}(\tilde{Y} | \tilde{X}Mm, \sigma^2(I_m + \tilde{X}M\tilde{X}^T)) \times I_G(\sigma^2 | a_*, b_*) d\sigma^2$$

$$= \int \text{NIIG}(\tilde{Y}, \sigma^2 | \underset{\mu}{\tilde{X}Mm}, \underset{\Sigma}{I_m + \tilde{X}M\tilde{X}^T}, \underset{a}{a_*}, \underset{b}{b_*}) d\sigma^2$$



Multivariate (non-central)
Student's t distribution



Another look @ normal densities & Linear Models.

$$y \sim \mathcal{N}(\mu, \Sigma)$$

$n \times 1$

$$y = (Y_1, Y_2, \dots, Y_n)^T$$

$$\bullet P(Y) = P(Y_1) P(Y_2 | Y_1) P(Y_3 | Y_1, Y_2) \dots P(Y_n | Y_1, \dots, Y_{n-1})$$

$$[Y_1] \rightarrow Y_1 = \mu_1 + \epsilon_1; \epsilon_1 \sim \mathcal{N}(0, d_1)$$

$$[Y_2 | Y_1] \rightarrow Y_2 - \mu_2 = a_{21}(Y_1 - \mu_1) + \epsilon_2; \epsilon_2 \sim \mathcal{N}(0, d_2)$$

$$[Y_3 | Y_2, Y_1] \rightarrow Y_3 - \mu_3 = a_{31}(Y_1 - \mu_1) + a_{32}(Y_2 - \mu_2) + \epsilon_3; \epsilon_3 \sim \mathcal{N}(0, d_3)$$

⋮

$$[Y_i | Y_{i-1}, \dots, Y_1] \rightarrow$$

$$Y_i - \mu_i = \sum_{j=1}^{i-1} a_{ij}(Y_j - \mu_j) + \epsilon_i; \epsilon_i \sim \mathcal{N}(0, d_i)$$

$i = 2, \dots, n$

⇒

$$\begin{bmatrix} Y_1 - \mu_1 \\ Y_2 - \mu_2 \\ \vdots \\ Y_n - \mu_n \end{bmatrix}$$

≡

$$\begin{bmatrix} 0 & 0 & \dots & 0 & 0 \\ a_{21} & 0 & \dots & 0 & 0 \\ a_{31} & \vdots & \ddots & \vdots & \vdots \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{(n-1)1} & \dots & a_{(n-1)n-1} & 0 & 0 \\ a_{n1} & a_{n2} & \dots & a_{nn-1} & 0 \end{bmatrix} \begin{bmatrix} Y_1 - \mu_1 \\ Y_2 - \mu_2 \\ \vdots \\ Y_n - \mu_n \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$\bullet Y - \mu = A(Y - \mu) + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, D), D = \text{Diag}(d_i)$$

$$\Rightarrow (I - A)(Y - \mu) = \epsilon$$

$$(I - A)x = 0 \Rightarrow$$

$$\begin{bmatrix} 1 & 0 & \dots & 0 \\ a_{21} & 1 & \dots & 0 \\ a_{31} & a_{22} & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{(n-1)1} & a_{(n-1)2} & \dots & a_{(n-1)n-1} & 1 \\ a_{n1} & a_{n2} & \dots & a_{nn-1} & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \Rightarrow x_1 = x_2 = \dots = x_n = 0!$$

⇒ $I - A$ is nonsingular!

$\Rightarrow (I-A)^{-1}$ exists!

So $(I-A)(Y\mu) = \epsilon$;

$$(I-A)\text{Var}(Y\mu)(I-A)^T = \text{Var}(\epsilon)$$

or

$$(I-A)\sum(I-A)^T = D$$

$$\Rightarrow \Sigma = \underbrace{(I-A)^{-1}D(I-A)^{-T}}_{\text{Cholesky}}$$

• Notation:

$$(I-A)^{-T} = \left[(I-A)^{-1} \right]^T = \left[(I-A)^T \right]^{-1} \\ = \left[(I-A^T) \right]^{-1}$$

upper triangular

$$\Sigma^{-1} = (I-A)^T D^{-1} (I-A)$$

$$\text{Cholesky: } \Sigma = L D L^T$$

L is unit lower triangular, D has diagonal element.

$$D = D^{\frac{1}{2}} D^{\frac{1}{2}}$$

$$L D L^T = L D^{\frac{1}{2}} D^{\frac{1}{2}} L^T$$

$$= \underbrace{(LD^{\frac{1}{2}})}_{L} \underbrace{(LD^{\frac{1}{2}})^T}_{L^T}$$

We have derived the Cholesky decom.
using a probability factorization.

Let K be a p.d. matrix.

Spectral decomp. $K = P \Lambda P^T$

$$\cdot z \sim N(0, I)$$

$$\cdot y = \mu + P \Lambda^{\frac{1}{2}} z$$

$$\Rightarrow y \sim N(\mu, K)$$

so if p.d. K leads to

$$y \sim N(\mu, K)$$

So every p.d. K has a Cholesky decomp.

WEEK 7 LEC 13

May 10. 2021

Recap of Midterm

1) Brook's Lemma

2) Bayesian LR

$$y = (y_1^T, y_2^T, \dots, y_k^T)^T$$

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_k \end{bmatrix} \quad D_k = \{y_k, X_k\}$$

- $p(\beta, \sigma^2 | D_1, D_2, \dots, D_k)$

Step I:

$$\propto p(\beta, \sigma^2 | D_1) \propto p(\beta, \sigma^2) \times p(D_1 | \beta, \sigma^2)$$

$$\propto NIG(\beta, \sigma^2 | M_\beta, V_\beta, a, b) \times$$

$$\mathcal{N}(y_1 | X_1 \beta_1, \sigma^2 I_{m_1})$$

$$\propto NIG(\beta, \sigma^2 | M_1, m_1, M_1, a_1, b_1)$$

where

$$\begin{cases} m_1 = V_\beta^\top \mu_\beta + X_1^\top y_1 \\ M_1^{-1} = V_\beta^{-1} + X_1^\top X_1 \end{cases} \quad \begin{cases} a_1 = a + \frac{m_1}{2} \\ b_1 = b + \frac{1}{2} \{ y_1^\top y_1 + \mu_\beta^\top V_\beta^{-1} \mu_\beta \\ - m_1^\top M_1 m_1 \} \end{cases}$$

Step II :

$$\propto p(\beta, \sigma^2 | D_1, D_2) \propto p(\beta, \sigma^2) p(D_1, D_2 | \beta, \sigma^2)$$

$$\propto \underbrace{p(\beta, \sigma^2 | D_1)}_{p(\beta, \sigma^2 | D_1)} p(D_2 | \beta, \sigma^2)$$

$$\propto NIG(\beta, \sigma^2 | M_1, m_1, M_1, a_1, b_1) \\ \propto \mathcal{N}(y_2 | X_2 \beta, \sigma^2 I_{m_2})$$

L

$$NIG(\beta, \sigma^2 | M_2, m_2, M_2, a_2, b_2)$$

where

$$\begin{cases} m_2 = m_1 + X_2^\top y_2 \\ M_2^{-1} = M_1^{-1} + X_2^\top X_2 \end{cases} \quad \begin{cases} a_2 = a_1 + \frac{m_1}{2} \\ b_2 = b_1 + \frac{1}{2} \{ y_2^\top y_2 + m_1^\top M_1 m_1 \\ - m_2^\top M_2 m_2 \} \end{cases}$$

We can go for K steps:

- $M_K = M_{K-1} + X_K^\top Y_K$
 $= V_\beta^{-1} \mu_\beta + \sum_{i=1}^{K-1} X_i^\top Y_i + X_K^\top Y_K$
 $= V_\beta^{-1} \mu_\beta + \sum_{i=1}^K X_i^\top Y_i$

- $M_K^{-1} = M_{K-1}^{-1} + X_K^\top X_K$
 $= V_\beta^{-1} + \sum_{i=1}^{K-1} X_i^\top X_i + X_K^\top X_K$

- $\alpha_* = \alpha_K = a + \frac{1}{2} \sum_{k=1}^K m_k = a + \frac{n}{2}$
- $b_* = b_K = b + \frac{1}{2} \{ \sum_{k=1}^K y_i^\top y_i + \mu_\beta^\top V_\beta^{-1} \mu_\beta \\ - m_K^\top M_K^{-1} m_K \}$

So two approaches give the same result.

$$\int p(\theta|x_1)d\theta = \int \frac{p(\theta)f(x_1|\theta)}{\left[\int p(\theta)f(x_1|\theta)d\theta \right]} d\theta = 1$$

Homework week 7

3) Augmented linear model

4) Density function we have:

$$X = (x_1, \dots, x_n)^T$$

$$p(x) = \prod_{i=1}^n f(x_i|\theta)$$

$$x_i|\theta \stackrel{iid}{\sim} f(\cdot|\theta)$$

$$\int p(\theta)d\theta = +\infty \text{ and}$$

$$\int p(\theta)f(x_1|\theta)dx_1 < +\infty$$

What can we say about

$$\int p(\theta) \prod_{i=1}^n f(x_i|\theta) d\theta = ?$$

$$p(\theta)f(x_1|\theta) \propto p(\theta|x_1)$$

$$\int p(\theta|x_1)d\theta = 1 -$$

$$\text{NOTE } p(\theta|x_1) = C \cdot p(\theta) \cdot f(x_1|\theta)$$

If the prior \times likelihood integrates to finite number, the posterior is well-defined.

$$\text{Also if } \int p(\theta|x)d\theta = C \Rightarrow$$

$$\frac{1}{C} \int p(\theta|x)d\theta = 1$$

So $\frac{1}{C} p(\theta|x)$ is a proper prior.

Look @ $p(\theta|x_1, x_2)$

$$\begin{aligned} p(\theta|x_1, x_2) &\propto p(\theta) \cdot f(x_1, x_2|\theta) \\ &\propto \underbrace{p(\theta)f(x_1|\theta)}_{\propto p(\theta|x_1)} f(x_2|\theta) \end{aligned}$$

$$\propto \underbrace{p(\theta|x_1)f(x_2|\theta)}_{\text{proper prior proper density}}$$

$\Rightarrow p(\theta|x_1, x_2)$ must be well-defined.

$\Rightarrow \int p(\theta)f(x_1, x_2|\theta)d\theta$ is well-defined

□

WEEK 7 LEC 14 · May 12

Recall joint prob. model:

$$p(Y_1) p(Y_2|Y_1) p(Y_3|Y_1, Y_2) \dots p(Y_n|Y_1, \dots, Y_{n-1})$$

We saw the linear model representation

$$Y_i = \mu_i + \epsilon; \epsilon \sim N(0, d_1)$$

$$Y_i = \sum_{j=1}^{i-1} a_{ij} Y_j + \epsilon_i; \epsilon_i \sim N(0, d_2) \quad i=2, \dots, n.$$

$$\Rightarrow y = Ay + \epsilon,$$

$A = \{a_{ij}\}$ strictly lower- Δ

$$\Rightarrow \underbrace{(I-A)y}_{\text{unit lower-}\Delta} = \epsilon$$

unit lower- Δ

$$y \sim N(0, (I-A)^{-1} D (I-A)^{-1})$$

If $y \sim N(0, \Sigma)$, then

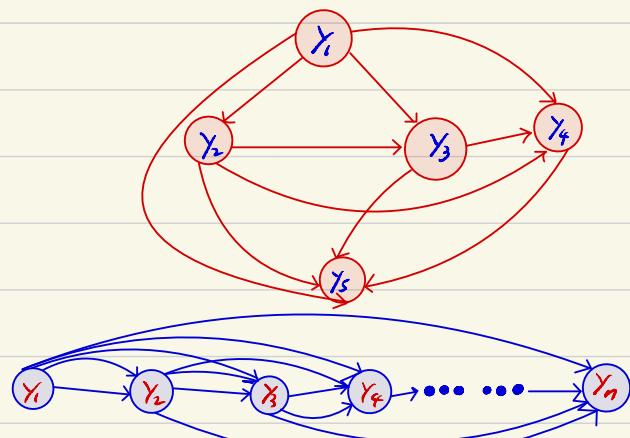
$$\Sigma = \underbrace{(I-A)^{-1}}_{\text{P.d.}} \underbrace{D}_{L} \underbrace{(I-A)^{-1}}_{L^T}$$

Cholesky Decomp.

DAGS

$$p(Y_1) p(Y_2|Y_1) \dots p(Y_n|Y_1, \dots, Y_{n-1})$$

Directed Acyclic Graphical Model



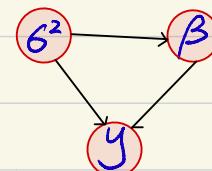
DAGs: Bayesian Networks

All Bayesian models are DAGs

- DAGs because they (~~at the core~~) specify joint distributions.

$$\begin{aligned} \text{Ex } Y_i | \beta, \sigma^2 &\sim N(x_i^\top \beta, \sigma^2) \\ \beta | \sigma^2 &\sim N(\mu_\beta, \sigma^2 V_\beta) \\ \sigma^2 &\sim IG(a, b) \end{aligned}$$

How to represent it as a DAG?



Suppose:

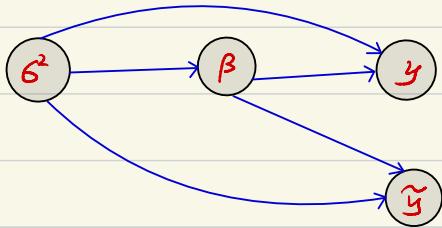
$$y | \beta, \sigma^2 \sim N(x\beta, \sigma^2 I_n)$$

$$\beta | \sigma^2 \sim N(\mu_\beta, \sigma^2 V_\beta)$$

$$\sigma^2 \sim IG(a, b)$$

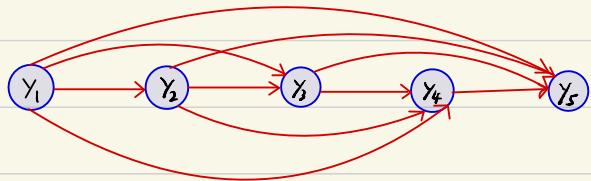
$$\tilde{y} | \beta, \sigma^2 \sim N(\tilde{x}\beta, \sigma^2 I_n)$$

What is the DAG?



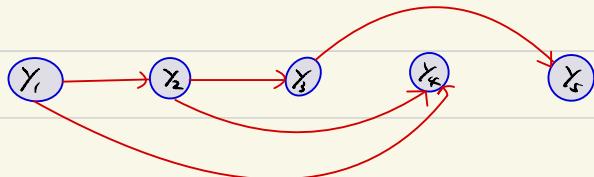
DAGs are easily modified to incorporate conditional independence (and hence sparsity).

Consider $n=5$, y_1, y_2, \dots, y_5 , complete DAG:



$$P(Y_1) P(Y_2|Y_1) P(Y_3|Y_1, Y_2) P(Y_4|Y_1, Y_2, Y_3) P(Y_5|Y_1, Y_2, Y_3)$$

Drop edges:



$$\begin{aligned} & P(Y_1) P(Y_2|Y_1) P(Y_3|Y_1) P(Y_4|Y_1, Y_2) P(Y_5|Y_1, Y_2) \\ \Rightarrow \quad & Y_3 \perp\!\!\!\perp Y_1 | Y_2 ; \quad Y_4 \perp\!\!\!\perp Y_3 | Y_1, Y_2 \\ & Y_5 \perp\!\!\!\perp (Y_1, Y_2, Y_4) | Y_3 \end{aligned}$$

Conditional indep \rightarrow Sparsity

Return to Linear model representation

Suppose $\mu=0$,

$$y_1 = 0 + \epsilon_1, \quad \epsilon_1 \sim \mathcal{N}(0, d_1)$$

$$y_2 = \alpha_{21} y_1 + \epsilon_2, \quad \epsilon_2 \sim \mathcal{N}(0, d_2)$$

$$y_3 = \alpha_{31} y_1 + \alpha_{32} y_2 + \epsilon_3, \quad \epsilon_3 \sim \mathcal{N}(0, d_3)$$

$$Y_3 \perp\!\!\!\perp Y_1 | Y_2 \Rightarrow \alpha_{31} = 0$$

$$y_4 = \alpha_{41} y_1 + \alpha_{42} y_2 + \alpha_{43} y_3 + \epsilon_4; \quad \epsilon_4 \sim \mathcal{N}(0, d_4)$$

$$Y_4 \perp\!\!\!\perp Y_3 | Y_1, Y_2 \Rightarrow \alpha_{43} = 0$$

$$y_5 = \alpha_{51} y_1 + \alpha_{52} y_2 + \dots + \alpha_{54} y_4 + \epsilon_5; \quad \epsilon_5 \sim \mathcal{N}(0, d_5)$$

$$Y_5 \perp\!\!\!\perp Y_1, Y_2, Y_4 | Y_3 \Rightarrow$$

$$\alpha_{51} = \alpha_{52} = \alpha_{54} = 0$$

$$P(Y_5|Y_1, Y_2, Y_3, Y_4) = P(Y_5|Y_3)$$

$$P(Y_4|Y_1, Y_2, Y_3) = P(Y_4|Y_1, Y_2)$$

$$P(Y_3|Y_1, Y_2) = P(Y_4|Y_3)$$

Parents shrinking $\Leftrightarrow \alpha_{ij}'s = 0$

$$\Sigma = (I - A)^{-1} D (I - A^T)^{-1}$$

$$\Sigma^{-1} = (I - A^T) D^{-1} (I - A)$$

Sparse $A \mapsto$ sparsity in Σ^{-1} .

$$y^T \Sigma^{-1} y = \underbrace{y^T (I - A^T)}_{z^T} \underbrace{D^{-1}}_{z} \underbrace{(I - A)}_{z} y$$

$$z = (I - A) y.$$

Dense $y \xrightarrow{?} \text{sparse } \underline{z}$.

Statistical interpretation comes from

$$\sum^{-1} \longleftrightarrow A$$

precision matrix

General DAG model

$$\prod_{i=1}^n P(y_i | y_{\text{Pa}[i]})$$

$\text{Pa}[1] = \{\emptyset\}$. empty set so:

$$P(y_1 | y_{\text{Pa}[0]}) = P(y_1)$$



Full conditionals from a DAG

$$\textcircled{E} \quad \tilde{P}(y_1, \dots, y_n) = P(y_1)P(y_2 | y_1)P(y_3 | y_2) \\ P(y_4 | y_1, y_2)P(y_5 | y_3)$$

- $\tilde{P}(y_1 | \cdot) = \tilde{P}(y_1 | y_2, y_3, y_4, \dots)$

$$\propto P(y_1)P(y_2 | y_1)P(y_4 | y_1, y_2)$$

$$\propto P(y_1 | y_{\text{Pa}[1]}) \prod_{j: y_j \in \text{ch}(y_1)} P(y_j | y_{\text{Pa}[j]})$$

$$\tilde{P}(y_3 | \cdot) \propto P(y_3 | y_{\text{Pa}[3]}) \cdot \underbrace{P(y_4 | y_3)}_{\prod_{j: y_j \in \text{ch}(y_3)} P(y_j | y_{\text{Pa}[j]})}$$

In general,

$$\tilde{P}(y_i | \cdot) = P(y_i | \text{Pa}[i]) \prod_{j: y_j \in \text{ch}(y_i)} P(y_j | y_{\text{Pa}[j]})$$

Suppose $y = (y_1, \dots, y_n)^T \sim \mathcal{N}(0, Q^{-1})$

$Q = \Sigma^{-1}$ is the precision mat.

$$y_i \perp\!\!\!\perp y_j \mid \underbrace{y_{-c(i,j)}}_{\text{everything except } i \& j.}$$

$P(y_i | y_{-i})$ does not depend on y_j .

i.e. $P(y_i | y_{-i}) = P(y_i | y_{-c(i,j)})$? when?

Can we get a clue in terms of Q ?

Keep an eye here!

$$P(y) \propto e^{-\frac{1}{2} y^T Q y}$$

$$y^T Q y = \sum_{k=1}^n \sum_{l=1}^n q_{kk} y_k y_l$$

$$= \sum_{k=1}^n q_{kk} y_k^2 + \sum_{k \neq l} q_{kk} y_k y_l$$

$P(y_i | \cdot)$ for a fixed i :

need only terms involving i :

$$y^T Q y \Rightarrow q_{ii} y_i^2 + y_i \left(\sum_{l \neq i} q_{il} y_l \right)$$

$$= q_{ii} y_i^2 + y_i (2q_{ij} y_j + \sum_{l \neq i, j} q_{il} y_l)$$

$$= q_{ii} y_i^2 + 2q_{ij} y_i y_j + 2y_i f(y_{-c(i,j)})$$

WEEK 8 • LEC 15 • May 17

\Rightarrow

$$y^T Q y \Rightarrow a_{ii} y_i^2 + 2a_{ij} y_i y_j + 2y_i (\text{terms free of } y_j)$$

So $p(y_i | Y_{-i})$ will be free of y_j if & only if $Q_{ij} = 0$

$$\text{So } p(y_i | Y_{-i}) = p(Y_i | Y_{-i})$$

iff $Q_{ij} = 0$

$$Y_i \perp\!\!\!\perp Y_j | Y_{-i} \Leftrightarrow Q_{ij} = 0$$

Recall precision matrix

$$Y_i \perp\!\!\!\perp Y_j | Y_{-i,j} \text{ iff } Q_{ij} = 0$$

\times

Closer look @ multivariate normal with a precision matrix Q

$$X = \begin{matrix} n_A \\ n_B \end{matrix} \begin{bmatrix} x_A \\ x_B \end{bmatrix} \sim N \left(\begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix}, \underbrace{\begin{bmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{bmatrix}}_{\Sigma} \right)$$

$n_A + n_B = n$

$$\mu_A \quad n_A \times 1, \quad \mu_B \quad n_B \times 1, \quad \Sigma_{AA} \quad n_A \times n_A$$

$$\Sigma_{AB} \quad n_A \times n_B, \quad \Sigma_{BA} \quad n_B \times n_A, \quad \Sigma_{BB} \quad n_B \times n_B$$

$$\Sigma_{AB} = \Sigma_{BA}^T$$

$$Q = \Sigma^{-1} = \begin{bmatrix} Q_{AA} & Q_{AB} \\ Q_{BA} & Q_{BB} \end{bmatrix}$$

$$p(x) \propto e^{-\frac{1}{2}(x-\mu)^T Q (x-\mu)}$$

Closer look @ $(x-\mu)^T Q (x-\mu)$

$$[(x_A - \mu_A)^T : (x_B - \mu_B)^T] \begin{bmatrix} Q_{AA} & Q_{AB} \\ Q_{BA} & Q_{BB} \end{bmatrix} \begin{bmatrix} x_A - \mu_A \\ x_B - \mu_B \end{bmatrix}$$

$$= (x_A - \mu_A)^T Q_{AA} (x_A - \mu_A) + (x_B - \mu_B)^T Q_{BB} (x_B - \mu_B)$$

$$+ (x_A - \mu_A)^T Q_{AB} (x_B - \mu_B) + (x_B - \mu_B)^T Q_{BA} (x_A - \mu_A)$$

Scalars: transpose of each other
 \Rightarrow same or equal.

$$= (x_A - \mu_A)^T Q_{AA} (x_A - \mu_A) + (x_B - \mu_B)^T Q_{BB} (x_B - \mu_B)$$

$$+ 2 (x_A - \mu_A)^T Q_{AB} (x_B - \mu_B)$$

Mn-formula

$$\text{We know } p(x) = p(x_B)p(x_A|x_B)$$

Goal: identify $p(x_B)$ and $p(x_A|x_B)$.

Factorize $p(x) \propto f(x_B) g(x_A|x_B)$

s.t.

$$\int f(x_B) dx_B = C_1 < \infty$$

$$\int g(x_A|x_B) dx_A = C_2 < \infty$$

$$\Rightarrow f(x_B) \propto p(x_B)$$

$$g(x_A|x_B) \propto p(x_A|x_B)$$

$$\text{Pf: } p(x) = p(x_B)p(x_A|x_B) = \tilde{p}(x_B)\tilde{p}(x_A|x_B)$$

$$\Rightarrow p(x_B) = \tilde{p}(x_B) \text{ Because}$$

$$\begin{aligned} p(x_B) &= \int p(x_B)p(x_A|x_B) dx_A = \int \tilde{p}_1(x_B)\tilde{p}_2(x_A|x_B) dx_A \\ &= \tilde{p}_1(x_B) \underbrace{\int \tilde{p}_2(x_A|x_B) dx_A}_{1} \\ &= \tilde{p}_1(x_B) \end{aligned}$$

$$\text{Hence, } p(x_B)p(x_A|x_B) = \tilde{p}_1(x_B)\tilde{p}_2(x_A|x_B)$$

$$\Rightarrow p(x_A|x_B) = \tilde{p}_2(x_A|x_B) \wedge$$

$$x_B \in \{x: p(x) > 0\} \quad \square.$$

Return to our case:

$$p(x) \propto e^{-\frac{1}{2}x^T Q x}$$

$$x^T Q x = (x_A - \mu_A)^T Q_{AA} (x_A - \mu_A) + (x_B - \mu_B)^T Q_{BB} (x_B - \mu_B)$$

$$+ 2(x_A - \mu_A)^T Q_{AB} (x_B - \mu_B)$$

$$= x_A^T Q_{AA} x_A - 2x_A^T Q_{AB} \mu_A + 2x_A^T Q_{AB} (x_B - \mu_B)$$

$$- 2x_B^T Q_{BA} \mu_A + x_B^T Q_{BB} x_B - 2x_B^T Q_{BB} \mu_B$$

+ const.

$$= G(x_A, x_B) + F(x_B) + \text{const.}$$

Where

$$\begin{aligned} G(x_A, x_B) &= x_A^T Q_{AA} x_A - 2x_A^T Q_{AB} \mu_A \\ &\quad + 2x_A^T Q_{AB} (x_B - \mu_B) \end{aligned}$$

$$F(x_B) = x_B^T Q_{BB} x_B - 2x_B^T (Q_{BA} \mu_A + Q_{BB} \mu_B)$$

$$\therefore p(x) \propto e^{-\frac{1}{2}G(x_A, x_B) - \frac{1}{2}F(x_B)}$$

↑
luxury

$$G(x_A, x_B) = x_A^T M^{-1} x_A - 2x_A^T m$$

$$\begin{cases} M^{-1} = Q_{AA} \\ m = Q_{AA} \mu_A - Q_{AB} (x_B - \mu_B) \end{cases},$$

\Rightarrow

$$G(x_A, x_B) = (x_A - Mm)^T M^{-1} (x_A - Mm)$$

$$p(x_A|x_B) = N(x_A | Mm, M)$$

$$\text{Var}(x_A|x_B) = Q_{AA}^{-1} = M$$

$$\mathbb{E}(x_A|x_B) = Mm = Q_{AA}^{-1} (Q_{AA} \mu_A - Q_{AB} (x_B - \mu_B))$$

$$= \mu_A + (-Q_{AA}^{-1} Q_{AB} (x_B - \mu_B))$$

$$G(x_A, x_B) + F(x_B)$$

$$\begin{aligned} &= (x_A - Mm)^T M^{-1} (x_A - Mm) - m^T M m \\ &\quad + F(x_B) \end{aligned}$$

HW: Show $p(x_B) = \mathcal{N}(\mu_B, \star)$

$$\star = \sum_{BB}^{-1}$$

$$x_B = \mu_B + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \Sigma_{BB})$$

$$x_A = \mu_A + (-Q_{AA}^{-1} Q_{AB})(x_B - \mu_B)$$

Also from standard normal:

$$\mathbb{E}[x_A | x_B] = \mu_A + \sum_{AB} \sum_{BB}^{-1} (x_B - \mu_B)$$

$$\sum_{AB} \sum_{BB}^{-1} = -Q_{AA}^{-1} Q_{AB}$$

Recall $\sum^{-1} = Q$, HW:

without explicitly solving a linear system,

Can you use the factorization

$$p(x_A, x_B) = p(x_B) p(x_A | x_B)$$

to express:

$$Q_{AA}^{-1} Q_{AB}^{-1} Q_{BB}^{-1}$$

$$\sum_{AA}, \sum_{AB}, \sum_{BB}$$

• $x_A \perp\!\!\!\perp x_B$ iff $Q_{AB} = 0$

$$p(x_A | x_B) = p(x_A)$$

$$\begin{aligned} m^T M m &= (Q_{AA}\mu_A - Q_{AB}(x_B - \mu_B))^T M m \\ &= (Q_{AA}\mu_A - Q_{AB}(x_B - \mu_B))^T \\ &\quad (\mu_A - Q_{AA}^{-1} Q_{AB}(x_B - \mu_B)) \\ &= \mu_A^T Q_{AA} \mu_A - (x_B - \mu_B)^T Q_{BA} \mu_A - \\ &\quad \mu_A Q_{AB} (x_B - \mu_B) + (x_B - \mu_B)^T Q_{BA} Q_{AA}^{-1} Q_{AB} (x_B - \mu_B) \\ &= (x_B - \mu_B)^T Q_{BA} Q_{AA}^{-1} Q_{AB} (x_B - \mu_B) - \\ &\quad 2(x_B - \mu_B)^T Q_{BA} \mu_A + \text{const.} \end{aligned}$$

$$F(x_B) - m^T M m$$

$$\begin{aligned} &= (x_B - \mu_B)^T (Q_{BB} - Q_{BA} Q_{AA}^{-1} Q_{AB})(x_B - \mu_B) - \\ &\quad 2(x_B - \mu_B)^T (Q_{BA} \mu_A - Q_{BA} \mu_A) \\ &= (x_B - \mu_B)^T [Q_{BB} - Q_{BA} Q_{AA}^{-1} Q_{AB}] (x_B - \mu_B) \end{aligned}$$

$$\Rightarrow p(x_B) \propto e^{-\frac{1}{2} \{ F(x_B) - m^T M m \}} \\ = \mathcal{N}(x_B | \mu_B [Q_{BB} - Q_{BA} Q_{AA}^{-1} Q_{AB}])$$

$$\text{So } \sum_{BB}^{-1} = Q_{BB} - Q_{BA} Q_{AA}^{-1} Q_{AB}.$$

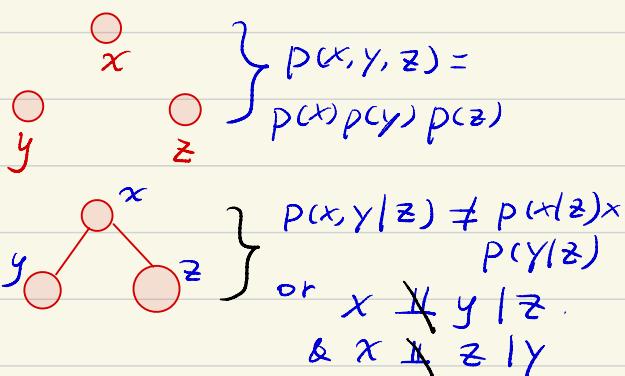
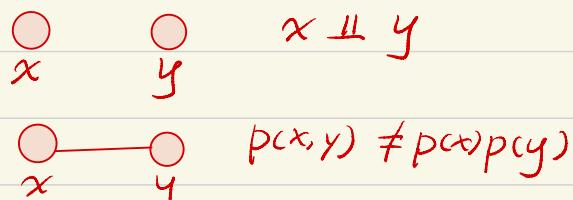
HW is DONE!

WEEK 9 • LEC 16

May 24, 2021

Recall Directed Acyclic graphical models \longleftrightarrow joint probability models

Another class of directed models arise from undirected graphical models.

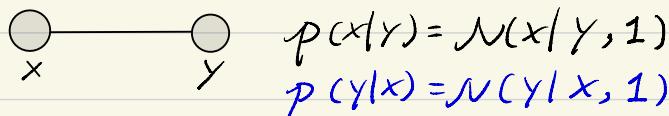


But: $y \perp\!\!\!\perp z \mid x$

In general, $G = \{V, E\}$ where

V = Vertices & RVS

$E = \text{edges}(\text{conditional indept. relation})$



$$\Rightarrow \frac{p(x, y)}{p(x_0, y_0)} = \frac{p(y|x)}{p(y_0|x)} \frac{p(x|y_0)}{p(x_0|y_0)}$$

Brook's Lemma

$$\Rightarrow p(x,y) \propto e^{-\frac{1}{2}(x-y)^2}$$

$$(x-y)^2 = x^2 - 2xy + y^2$$

$$= (x \ y) \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}$$

$$D = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad W = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

Diagonal matrix
with # of neighbors
along diagonal

Adjacency matrix

Generalize:

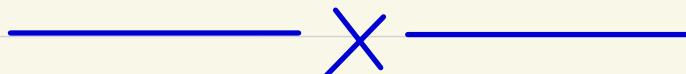
$$x_i | x_{-i} \sim \mathcal{N}\left(\rho \sum_{j=1}^n \frac{w_{ij} x_j}{\sum_{j=1}^n w_{ij}}, \frac{\sigma^2}{\sum_{j=1}^n w_{ij}}\right)$$

$$w_{ii} = 0 \quad \forall i=1, 2, \dots, n$$

$$w_{ij} = \begin{cases} 1 & \text{if } (i,j) \in E \\ 0 & \text{if } (i,j) \notin E \end{cases}$$

$$\sum_{j=1}^n w_{i,j} = m_i = \# \text{ of neighbors of } i.$$

(w_{it} often used)



Homework

Apply Brook's Lemma to the above framework, where

$$x_i | x_{-i}) \sim \mathcal{N} \left(\rho \frac{\sum w_{ij} x_j}{w_{it}}, \frac{\sigma^2}{w_{it}} \right)$$

to show that

$$P(x_1, x_2, \dots, x_n) \propto e^{-\frac{1}{2\sigma^2} x^T (D - \rho W) x}$$

Where $x = (x_1, \dots, x_n)^T$

$$D = \text{Diag}(w_{1t}, w_{2t}, \dots, w_{nt})$$

and W is the binary adjacency matrix.

Hint:

$$\frac{P(x)}{P(\bar{x})} = \prod_{i=1}^n P(x_i | \cdot) / \prod_{i=1}^n P(\bar{x}_i | \cdot)$$

X

As we saw for



the joint distribution $P(x, y)$ need not be well defined.

So: When will

$$P(x) \propto e^{-\frac{1}{2\sigma^2} x^T (D - \rho W) x}$$

be well-defined?

ANS: If $D - \rho W$ is p.d.

i.e. all eigen-values are > 0 .

So over question:

when is $D - \rho W \succ 0$?

Fact: $D - \rho W$ is symmetric because graph is undirected (Relation defined by the graph is a sym relation)
So $w_{ij} = w_{ji}$.

So W is symmetric & D is diagonal.

Can we analytically access the eigenvalues of $D - \rho W$?

$$D - \rho W = D^{\frac{1}{2}} \left(I - \rho D^{-\frac{1}{2}} W D^{-\frac{1}{2}} \right) D^{\frac{1}{2}}$$

McGregory's Decomposition.

$$\Rightarrow D - \rho W = D^{\frac{1}{2}} (I - \rho \tilde{W}) D^{\frac{1}{2}}$$

$$\tilde{W} = \underbrace{D^{-\frac{1}{2}} W D^{-\frac{1}{2}}}_{\text{symmetric}}$$

$\Rightarrow \tilde{W}$ has real eigen-values.

Spectral Decomposition for real symmetric matrices:

$$\tilde{W} = P \Lambda P^T$$

We wish to find when

$$x^T (D - \rho W) x > 0 \quad \forall x \neq 0$$

$$\Leftrightarrow \underbrace{x^T D^{\frac{1}{2}}}_{y^T} (I - \rho \tilde{W}) D^{\frac{1}{2}} x > 0 \quad \forall x \neq 0$$

$$\Leftrightarrow \underbrace{y^T}_{y^T} (I - \rho \tilde{W}) y > 0 \quad \forall y \neq 0$$

Suffices to see when

$$I - \rho \tilde{W} > 0.$$

$$\begin{aligned} I - \rho \tilde{W} &= PP^T - \rho P \Lambda P^T \\ &= \underbrace{P(I - \rho \Lambda)P^T}_{\text{spectral decomp. for } I - \rho \tilde{W}} \end{aligned}$$

$\Rightarrow \{1 - \rho \lambda_1, 1 - \rho \lambda_2, \dots, 1 - \rho \lambda_n\}$
are the eigen-vals of $I - \rho \tilde{W}$.

$I - \rho \tilde{W}$ is p.d. \Leftrightarrow

$$1 - \rho \lambda_i > 0 \quad \forall i = 1, \dots, n$$

Can we find an interval for ρ s.t.

$I - \rho \tilde{W}$ is p.d.?

Take the ordered eigen-vals:

$$\lambda_{(1)} < \lambda_{(2)} < \dots < \lambda_{(n)} : \lambda_i \text{ eigen-vals of } \tilde{W} = D^{-\frac{1}{2}} W D^{-\frac{1}{2}}$$

$$1 - \rho \lambda_{(1)} > 0 \quad \& \quad 1 - \rho \lambda_{(n)} > 0$$

Can we say that $\lambda_{(1)} < 0$ and $\lambda_{(n)} > 0$? If yes, then

$$\underbrace{1 - \rho \lambda_{(1)} > 0}_{\text{iff}} \quad \underbrace{1 + \rho |\lambda_{(1)}| > 0}$$

$$\Rightarrow \rho > -\frac{1}{|\lambda_{(1)}|} = \frac{1}{\lambda_{(1)}}$$

$$\text{Also, } \rho < \frac{1}{\lambda_{(n)}}$$

$$\Rightarrow (I - \rho \tilde{W}) > 0 \text{ iff } \rho \in (\frac{1}{\lambda_{(1)}}, \frac{1}{\lambda_{(n)}})$$

So is $\lambda_{(1)} < 0$ & $\lambda_{(n)} > 0$?

$$\bullet \quad \tilde{W} = D^{-\frac{1}{2}} W D^{-\frac{1}{2}}$$

$$\text{Tr}(\tilde{W}) = \text{Tr}(D^{-\frac{1}{2}} W D^{-\frac{1}{2}})$$

$$= \text{Tr}(D^{-1} W) \quad \begin{cases} W = \begin{pmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{pmatrix} \\ D = \begin{pmatrix} 1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 1 \end{pmatrix} \end{cases}$$

$$= 0 \quad (\text{verify easy})$$

$$\text{But } \text{Tr}(\tilde{W}) = \sum_{i=1}^n \lambda_i = \sum_{i=1}^n \lambda_{(i)}$$

\Rightarrow At least one of the $\lambda_{(i)s}$ is < 0 & at least one of the $\lambda_{(i)s}$ is > 0 .

$$\Rightarrow \lambda_{(1)} < 0 \quad \& \quad \lambda_{(n)} > 0$$



Suppose Prof. Horvath gives KM
a graph representing relationships in a
gene network. Undirected graph.

He asks KM to construct a probability
model as this network.

Response: y_i , $i=1, 2, \dots, n$

Covariates: $x_{i,p}$, $i=1, 2, \dots, n$

$$y_i = x_i^T \beta + \phi_i + \epsilon_i \stackrel{\text{could as well be}}{\equiv} N(0, \sigma^2)$$

a glm

$$\phi_i \sim N(0, \Sigma_\phi) \quad \text{Var}(\phi_i)$$

$$\phi_i = \sum_{j=1}^n w_{ij} \phi_j / w_{ii} + \eta_i, \quad \eta_i \sim N(0, \frac{\sigma^2}{w_{ii}})$$

$$p(\phi) \propto e^{-\frac{1}{2\sigma^2} \phi^T (D-W) \phi}$$

$D-W$ is singular! Because

$$(D-W)\mathbf{1} = \begin{bmatrix} w_{11}-w_{11} \\ \vdots \\ w_{nn}-w_{nn} \end{bmatrix} = \mathbf{0}$$

$$\text{or } \mathbf{1} \in \mathcal{N}(D-W)$$

$$\Rightarrow \text{rank}(D-W) < n$$

$$\text{modify: } \phi_i = P \sum_{j=1}^n \frac{w_{ij} \phi_j}{w_{ii}} + \eta_i$$

$$\Rightarrow p(\phi) \propto e^{-\frac{1}{2\sigma^2} (D-PW) \phi}$$

$D^{-\frac{1}{2}} W D^{-\frac{1}{2}} \leftarrow \text{eigen analysis}$

$$\text{get } \lambda_{(1)}, \lambda_{(n)}, \quad P \sim \text{Unif}(\frac{1}{\lambda_{(1)}}, \frac{1}{\lambda_{(n)}})$$

$$\phi_i \sim N(0, \sigma^2 (D-PW)^{-1})$$

$D-PW$ is called the
Laplacian of the graph.
precision (inverse of the covariance)

WEEK 9 · LEC 17 · May 26 · 2021

Central to modern multivariate statistics & machine learning.

① Simple Linear Reg.

one outcome Y

one explanatory variable X

$$y = \beta_0 + \beta_1 X + \text{Error}$$

② multiple Linear Regression

one outcome: Y

$p+1$ explanatory variables: X_1, \dots, X_p

$$y = \beta_0 + \sum_{j=1}^p X_j \beta_j + \text{Error}$$

\Rightarrow with n measurements:

$$\underset{n \times 1}{y} = \underset{n \times (p+1)}{X} \underset{n \times 1}{\beta} + \underset{n \times 1}{e}$$

③ multivariate (matrix-valued) linear Regression

m outcomes: Y_1, \dots, Y_m

common design matrix X

$$\underset{\text{correlated}}{\underbrace{[Y_1 \dots Y_m]}} = X \underset{\text{correlated}}{\underbrace{[\beta_1 \dots \beta_m]}} + \underset{m \times 1}{[E_1 \dots E_m]}$$

Take a closer look:

m = # of outcomes: Y_1, \dots, Y_m

n = # of observations for each outcome

Let $y_{ij} = i^{\text{th}}$ observation on outcome j ,

let X_i be a $p \times 1$ vector of explanatory variables corresponding to observation i .

$$y_{ij} = X_i^T \beta_j + e_{ij}, \quad i=1, \dots, n, \quad j=1, \dots, m$$

$$\Rightarrow \underset{n \times 1}{y_j} = \underset{n \times 1}{X} \underset{n \times 1}{\beta_j} + \underset{n \times 1}{e_j}, \quad j=1, \dots, m$$

$$X = \begin{bmatrix} X_1^T \\ \vdots \\ X_n^T \end{bmatrix}$$

y_j 's are NOT indept.

So how do we model e_j 's jointly?

If e_j 's were indept., then

$$\text{Cov}(e_j, e_j) = \begin{cases} 0 & \forall j \neq j' \\ \underset{n \times n}{\underbrace{S_j V}} & \forall j = j' \end{cases}$$

Dependence can be introduced as follows:

$$\text{Cov}(e_j, e_j) = \underset{\text{scalar}}{U_{jj}} \underset{\text{non}}{V} \quad \forall j, j'$$

$$y = \underset{m \times 1}{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}} = \underset{m \times n}{\begin{bmatrix} X & \cdots & X \end{bmatrix}} \underset{n \times p}{\begin{bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}} + \underset{m \times 1}{\begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_m \end{bmatrix}}$$

$$\bullet \quad \boxed{y = (I \otimes X) \beta + e} \quad \begin{matrix} m \times 1 \\ \underbrace{m \times m & n \times p}_{m \times m \times p} & m \times 1 \\ m \times m \times p & & m \times 1 \end{matrix}$$

$$\text{Cov}(e) = \text{Cov} \left(\begin{bmatrix} e_1 \\ \vdots \\ e_m \end{bmatrix} \right)$$

$$= [\text{Cov}(e_i, e_j)]_{\substack{i=1, \dots, m \\ j=1, \dots, m}}$$

$$= [U_{ij} V]$$

(Ex) $m=2$:

$$\text{Cov} \left(\begin{bmatrix} e_1 \\ e_2 \end{bmatrix} \right) = \begin{bmatrix} \text{Var}(e_1) & \text{Cov}(e_1, e_2) \\ \text{Cov}(e_2, e_1) & \text{Var}(e_2) \end{bmatrix}$$

$$= \begin{bmatrix} U_{11} V & U_{12} V \\ U_{21} V & U_{22} V \end{bmatrix}$$

$$\stackrel{\text{def}}{=} U \otimes V, \quad U = \begin{bmatrix} U_{11} & U_{12} \\ U_{21} & U_{22} \end{bmatrix}$$

In general,

$$\text{Cov}(e) = U \otimes V$$

$m \times m$ outcomes $n \times n$ measurements

$$y \sim \mathcal{N}((I \otimes X)\beta, U \otimes V)$$

~~X~~ Digression ~~X~~

Kronecker Products

Def: $A \otimes B = [a_{ij} B]$

$(ij)^{\text{th}}$ element of A

$$= \begin{bmatrix} a_{11}B & a_{12}B & \cdots & a_{1p}B \\ \vdots & \ddots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \cdots & a_{mp}B \end{bmatrix}$$

$m \times n \times p \times q$

Start simple:

$$(x^T \otimes B) = \underbrace{[x_1 B : x_2 B : \cdots : x_p B]}_{n \times p \times q}$$

$$\begin{bmatrix} x^T \\ y^T \end{bmatrix} \otimes B = \begin{bmatrix} x_1 B & \cdots & x_p B \\ y_1 B & \cdots & y_p B \end{bmatrix} = \begin{bmatrix} x^T \otimes B \\ y^T \otimes B \end{bmatrix}$$

$$A \otimes B = \begin{bmatrix} a_{1*}^T \\ \vdots \\ a_{m*}^T \end{bmatrix} \otimes B$$

$$= \begin{bmatrix} a_{1*}^T \otimes B \\ \vdots \\ a_{m*}^T \otimes B \end{bmatrix}$$

- Understanding products of Kronecker products.

$$(A \otimes B) \cdot (C \otimes D) = ?$$

Start simple: $m=1, r=1$,

$$(x^T \otimes B) \cdot (y \otimes D)$$

$$= [x_1 B : x_2 B : \cdots : x_p B] \begin{bmatrix} y_1 D \\ y_2 D \\ \vdots \\ y_p D \end{bmatrix}$$

$$= \sum_{i=1}^p x_i y_i B D = \left(\sum_{i=1}^p x_i y_i \right) B D$$

$$= x^T y \otimes BD$$

Return to general case:

$$A \otimes B = \begin{bmatrix} a_{1*}^T \otimes B \\ \vdots \\ a_{m*}^T \otimes B \end{bmatrix}$$

$$C \otimes D = [c_{*1} \otimes D : \cdots : c_{*r} \otimes D]$$

$$(A \otimes B) \cdot (C \otimes D)$$

$$= \begin{bmatrix} a_{1*}^T \otimes B \\ \vdots \\ a_{m*}^T \otimes B \end{bmatrix} \begin{bmatrix} c_{*1} \otimes D : c_{*2} \otimes D : \cdots : c_{*r} \otimes D \end{bmatrix}$$

$$= \left[(a_{1*}^T \otimes B) \cdot (c_{*j} \otimes D) \right]_{i,j=1,\dots,m}$$

$$= \left[\underbrace{(a_{1*}^T c_{*j})}_{\text{scalar}} B D \right]_{i,j=1,2,\dots,m}$$

$$= AC \otimes BD$$

Thus,

$$(A \otimes B) \cdot (C \otimes D) = AC \otimes BD$$

If A, B are non-singular (square & invertible), then

$$(A \otimes B) \cdot (A^{-1} \otimes B^{-1}) = (AA^{-1}) \otimes (BB^{-1})$$

$\Rightarrow A \otimes B$ is non-singular with

$$(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$$

X

Return to $\text{Cov}(y) = \text{Cov}(\epsilon) = U \otimes V$

$$y \sim \mathcal{N}[(I \otimes X)\beta, U \otimes V]$$

$$Z = y - (I \otimes X)\beta$$

$$P(Z) = \frac{1}{(2\pi)^{\frac{mn}{2}}} \frac{1}{|U \otimes V|^{\frac{1}{2}}} e^{-\frac{1}{2} Z^T (U \otimes V)^{-1} Z}$$

- U and V are 2 p.d. matrices.

Then is $U \otimes V$ p.d.?

Ans: Yes.

Proof: $U = P_u \Lambda_u P_u^T$ (spectral)

$$V = P_v \Lambda_v P_v^T$$
 (spectral)

$$P_u P_u^T = I_m, P_v P_v^T = I_n,$$

$$U \otimes V = (P_u \Lambda_u P_u^T) \otimes (P_v \Lambda_v P_v^T)$$

$$\stackrel{\text{by lemma}}{=} (P_u \Lambda_u \otimes P_v \Lambda_v) \cdot (P_u^T \otimes P_v^T)$$

$$= (P_u \otimes P_v) \cdot (\Lambda_u \otimes \Lambda_v) \cdot (P_u^T \otimes P_v^T)$$

Quick result:

$$(A \otimes B)^T = [a_{ij}B]_j^T = [a_{ij}B^T]_j = A^T \otimes B^T$$

\Rightarrow

$$U \otimes V = \underbrace{(P_u \oplus P_v)}_{\widetilde{P}} \underbrace{(\Lambda_u \otimes \Lambda_v)}_{\widetilde{\Lambda}} \underbrace{(P_u \otimes P_v)}_{\widetilde{P}^T}$$

diag($\lambda_u^{(c_i)}, \lambda_v^{(c_j)}$)

Is $\widetilde{P} \widetilde{\Lambda} \widetilde{P}^T = U \otimes V$ a spectral decomposition?

Ans: only if \widetilde{P} is orthogonal.

$$\begin{aligned} & (P_u \otimes P_v) \cdot (P_u^T \otimes P_v^T) \\ &= P_u P_u^T \otimes P_v P_v^T \\ &= I_{mn} \end{aligned}$$

So $U \otimes V = \widetilde{P} \widetilde{\Lambda} \widetilde{P}^T$, where $\widetilde{\Lambda}$'s diagonal elements are the eigen-vals of $U \otimes V$.

\therefore the eigen-vals of $U \otimes V$ are the set

$$\{ \lambda_u^{(i)} \lambda_v^{(j)} : i=1, \dots, m, j=1, \dots, n \}$$

\Rightarrow all > 0 since $\begin{cases} \lambda_u^{(i)} > 0 \\ \lambda_v^{(j)} > 0 \end{cases}$

$\Rightarrow U \otimes V$ is positive definite.

$\therefore \mathcal{N}((I \otimes X)\beta, U \otimes V)$ is a valid distribution!

X

WEEK 10 • LEC 18 • June 02, 2021

$$\begin{aligned}
 \bullet \det(U \otimes V) &= \det(\tilde{P} \tilde{\Lambda} \tilde{P}^T) \\
 &= \det(\tilde{\Lambda}) \\
 &= \det(I_u \otimes I_v) \\
 &= \det \begin{pmatrix} \lambda_u^{(1)} I_v & 0 \\ 0 & \ddots \\ 0 & \lambda_v^{(m)} I_u \end{pmatrix} \\
 &= \prod_{i=1}^m \det(\lambda_u^{(i)} I_v) \\
 &= \prod_{i=1}^m \left\{ [\lambda_u^{(i)}]^n \prod_{j=1}^n \lambda_v^{(j)} \right\} \\
 &= (\lambda_u^{(1)} \cdots \lambda_u^{(m)})^n [\det(I_v)]^m \\
 &= (\det(P_u \Lambda_u P_u^T))^n (\det(P_v \Lambda_v P_v^T))^m \\
 &= (\det I_u)^n (\det I_v)^m
 \end{aligned}$$

• $Z = Y - (I \otimes X)\beta$ has density

$$p(z) = \frac{1}{(2\pi)^{\frac{mn}{2}}} \frac{1}{|U|^{\frac{n}{2}} |V|^{\frac{m}{2}}} e^{-\frac{1}{2} z^T (U' \otimes V') z}$$

□.

Next step: Construct a density for the $m \times n$ matrix

$$\begin{aligned}
 Y &= [y_{ij}]_{n \times m} \\
 &= [y_1, \dots, y_m]
 \end{aligned}$$

Recall: Matrix-variate regression

$$y_{ij} = x_i^T \beta_j + e_{ij}; \quad i=1, \dots, n$$

$$y_j = X \beta_j + e_j; \quad j=1, 2, \dots, m$$

$$\underbrace{[y_1 : y_2 : \dots : y_m]}_{Y_{n \times m}} = \underbrace{X}_{p \times m} \underbrace{[\beta_1 : \beta_2 : \dots : \beta_m]}_{B_{p \times m}} + \underbrace{[e_1 : e_2 : \dots : e_m]}_{E_{n \times m}}$$

Matrix-variate regression:

$$Y = XB + E$$

Find a distribution for E , hence for Y .

Recall from last Lecture:

$$\begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_m \end{bmatrix}$$

$$\sim \mathcal{N}(0, U \otimes V)$$

- $U \otimes V$ is p.d. as long as U and V are p.d.
- $(U \otimes V)^{-1} = U^{-1} \otimes V^{-1}$
- $\det(U \otimes V) = (\det U)^n (\det V)^m$
- $(A \otimes B)(C \otimes D) = (AC \otimes BD)$

Define: $\text{Vec}(\cdot)$ operator

$$\text{Vec}(\cdot) : \underbrace{\mathcal{M}_{n \times m}}_{\text{matrices}} \mapsto \mathbb{R}^{nm \times 1}$$

Let A be an $n \times m$ matrix

$$A = [a_{*1} : a_{*2} : \dots : a_{*m}]$$

$$\text{Vec}(A) = \begin{bmatrix} a_{*1} \\ a_{*2} \\ \vdots \\ a_{*m} \end{bmatrix} \in \mathbb{R}^{nm \times 1}$$

Back to matrix-variate linear model:

$$Y = X\beta + E$$

$$\text{Vec}(E) = \begin{bmatrix} e_1 \\ \vdots \\ e_m \end{bmatrix}$$

$$\text{Vec}(E) \sim \mathcal{N}(0, U \otimes V)$$

$p(\text{Vec}(E))$ is

$$\frac{1}{(2\pi)^{\frac{mn}{2}}} \frac{1}{|U \otimes V|^{\frac{1}{2}}} e^{-\frac{1}{2} \text{Vec}(E)^T (U \otimes V)^T \text{Vec}(E)}$$

$$\Rightarrow \frac{1}{(2\pi)^{\frac{mn}{2}}} \frac{1}{|U|^{\frac{n}{2}} |V|^{\frac{m}{2}}} e^{-\frac{1}{2} \text{Vec}(E)^T (U^T \otimes V^T) \text{Vec}(E)}$$

want to replace with matrix E

Goal: write the above density as a function of E (not $\text{Vec}(E)$) with parameters U and V , no \otimes .

Need to simplify:

$$\underbrace{\text{Vec}(E)^T}_{\text{long vec}} (U^{-1} \otimes V^{-1}) \underbrace{\text{Vec}(E)}_{\text{long vec}}$$

Digression vec algebra

$$\text{Suppose } Y = [y_1 : \dots : y_m]^{nm \times 1}$$

$$X = [x_1 : \dots : x_m]^{n \times m}$$

$$\text{Vec}(Y)^T \text{Vec}(X) = \text{Tr}(Y^T X) \\ = \text{Tr}(X Y^T)$$

Next question:

$$(U^T \otimes V^T) \text{Vec}(E) = \text{Vec}(?)$$

• Second digression

$$(A \otimes B) \underbrace{\text{Vec}(X)}_{np \times 1} = \text{Vec}(?)$$

$$X = [x_1 : x_2 : \dots : x_n]^{pxn}$$

$$(A \otimes B) \text{Vec}(X) = \begin{bmatrix} a_{11}B & a_{12}B & \dots & a_{1n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \dots & a_{mn}B \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$= \begin{bmatrix} \sum_{j=1}^n a_{1j} B x_j \\ \vdots \\ \sum_{j=1}^n a_{mj} B x_j \end{bmatrix}_{mn \times 1}$$

$$= B \left(\begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} a_{11} \\ a_{12} \\ \vdots \\ a_{1n} \end{bmatrix} \right)$$

$$\Rightarrow B(X a_{1*})$$

$\Rightarrow i^{\text{th}}$ row(block) of $(A \otimes B) \text{Vec}(X)$

$E: n \times m$

Hence

$$(A \otimes B) \text{vec}(X)$$

$$= \begin{bmatrix} BX\alpha_{1*} \\ BX\alpha_{2*} \\ \vdots \\ BX\alpha_{m*} \end{bmatrix}$$

$$= \text{vec}(BXA^T)$$

(Since α_{i*} is the i^{th} row of A)

\Rightarrow

$$(A \otimes B) \text{vec}(X) = \text{vec}(BXA^T)$$

— Digression Ends —

Back to our question:

$$(\bar{U}^T \otimes \bar{V}^T) \text{Vec}(E) = \text{Vec}(\bar{V}^T E \bar{U})$$

note: $(\bar{U}^T)^T = \bar{U}^T$ since U sym.

So using first digression:

$$\text{Vec}(E)^T \text{Vec}(\bar{V}^T E \bar{U})$$

$$= \text{Tr}(E^T \bar{V}^T E \bar{U})$$

$$= \text{Tr}(\bar{V}^T E \bar{U}^T E^T)$$

Thus,

$$P(E) = \frac{1}{(2\pi)^{\frac{nm}{2}}} \frac{1}{|\bar{U}|^{\frac{n}{2}} |\bar{V}|^{\frac{m}{2}}} e^{-\frac{1}{2} \text{vec}(E)^T (\bar{U}^T \bar{V}^T) \text{vec}(E)}$$

$$\Rightarrow \frac{1}{(2\pi)^{\frac{nm}{2}}} |\bar{U}|^{-\frac{n}{2}} |\bar{V}|^{-\frac{m}{2}} e^{-\frac{1}{2} \text{Tr}(E^T \bar{V}^T E \bar{U})}$$

$$= MN(E | 0, U, V)$$

$\Rightarrow E = Y - XB$ leads to

$$P(Y) = (2\pi)^{-\frac{nm}{2}} |\bar{U}|^{-\frac{n}{2}} |\bar{V}|^{-\frac{m}{2}} e^{-\frac{1}{2} \text{Tr}[(Y - XB)^T \bar{V}^T (Y - XB) \bar{U}]}$$

$$= MN(Y | XB, U, V)$$

matrix normal

$U \leftarrow$ Covariance matrix among
columns of Y .

$V \leftarrow$ Covariance matrix among
rows of Y .

Homework:

Q1: ① $\text{Vec}(AB) = (?) \text{vec}(B)$

$$= (?) \text{vec}(A)$$

② $\text{Vec}(ABC) = (?) \text{vec}(C)$

$$= (?) \text{vec}(A)$$

Q2: Generating from a Matrix-normal.

- Devise an algorithm to generate

$$Y \sim MN(M, U, V)$$

WITHOUT using Kronecker product.

- Matrix Normal - Inverse Wishart

MNIW

Inverse Wishart: Family of P & q
on positive definite covariance matrices.

• Extends inverse-Gamma to multivariate.
 Σ is $m \times m$ covariance matrix

$$p(\Sigma | v, S) = IW(\Sigma | v, S)$$

$$\propto |\Sigma|^{-\frac{v+m+1}{2}} e^{-\frac{1}{2} \text{Tr}(S\Sigma^{-1})}$$

$MNIW(B, \Sigma | C, U, v, S)$

$$= IW(\Sigma | v, S) \times MN(B | C, U, \Sigma)$$

$\underbrace{p(\Sigma)}_{p(\Sigma)}$ $\underbrace{MN(B | C, U, \Sigma)}_{p(B | \Sigma)}$

General Framework:

$$Y | B, \Sigma \sim MN(XB, \Sigma, I_n)$$

$n \times m$ $n \times p$ $p \times m$

$$B | \Sigma \sim MN(C, U, \Sigma)$$

$$\Sigma \sim IW(v, S)$$

- $p(B, \Sigma | Y) \propto$

$$MNIW(B, \Sigma | C, U, v, S) \times$$

$$MN(Y | XB, \Sigma, I_n)$$

$$\propto IW(\Sigma | v, S) \times MN(B | C, U, \Sigma)$$

$$MN(Y | XB, \Sigma, I_n)$$

$$\propto |\Sigma|^{-\frac{v+m+1}{2}} e^{-\frac{1}{2} \text{Tr}(S\Sigma^{-1})} \times |U|^{-\frac{m}{2}} |\Sigma|^{-\frac{p}{2}} \times$$

$$e^{-\frac{1}{2} \text{Tr}[(B-C)^T \Sigma^{-1} (B-C) U^{-1}]} \times$$

$$|\Sigma|^{-\frac{n}{2}} |I_n|^{-\frac{m}{2}} e^{-\frac{1}{2} \text{Tr}[(Y-XB)^T (Y-XB) \Sigma^{-1}]}$$

keeping track of this!

$$\propto IW(\Sigma | v^*, S^*)$$

$$\times MN(B | M, V, \Sigma)$$

Find M & V !

- Question: $p(B, \Sigma | Y) ??$

It turns out that

$$p(B, \Sigma | Y) = MNIW(B, \Sigma | M, V, v, S)$$