# Mathematical Statistics II

## STA2212H S LEC9101

Week 11 March

31 2021 Start

recording!



link

1.  HW 10 (current) not due until Monday April 5
    HW 11 posted April 2 due April 9
    Take-home posted April 9 due April 19                    April 7 if I can
    No class Friday April 2; no office hour Monday April 5

2.  Course evaluations available until April 12
3.  Bayesian methods for text classification; discriminant analysis
4.  Intro to graphical models and causality

- Multivariate normal distribution $f(x; \mu, \Omega) = \frac{1}{(\sqrt{2\pi})^k} |\Omega|^{-\frac{1}{2}} e^{-\frac{1}{2}(x-\mu)^T \Omega^{-1}(x-\mu)}$

- $X \sim N_k(\mu, \Omega)$ $\quad \hat{\mu} = \bar{x} \quad \hat{\Omega} = \frac{1}{n} S = \cdots$

  $x_1, \cdots, x_n \, iid$

- correlation $\rho_{ij} \;\; rs \;\; = \;\; cov(X_r, X_s) / \sqrt{var(x_n) var(x_s)}$

  $-1 \leq \rho_{rs} \leq 1$

  $\rho_{rs} = 0 = X_r \perp\!\!\!\perp X_s$
  
  $(Normality)$

- partial correlation $\rho_{ij \mid -(i,j)}$
  
  $\sqrt{s} \quad r_{,s}$

$\rho_{rs \mid -(r,s)} \qquad corr^n$ between $X_n, X_s,$

$\qquad\qquad\qquad conditional$ on all other $X's \qquad \Omega_{22}$

$\downarrow$

$\Omega^{rs} = (\Omega^{-1})_{rs} = cov(X_n, X_s \mid X_{-(r,s)})$

$\qquad\qquad\qquad\qquad\qquad f(x) = f(x_{(1)} \mid x_{(2)}) f(x_{(2)})$

$\qquad\qquad \uparrow Mar 24 \; \Omega_{rs} - \cdots \; N \uparrow \qquad\qquad\qquad \uparrow N$

- $X \sim Mult_{(k)}(n; p)$    $j = 1, \ldots, k$   categories     $X_j$ = number of obs in category $j$

- $\mathrm{pr}(X_1 = x_1, \ldots, X_k = x_k; p) = \dfrac{n!}{x_1! \cdots x_k!} \, P_1^{x_1} \cdots P_k^{x_k}$    $0 \le p_j \le 1$   $\sum_{j=1}^{k} p_j = 1$

                                                        $\sum x_j = n$

- $\mathrm{E}(X) = np_k$

- $\mathrm{cov}(X) =$

$$\begin{bmatrix} np_1(1-p_1) & -np_1p_2 & \cdots & -np_1 p_k \\ & & & \\ & & & \\ -np_1p_k & \cdots & & np_k(1-p_k) \end{bmatrix}$$

                                        AoS Thm 14.4

- $\hat{p} = \dfrac{X}{n}$

- $\mathrm{cov}(\hat{p}) = \dfrac{\mathrm{cov}(X)}{n^2}$

- $X \sim Mult_k(n; p)$

- $\mathrm{pr}(X_1 = x_1, \ldots, X_k = x_k; p) =$

- $\mathrm{E}(X) =$

- $\mathrm{cov}(X) =$

- $\hat{p} = \dfrac{X}{n}$

- $\mathrm{cov}(\hat{p}) =$   ntbc

$\left( \text{loglinear models} \right) \longrightarrow$

$[ \text{see Ch 15 of AoS} ]$

$X_j =$ number of obs in category $j$

$P_k = P_k(\theta) \quad \dim \theta < k - 1$   Contingency tables

$1 \qquad\qquad \cdots\cdots \qquad\qquad k$

constrained

$\boxed{n P_1(\theta)} \quad \cdots \cdots \quad \boxed{n P_k(\theta)}$

AoS Thm 14.4

$\dfrac{X_1}{n\hat{p}_1} \qquad\qquad \dfrac{X_k}{n\hat{p}_k}$

$\downarrow \quad \chi^2$ test compares these

| | 1 $\cdots$ m |
|---|---|
| 1 : : $\ell$ | $X_{1k}$ |

$X_{\ell m}$

$X = (X_{ij}) \qquad k = m \times \ell$

$i = 1, \ldots, \ell \quad j = 1, \ldots, m$

# Overview

- multivariate and multinomial distributions used to study the joint distribution

- analogous to unsupervised learning ← learning relationships among variables on an equal footing

- AoS §15.1,2: 2 binary variables; 2 discrete variables          multinomial

- AoS §14.2: pairs of normal variables

- AoS §15.4 one discrete, one continuous variable

  Compare    $\text{cont}^s$
  $$F_1(\underline{x}) \text{ to } F_2(\underline{x})$$
  discrete

- AoS Ch.15 Inference about independence

$h(x): \; x \rightarrow y$

$$\hat{h}(\underline{x}) = \begin{matrix} 1 \\ 0 \end{matrix} \; \text{if} \quad \left. \begin{matrix} \hat{r}(x) > \frac{1}{2} \\ < \frac{1}{2} \end{matrix} \right\} \text{binary}$$

Bayes rule is

$\underset{k}{max} \quad \dfrac{\boxed{f_k(\underline{x})\pi_k}}{\sum_{r=1} f_r(\underline{x})\pi_r}$

$\left( \hat{r}(x=k) \text{ if} \atop \begin{matrix} \text{are predict} \\ y \mid X = k \end{matrix} \right)$

could use parametric models for $\hat{f}_k \overset{(\text{2-class})}{\longrightarrow} MVN$

$\longrightarrow (\overset{LDA}{\&DA})$

or non par models $\overset{\text{to get}}{\;} \hat{f}_k$

$\underset{k}{arg\,max} \quad \dfrac{\hat{f}_k(x)\hat{\pi}_k}{\sum \hat{f}_r(x)\hat{\pi}_r} \qquad \hat{\pi}_r = \dfrac{\#(Y's \text{ in class } r)}{n}$

BCE loss $\longleftarrow$ binary cross-entropy

$-\frac{1}{n} \sum_{i=1}^{n} \{ y_i \log p_i + (1-y_i)\log(1-p_i) \} \qquad y_i \begin{matrix} 0 \\ 1 \end{matrix}$

$= -\frac{1}{n} \ell(p; y)$ under $Y_i \sim Ber(p_i)$ $\quad p_i = p(y_i=1)$

"log-likelihood loss" $\qquad \hat{p}_i = y_i$

supervised vs unsupervised        M/L        Stats        categorical vs continuous

some labelled data $y_1, \cdots, y_n$

some features $\underline{x}_1, \cdots, \underline{x}_n$

use features to predict $y$

$$\zeta(x) = P_n(y = 1 \mid \underline{x}) = \frac{e^{\underline{x}^T \beta}}{1 + e^{\underline{x}^T \beta}}$$

$$(e.g.)$$

text analysis

$$E(y \mid x) = \underline{x}^T \beta$$

or

$$= e^{\underline{x}^T \beta}$$

or

$$= e^{\underline{P(\underline{x})}^T \underline{\theta}}$$

map $h: \mathcal{X} \to \mathcal{Y}$        classifier

$$L(h) = P\{h(x) \neq y\} \qquad \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}\{h(x_i) \neq y_i\}$$

$$\hat{L}_n(h)$$

- Markov chains      SM §6.1
- continuous time Markov models     finite state space $\mathcal{S}$     SM §6.2
- Markov random fields; directed acyclic graphs     SM §6.2
- Multivariate normal     SM §6.3
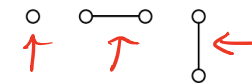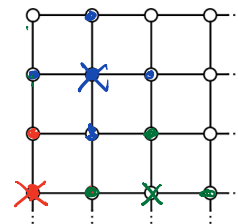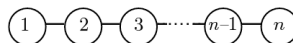- Time series     SM §6.4
- Point processes     SM §6.5

- $X_1, \ldots, X_n$ a random sequence
- Markov property:

$$\mathrm{pr}(X_{j+1} = x_{j+1} \mid \underbrace{x_j}_{X_j =}, \ldots, \underbrace{x_1}_{X_1 =}) = P_n(X_{j+1} = x_{j+1} \mid \underbrace{x_j}_{X_j =})$$

- set of sites $\mathcal{J} = \{1, \ldots, n\}$  $\quad p_n(X_j = x_j \mid X_{-j} = x_{-j}) = P_n(X_j = x_j \mid X_{\mathcal{N}_j} = x_{\mathcal{N}_j})$
- neighbourhood system $\mathcal{N} = \{\mathcal{N}_j, j \in \mathcal{J}\}$

**Figure 6.4** Markov random fields. Left: neighbourhood structure for first-order Markov chain and its cliques and their subsets. Right: first-order neighbourhood structure, cliques and their subsets for rectangular grid of sites.
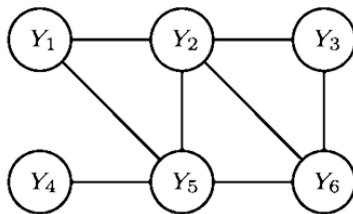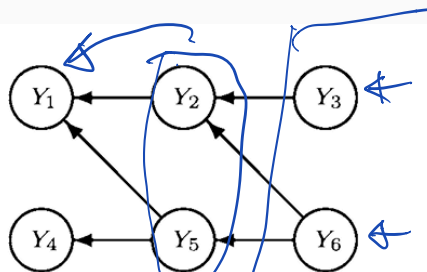
labels $y_1, \ldots, y_n$

$\not{y}_{\mathfrak{d}_1} \quad \not{y}_{\mathfrak{d}_2} \quad \cdots \quad \not{y}_{\mathfrak{d}^{m \times m}}$

$i = 1, \ldots, n$

$\not{x} \in \{1, \ldots, k\}$
$\not{x} \in \mathbb{R}$   measures intensity

e.g.

- can be convenient for studying relationships between variables $\leftarrow$ "mechanistic"
- through a probability distribution on the graph
- that is specified by a factorization of the joint density

parents          nodes
                 edges
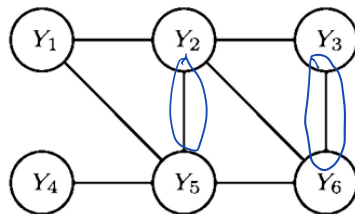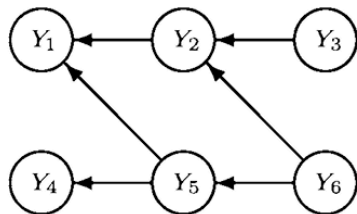                 directed edges



**Figure 6.6** Directed acyclic and moral graphs. Left: directed acyclic graph representing (6.17). Right: moral graph, formed by moralizing the directed acyclic graph, that is, 'marrying' parents and dropping arrowheads.

cycles
cliques
colliders

$$f(y_1,...,y_6) = f(y_6)\,f(y_5|y_6)\,f(y_2|y_3,y_6)\,f(y_1|y_2,y_5)\,f(y_2|y_3)$$
$$\rightarrow f(y) = \prod_{j\in1} f(y_j|\text{parents of } y_j)\qquad f(y_4|y_5)\,f(y_3)\qquad (\text{check SM})$$

**Figure 6.6** Directed acyclic and moral graphs. Left: directed acyclic graph representing (6.17). Right: moral graph, formed by moralizing the directed acyclic graph, that is, 'marrying' parents and dropping arrowheads.

31/2 and p.250

graph ← maybe motivated by an app⁼
       & cond'l independence encoded
need a prob. dist⁼ on graph ← modelling
prep. MRF
we convert a DAG to a MRF by: removing
                                        arrows
                                joining parents

**Figure 6.7** Directed acyclic graph representing the incidence and presentation of six possible diseases that would lead to a 'blue' baby (Spiegelhalter *et al.*, 1993). LVH means left ventricular hypertrophy.

Handwritten annotations:

(6.18)

$$f(y) = \prod_{j \in J} f(y_j \mid \text{parents } y_j)$$

These models encode structured dependence

node

$Y_j \quad j = 1, \ldots, \#\text{nodes}$

dist$^n$ for $Y_j$

$\text{pr}(Y_j \in \text{class } k)$

Mult( ... )

Mult( ... )

(MV Normal)

$\Omega^{-1}$

parent

pa

Nodes in figure:
1: Birth asphyxia?
2: Disease?
3: Age at presentation?
4: LVH?
5: Duct flow?
6: Cardiac mixing?
7: Lung parenchema?
8: Lung flow?
9: Sick?
10: Hypoxia distribution?
11: Hypoxia in O2?
12: CO2?
13: Chest X-ray?
14: Grunting?
15: LVH report?
16: Lower body O2?
17: Right up. quad. O2?
18: CO2 report?
19: X-ray report?
20: Grunting report?

264

6 · Stochastic Models

**Figure 6.10** Graphs for the full model (left) and a reduced model (right) for the maths marks data. The interpretation of the reduced model is that given the result for algebra, results for vectors and mechanics are independent of those for analysis and statistics.
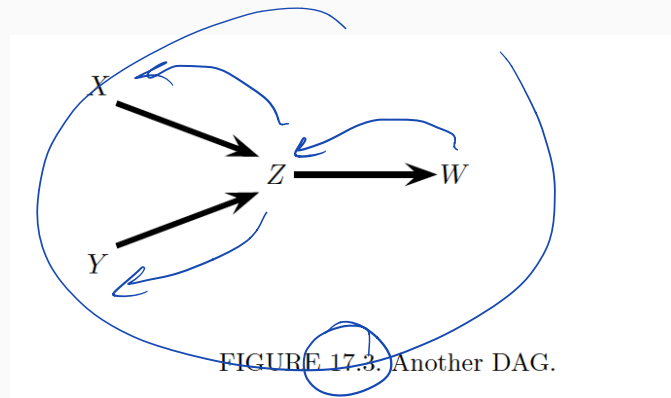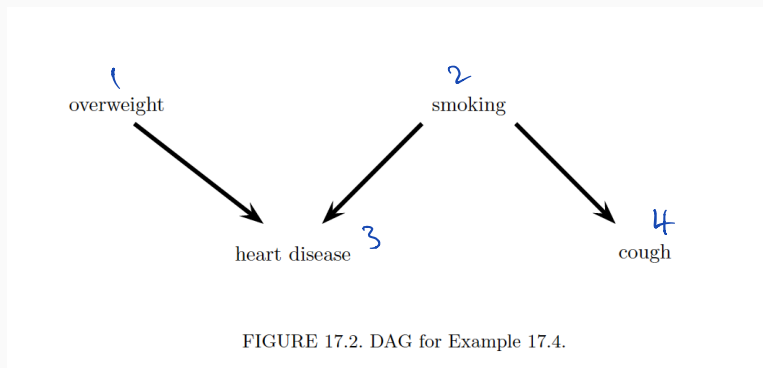


$$\Omega^{rs} \neq \Omega^{-1}_{rs}$$
$$\|$$
$$(\Omega^{-1})_{rs} \neq \Omega^{-1}_{rs}$$

complete indirected graph

simplified graph

if $\left(\Omega^{rs}\right) = 0$ then $X_r \perp\!\!\!\perp X_s \mid X_{-(r,s)}$
$$= (\Omega^{-1})_{r,s}$$

- Fig 17.2, Example 17.4
- Fig 17.3



FIGURE 17.2. DAG for Example 17.4.



FIGURE 17.3. Another DAG.

$$f(y) = f(y_3 | y_2, y_1) \cdot f(y_1) f(y_2) f(y_4 | y_2)$$

encodes dep. in graph.

$$f(x,y,z,w) = f(w|z) f(z|x,y) f(x) f(y)$$

FIGURE 17.2. DAG for Example 17.4.

**17.4 Example.** Figure 17.2 shows a DAG with four variables. The probability function for this example factors as

$$f(\text{overweight}, \text{smoking}, \text{heart disease}, \text{cough})$$
$$= \quad f(\text{overweight}) \times f(\text{smoking})$$
$$\times \quad f(\text{heart disease} \,|\, \text{overweight}, \text{smoking})$$
$$\times \quad f(\text{cough} \,|\, \text{smoking}). \quad \blacksquare$$

**17.5 Example.** For the DAG in Figure 17.3, $\mathbb{P} \in M(\mathcal{G})$ if and only if its probability function $f$ has the form

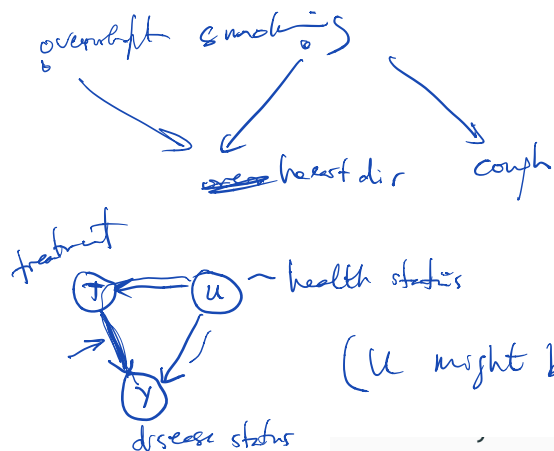$$f(x, y, z, w) = f(x)f(y)f(z \,|\, x, y)f(w \,|\, z). \quad \blacksquare$$

A₀S (17.8)
SM 9.1

DAGs are related to causality
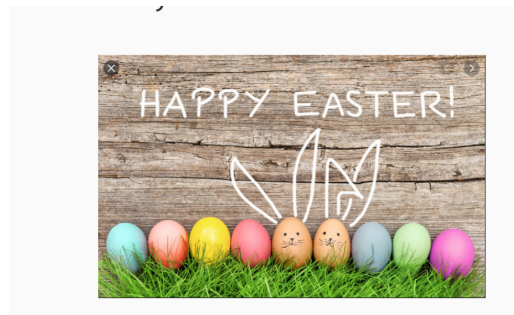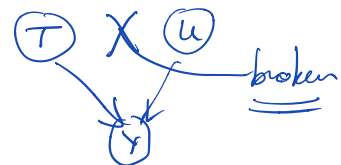value of y at each parent node is a potential cause
of the response at a child node

Causality Ch16
AoS
"counterfactuals"

overweight  smoking

heart dir     cough

treatment
T → U  ~ health status
↓ ↓
Y
disease status

(U might be unseen)

probs. model to encode all
possible dep. of interest;
see if data supports this,
or a simpler model

if we randomize but assign t

T   X   U
↘   ↓   broken
  Y



April 7    — causality + (DAG)
      9    — visualization

April 12   office hours