

# Mathematical Statistics II

STA2212H S LEC9101

Week 9

March 11 2025

## THE CONVERSATION

Academic rigour, journalistic flair



Microplastics are tiny bits of plastic that show up in the environment. Svetlozar Hristov/iStock via Getty Images Plus

**What's that microplastic? Advances in machine learning are making identifying plastics in the environment more reliable**

Published: March 6, 2025 8.35am EST

**Ambuj Tewari**

Professor of Statistics, University of Michigan



Microplastics are tiny bits of plastic that show up in the environment. Svetlozar Hristov/iStock via Getty Images Plus

## What's that microplastic? Advances in machine learning are making identifying plastics in the environment more reliable

Published: March 6, 2025 8.35am EST

**Ambuj Tewari**

Professor of Statistics, University of Michigan

Microplastics – the tiny particles of plastic shed when litter breaks down – are everywhere, from the deep sea to Mount Everest, and many researchers worry that they could harm human health.

I am a machine learning researcher. With a team of scientists, I have developed a tool to make identification of microplastics using their unique chemical fingerprint more reliable. We hope that this work will help us learn about the types of microplastics floating through the air in our study area,

# Today

1. Recap Mar 4 choosing test stats, hypothesis/significance testing, multiple testing
2. Nonparametric tests, goodness-of-fit
3. Introduction to causal inference
4. Reviewing project guidelines
5. Conformal prediction

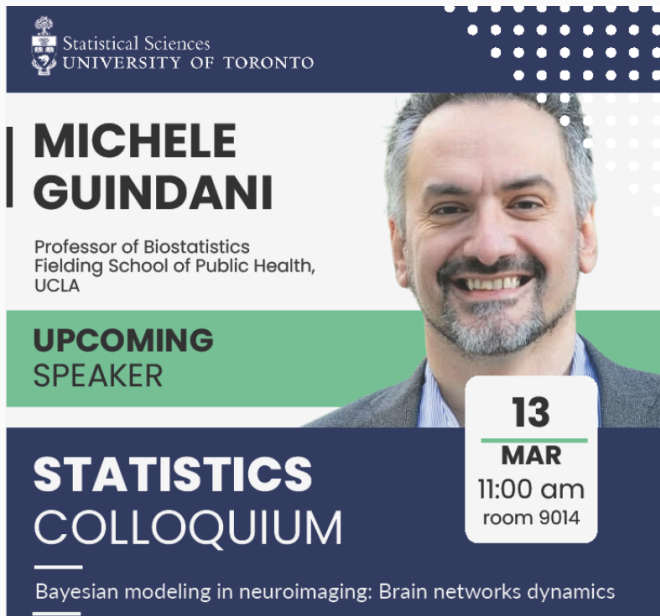
## Upcoming seminar

Department Seminar Thursday March 6 11.00 – 12.00

Hydro Building, Room 9014

Bayesian modelling in neuroimaging

Michele Guindani, UCLA



Statistical Sciences  
UNIVERSITY OF TORONTO

**MICHELE  
GUINDANI**

Professor of Biostatistics  
Fielding School of Public Health,  
UCLA

**UPCOMING  
SPEAKER**

**13  
MAR**  
11:00 am  
room 9014

**STATISTICS  
COLLOQUIUM**

Bayesian modeling in neuroimaging: Brain networks dynamics

link

## Project Guidelines

STA 2212S: Mathematical Statistics II 2025

Presentation on April 1, 2025.

Report submission due April 16, 2025.

### Part 1: Presentation [10 points]

On the last day of class (April 1), you will present your final project. This includes:

- Emailing a .pdf version of your team's slide deck pdf to `nancym.reid@utoronto.ca` by **09.00 April 1**. You are responsible for the slides corresponding to your sections of the write-up. Please email **one** complete version for each team.
- Presenting the slides in no more than 10 minutes; each team member to present for no more than 5 minutes.

# Recap

$$X_1, \dots, X_n \sim f(\mathbf{x}; \theta), \theta \in \Theta \subset \mathbb{R}^p$$

- testing  $H_0 : \theta \in \Theta_0$



Hyp. f.

- rejection region  $\{\mathbf{x} : t(\mathbf{x}) > c_\alpha\}$

$$\Pr_{H_0}\{t(\mathbf{X}) \geq c_\alpha\} \leq \alpha$$

Sign. f.

- p-value:

$$\Pr_{H_0}\{t(\mathbf{X}) \geq t(\mathbf{x}^{obs})\}$$

against some alternative  
simple or composite H

$t \geq c_\alpha$  "reject"  $H_0$   
 $t < c_\alpha$  don't

large values

- significance function ( $\theta \in \mathbb{R}$ )

$$p(\theta) = \Pr_\theta\{t(\mathbf{X}) \geq t(\mathbf{x}^{obs})\}$$

$$\theta \in \Theta$$

# Recap: Choosing test statistics $t(\cdot)$

1. Optimal choice – Neyman-Pearson lemma

Might be UMP (HW 7)

2. Pragmatic choice – likelihood-based test statistics

score test  
Wald  
LRT

3. Pragmatic choice – nonparametric test statistics

(a) Need to know distribution of test statistic under  $H_0$

$C_\alpha$   
p-value

(b) Test statistic should be large when  $H_0$  is not true

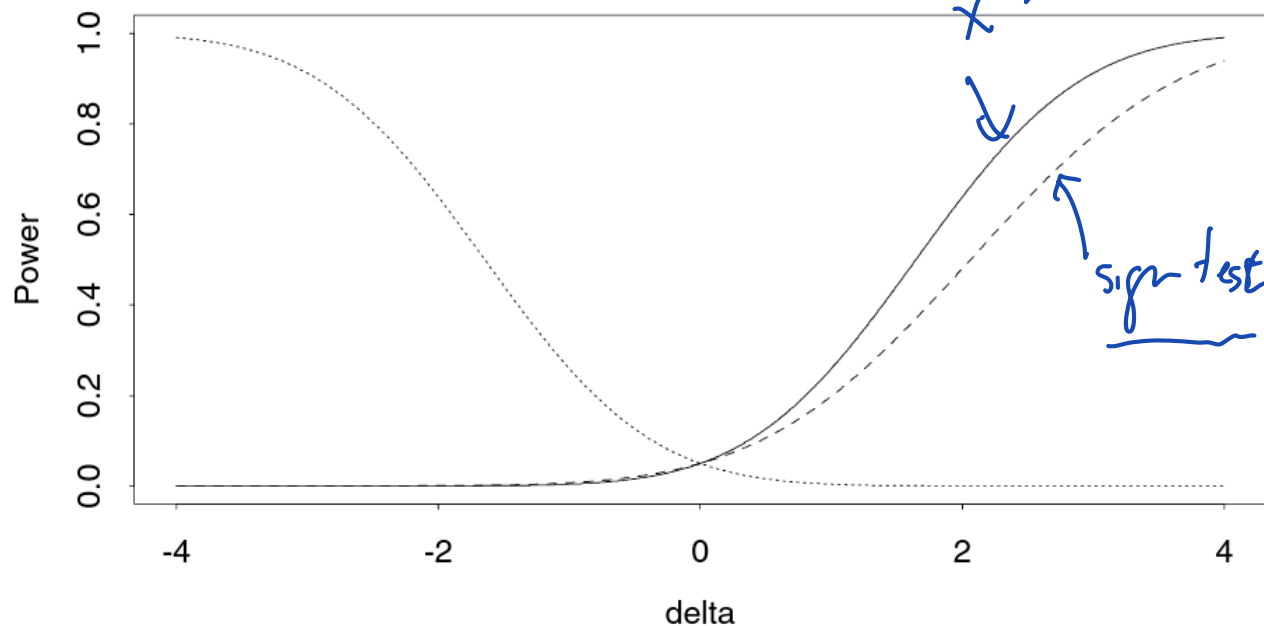
in probability

(c) Test statistic should have maximum power to detect departures from  $H_0$

334

7 · Estimation and Hypothesis Testing

VMP for 1-sided



**Figure 7.6** Power functions for a test of whether the mean of a  $N(\mu, \sigma^2)$  random sample of size  $n$  equals  $\mu_0$  against the alternative  $\mu = \mu_1$ , as a function of  $\delta = n^{1/2}(\mu_1 - \mu_0)/\sigma$ . The test size is  $\alpha = 0.05$ . The solid curve is the power function for a test of  $\mu_1 > \mu_0$  based on  $\bar{y}$ , and the dashed line is the power function for the sign test. Both critical regions are of form  $\bar{y} > t_\alpha$ . The dotted curve is the power function for  $\bar{y}$  when the critical region is  $\bar{y} < t_\alpha$ .

# Recap: Hypothesis tests and significance tests

- **Hypothesis tests** typically means:

- $H_0, H_1$
- critical/rejection region  $R \subset \mathcal{X}$ ,
- level  $\alpha$ , power  $1 - \beta$
- conclusion: “reject  $H_0$  at level  $\alpha$ ” or “do not reject  $H_0$  at level  $\alpha$ ”
- planning: maximize power for some relevant alternative

minimize type II error

- **Significance tests** typically means:

- $H_0$ ,
- test statistic  $T$
- observed value  $t^{obs}$ ,
- $p$ -value  $p^{obs} = \Pr(T \geq t^{obs}; H_0)$
- alternative hypothesis often only implicit

$H_1$  vague or unspecified

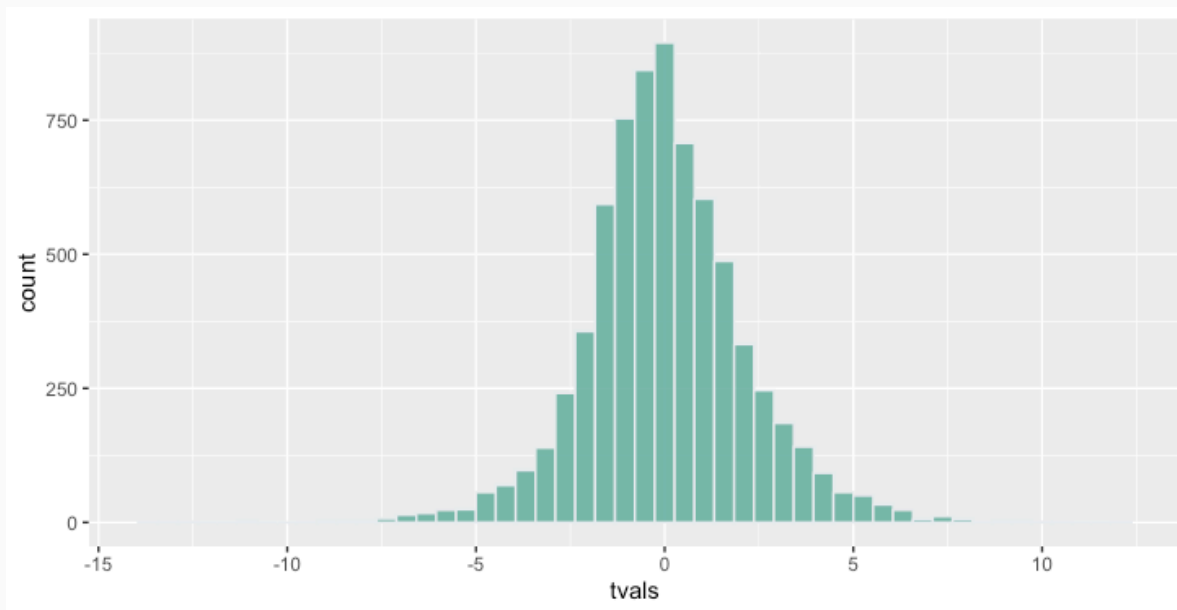
large  $T$  points to alternative



```
leukemia_big <- read.csv  
  ("http://web.stanford.edu/~hastie/CASI_files/DATA/leukemia_big.csv")  
dim(leukemia_big)  
[1] 7128    72
```

- each row is a different gene; 47 AML responses and 25 ALL responses
- we could compute 7128  $t$ -statistics for the mean difference between AML and ALL

```
tvals <- rep(0,7128)  
for (i in 1:7128){  
  leukemia_big[i,] %>% select(starts_with("ALL")) %>% as.numeric() -> x  
  leukemia_big[i,] %>% select(starts_with("AML")) %>% as.numeric() -> y  
  tvals[i] <- t.test(x,y,var.equal=T)$statistic  
}
```



```
summary(tvals)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
-13.52611	-1.20672	-0.08406	0.02308	1.20886	12.26065

$\underline{x}, \underline{y}$  permute

45  $\leftarrow$  Gp 1

27  $\leftarrow$  Gp 2

t-stat. for each gene

$H_0: F_x(x) = F_y(y)$  7,128

$H_1: \text{not } H_0$

- order the  $p$ -values  $p_{(1)}, \dots, p_{(m)}$
- find  $i_{max}$ , the largest index for which

$$p_{(i)} \leq \frac{i}{m} q$$

$\approx \alpha$

Bonferroni use  $\frac{\alpha}{m}$   
cutoff  
for 'sign'

- Let  $BH_q$  be the rule that rejects  $H_{0i}$  for  $i \leq i_{max}$ , not rejecting otherwise
- **Theorem:** If the  $p$ -values corresponding to valid null hypotheses are independent of each other, then

$$E \frac{\# \text{ rejects incorrect}}{\# \text{ rejects}} = FDR(BH_q) = \pi_0 q \leq q, \quad \text{where } \pi_0 = m_0/m$$

$\pi_0$  unknown but close to 1

- change the bound under dependence

$$p_{(i)} \leq \frac{i}{m C_m} q$$

$$C_m = \sum_{i=1}^m \frac{1}{i}$$

# Example

$$.05/10 = .005$$

AoS Ex.10.28

Bonferroni

index	1	2	3	4	5	6	7	8	9	10
pval	0.00017	0.00448	0.00671	0.00907	0.01220	0.33626	0.3934	0.5388	0.5813	0.9862
B-H cut1	0.00500	0.01000	0.01500	0.02000	0.02500	0.03000	0.0350	0.0400	0.0450	0.0500
cut2	0.00171	0.00341	0.00512	0.00683	0.00854	0.01024	0.0119	0.0137	0.0154	0.0171

dependence

to composite

$$H_0: F_x = F_y$$

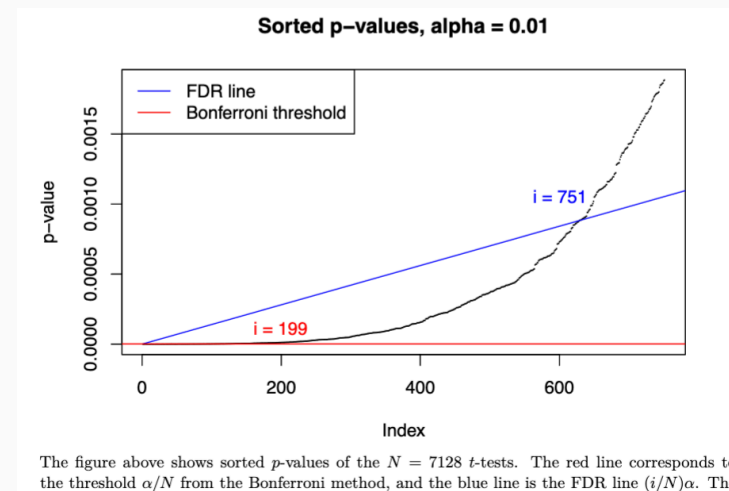
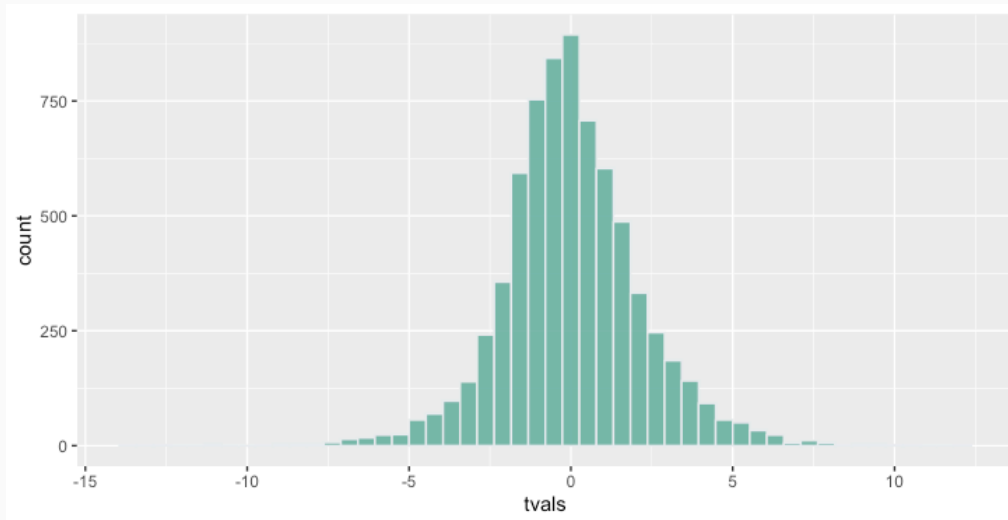
then all  $\binom{72}{27}$

permutation of data have = prob.

convert comp. H<sub>0</sub> to simple (by cond'g)

$H_0: \checkmark$  sample composite  
 $H_1: \checkmark$  completely spec. dist.

← p values are often easy to compute



```
> summary(ttest)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
-13.52611	-1.20672	-0.08406	0.02308	1.20886	12.26065

$X_1, \dots, X_{10000}$

$\hat{\beta}_5(X_5)$

"large" relative to  
1st. s.e.

(Lasso)  $\hat{\beta}_5 > 0$ , for  $\lambda_{opt}$

- $X_1, \dots, X_n$  i.i.d.
- $H_0 : X_i \sim f(x; \theta); \quad H_1 : X_i \text{ arbitrary distribution}$
- Define  $k$  sets  $A_1, \dots, A_k$  s.t.

*composite*

- Define

$$\text{pr}(X_i \in \cup_{j=1}^k A_j) = 1$$

$$Y_j = \sum_{i=1}^n 1\{X_i \in A_j\}$$

$\bigcup A_j = \mathcal{X}$   
sample  
sp

number of obs in category  $j$

- $X_1, \dots, X_n$  i.i.d.
- $H_0 : X_i \sim f(x; \theta); H_1 : X_i$  arbitrary distribution
- Define  $k$  sets  $A_1, \dots, A_k$  s.t.

$$p_j(\theta) = \int_{\{x \in A_j\}} f(x; \theta) dx$$

$$\text{pr}(X_i \in \cup_{j=1}^k A_j) = 1$$

- Define

$$Y_j = \sum_{i=1}^n 1\{X_i \in A_j\}$$

number of obs in category  $j$

- $Y = (Y_1, \dots, Y_k) \sim \text{Mult}_k(n; p)$

$$p_j = \text{Pr}(Y_i \in A_j)$$

- $\text{pr}(Y_1 = y_1, \dots, Y_k = y_k; p) =$

$$f(y; p) = \frac{n!}{y_1! \dots y_k!} p_1^{y_1} \dots p_k^{y_k}, \quad \underline{\underline{\sum y_j = n}}$$

- $H_0 : p = p(\theta); H_1 : p$  arbitrary

$$\sum p_j = 1, \quad 0 \leq p_j \leq 1$$

- log-likelihood function

$$l(\mathbf{p}) = \sum_{j=1}^k y_j \log p_j \quad j=1, \dots, k$$

- generalized likelihood ratio test

$$\max_{\mathbf{p}} l(\mathbf{p}) \text{ over } H_0 \cup H_1$$

$$\max_{\mathbf{p}} l(\mathbf{p}) \text{ over } H_0$$

$$\frac{\sup_{\mathbf{p}} L(\mathbf{p}; \mathbf{y})}{\sup_{\mathbf{p} \in \Theta_0} L(\mathbf{p}; \mathbf{y})}$$

$\swarrow H_1$   
 $\uparrow H_0$

$$W = 2 \left\{ \underbrace{\sup_{\mathbf{p}} l(\mathbf{p}; \mathbf{y})}_{\text{over } H_1 \cup H_0} - \sup_{\mathbf{p} \in \Theta_0} l(\mathbf{p}; \mathbf{y}) \right\}$$

$$\hat{\mathbf{p}} \text{ under no constraint} = \frac{\mathbf{y}}{n} \Rightarrow \hat{p}_j = \frac{y_j}{n}$$



- log-likelihood function

$$y \sim \text{Mult}_k(n, p) = \frac{n!}{y_1! \dots y_k!} \pi_1^{y_1} \dots \pi_k^{y_k}$$

- generalized likelihood ratio test

$$= \sum_{\text{bins}} \left( O \log \frac{O}{E} \right)$$

- Theorem 9.1 (MS): Under  $H_0$

$$p = \dim(\theta)$$

$$W = 2 \sum_{j=1}^k Y_j \log \left( \frac{Y_j}{np_j(\tilde{\theta})} \right) \xrightarrow{d} \chi_{k-1-p}^2$$

usual LR theory

$$= 2 \sum Y_j \log \frac{Y_j}{n} - 2 \sum Y_j \log p_j(\tilde{\theta})$$

$$L(p) = \prod_{j=1}^k p_j^{y_j} \quad \sum p_j = 1, \sum y_j = n$$

$$\begin{aligned} l(p) &= \sum_{j=1}^k y_j \log p_j, \quad \sum p_j = 1, \sum y_j = n \\ &= \sum_{j=1}^{k-1} \{y_j \log p_j\} + \underbrace{(n - y_1 - \dots - y_{k-1})}_{\substack{\uparrow \\ y_k}} \log(1 - p_1 - \dots - p_{k-1}) \end{aligned}$$

$$\begin{aligned} \frac{\partial l}{\partial p_j} &= \frac{y_j}{p_j} - \frac{n - y_1 - \dots - y_{k-1}}{(1 - p_1 - \dots - p_{k-1})}, \quad j = 1, \dots, k-1 \\ &= \frac{y_j}{p_j} - \frac{y_k}{p_k} \quad j = 1, \dots, k-1 \\ \int_{\hat{p}}^n = 0 &\implies \dots \frac{y_j}{n_j} = \hat{p}_j \end{aligned}$$

$$l(p(\theta)) = \sum_{j=1}^{k-1} y_j \log p_j(\theta) + (n - \sum_{j=1}^{k-1} y_j) \log \left( 1 - \sum_{j=1}^{k-1} p_j(\theta) \right)$$

$$\frac{\partial l}{\partial \theta_j} : \sum_{j=1}^k y_j \frac{p_j'(\theta)}{p_j(\theta)} \neq ; \quad \begin{aligned} \sum p_j(\theta) &= 1 \\ \sum y_j &= n \end{aligned}$$

$$\begin{aligned} \Theta_{k \times 1} \quad \frac{\partial l}{\partial \theta} \bigg|_{\hat{\theta}} &= 0 \quad \tilde{\theta} \text{ vector of length } k \\ \text{vector} & \quad ; \end{aligned}$$

- log-likelihood function

$$\log \frac{Y_j}{np_j(\tilde{\theta})} = \log \left( 1 + \underbrace{\frac{Y_j}{np_j(\tilde{\theta})} - 1}_{\varepsilon} \right)$$

- generalized likelihood ratio test

$$\log(1 + \varepsilon) \approx \varepsilon - \frac{1}{2}\varepsilon^2 + \frac{1}{3}\varepsilon^3 - \frac{1}{4}\varepsilon^4 + \dots$$

$p = \dim(\theta)$

- Theorem 9.1 (MS): Under  $H_0$

$$W = 2 \sum_{j=1}^k Y_j \log \left( \frac{Y_j}{np_j(\tilde{\theta})} \right) \xrightarrow{d} \chi_{k-1-p}^2$$

- Theorem 9.2. (MS): Under  $H_0$

$$Q = \sum_{j=1}^k \frac{\{Y_j - np_j(\hat{\theta})\}^2}{np_j(\hat{\theta})} \xrightarrow{d} \chi_{k-1-p}^2$$

dist of  $W$   
~~dist of  $Q$~~   
 for large  $n$

# Multinomial goodness-of-fit tests

AoS 10.8; MS 9.2

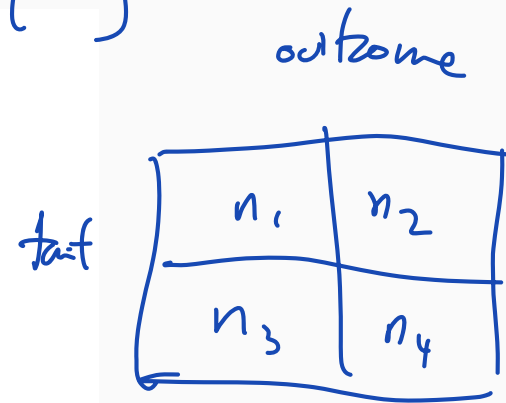
1, 2, 3 "Wald"

$$(\hat{\theta} - \theta)^T I(\hat{\theta}) (\hat{\theta} - \theta)$$

$$(p(\tilde{\theta}) - p) \underline{I(p(\tilde{\theta}))} (p(\tilde{\theta}) - p)$$

Table 9.1 Frequency of goals in First Division matches and "expected" frequency under Poisson model in Example 9.2

Goals	0	1	2	3	4	$\geq 5$
Frequency	252	344	180	104	28	16
Expected	248.9	326.5	214.1	93.6	30.7	10.2



$$p_0(\lambda) = 1 - \sum_{j=0}^4 p_j(\lambda); \quad p_j(\lambda) = e^{-\lambda} \lambda^j / j!, \quad \tilde{\lambda} = 1.3118$$

$$Q = 11.09; \quad W = 10.87; \quad \text{pr}(\chi_4^2 > [11.09, 10.87]) = [0.026, 0.028]$$

136

4 · Likelihood

		Antigen 'B'		Total
		Absent	Present	
Antigen 'A'	Absent	'O': 202	'B': 35	237
	Present	'A': 179	'AB': 6	185
Total		381	41	422

**Table 4.3** Blood groups in England (Taylor and Prior, 1938). The upper part of the table shows a cross-classification of 422 persons by presence or absence of antigens 'A' and 'B', giving the groups 'A', 'B', 'AB', 'O' of the human blood group system. The lower part shows genotypes and corresponding probabilities under one- and two-locus models. See Example 4.38 for details.

log-linear models for cross-classified

Group	Two-locus model		One-locus model	
	Genotype	Probability	Genotype	Probability
'A'	(AA; bb), (Aa; bb)	$\alpha(1 - \beta)$	(AA), (AO)	$\lambda_A^2 + 2\lambda_A\lambda_O$
'B'	(aa; BB), (aa; Bb)	$(1 - \alpha)\beta$	(BB), (BO)	$\lambda_B^2 + 2\lambda_B\lambda_O$
'AB'	(AA; BB), (Aa; BB), (AA; Bb), (Aa; Bb)	$\alpha\beta$	(AB)	$2\lambda_A\lambda_B$
'O'	(aa; bb)	$(1 - \alpha)(1 - \beta)$	(OO)	$\lambda_O^2$

$$Q = 15.73; W = 17.66 \text{ (two-locus)}$$

$$p < 10^{-5}$$

$$Q = 2.82; W = 3.17 \text{ (single locus)}$$

$$p = 0.09; 0.07$$

- $X_1, \dots, X_n$  i.i.d.  $F(\cdot)$ ;  $H_0 : F = F_0$  — simple e.g.  $U(0,1)$
- $\hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n 1\{X_i \leq t\}$  — comp e.g.  $N(\mu, \sigma^2) \leftarrow \hat{\mu}, \hat{\sigma}^2$
- three test statistics:
  - 1.  $\sup_t |\hat{F}_n(t) - F_0(t)|$  Kolmogorov-Smirnov
  - 2.  $\int \{\hat{F}_n(t) - F_0(t)\}^2 dF_0(t)$  Cramer-vonMises
  - 3.  $\int \frac{\{\hat{F}_n(t) - F_0(t)\}^2}{F_0(t)\{1 - F_0(t)\}} dF_0(t)$  Anderson-Darling
- SM Example 7.24 testing  $N(\mu, \sigma^2)$  distribution
- SM Example 7.23; 6.14 testing  $U(0, 1)$  distribution

cumulative d.f.

estimated

- Special case  $H_0 : F(t) = F_0(t) = t$
- Recall

$$E_0\{\hat{F}_n(t)\} = F_0(t) = t,$$

$$\text{var}\{\hat{F}_n(t)\} = t(1-t)/n$$

$$X_i \sim U(0, 1)$$

fixed  $t$

- What about distribution of

$$\sup_t |\hat{F}_n(t) - t|$$

$$\int \{\hat{F}_n(t) - t\}^2 dt$$

$$\int \frac{\{\hat{F}_n(t) - t\}^2}{F_0(t)\{1 - F_0(t)\}} dt$$

- need joint density of  $\hat{F}_n(t) \forall t$

- Special case  $H_0 : F(t) = F_0(t) = t$
- Recall

$$X_i \sim U(0, 1)$$

$$E_0\{\hat{F}_n(t)\} = F_0(t) = t, \quad \text{var}\{\hat{F}_n(t)\} = t(1-t)/n$$

- What about distribution of

$$\sup_t |\hat{F}_n(t) - t| \quad \int \{\hat{F}_n(t) - t\}^2 dt \quad \int \frac{\{\hat{F}_n(t) - t\}^2}{F_0(t)\{1 - F_0(t)\}} dt$$

- need joint density of  $\hat{F}_n(t) \forall t$

- define **stochastic process**

$$B_n(t) = \sqrt{n}(\hat{F}_n(t) - t)$$

0 at  $t=0$   
1 at  $t=1$

$$\{B_n(t) : t \geq 0\}$$

$$B_n(0) = 0 = B_n(1)$$

$$\rightarrow \text{vector } (B_n(t_1), \dots, B_n(t_k)) \xrightarrow{d} N_k(\mathbf{0}, \mathbf{C}), \quad C_{ij} = \min(t_i, t_j) - t_i t_j$$

MS 9.3

- a **Brownian bridge** is a continuous function on  $(0, 1)$

with all finite-dimensional distributions as above



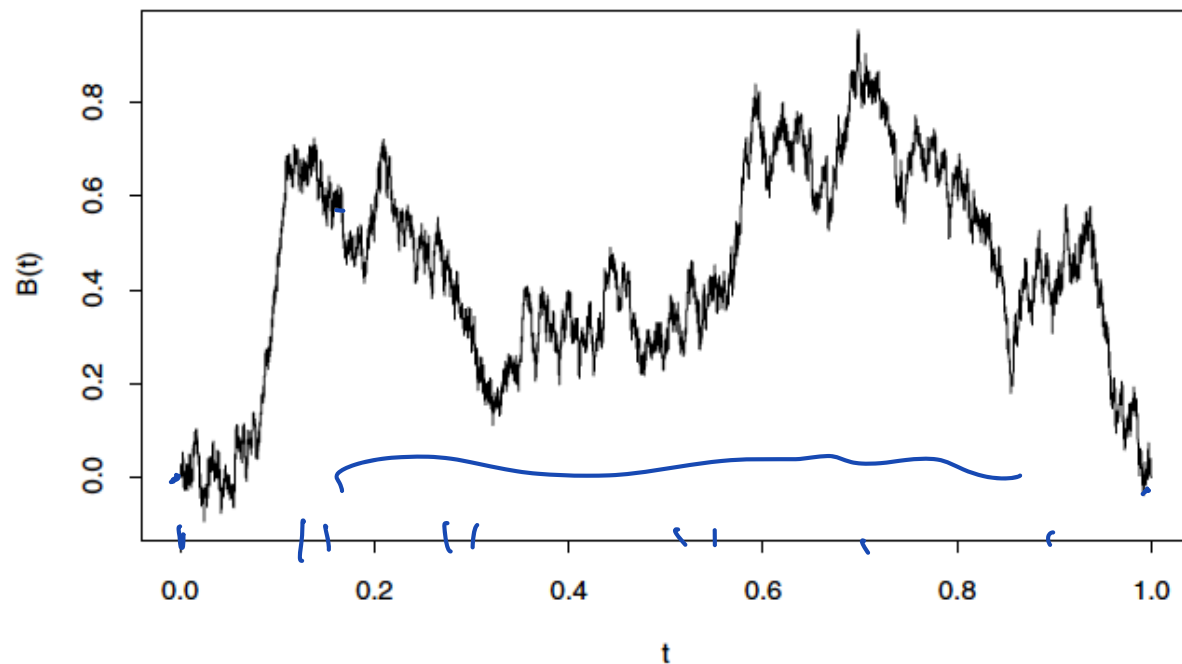


Figure 9.1 *A simulated realization of a Brownian bridge process.*

- Kolmogorov-Smirnov test
- Cramer-vonMises test
- Anderson-Darling test

$$K_n = \sup_{0 \leq t \leq 1} |B_n(t)|$$

$$W_n^2 = \int_0^1 B_n^2(t) dt$$

$$A_n^2 = \int_0^1 \frac{B_n^2(t)}{t(1-t)} dt$$

from  $\hat{F}_n(t) - F_0(t)$   
to  $B_n(t)$

- Kolmogorov-Smirnov test

$$K_n = \sup_{0 \leq t \leq 1} |B_n(t)|$$

very conservative

- Cramer-vonMises test

$$W_n^2 = \int_0^1 B_n^2(t) dt$$

- Anderson-Darling test

$$A_n^2 = \int_0^1 \frac{B_n^2(t)}{t(1-t)} dt$$

- limit theorems

$$\underbrace{K_n \xrightarrow{d} K,}$$

$$W_n^2 \xrightarrow{d} \sum_{j=1}^{\infty} \frac{Z_j^2}{j^2 \pi^2},$$

$$A_n^2 \xrightarrow{d} \sum_{j=1}^{\infty} \frac{Z_j^2}{j(j+1)}$$

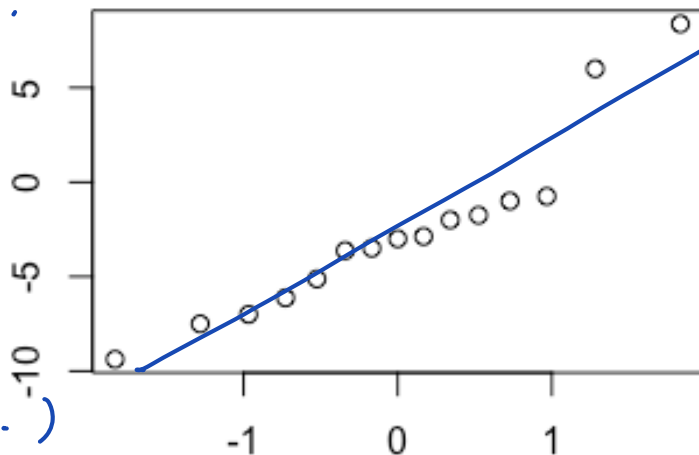
$$\text{pr}(K > x) = 2 \sum_{j=1}^{\infty} (-1)^{j+1} \exp(-2j^2 x^2)$$

Maize data SM Ex 7.24

obs.

Sample Quantiles

$F_n(\cdot)$



Theoretical Quantiles



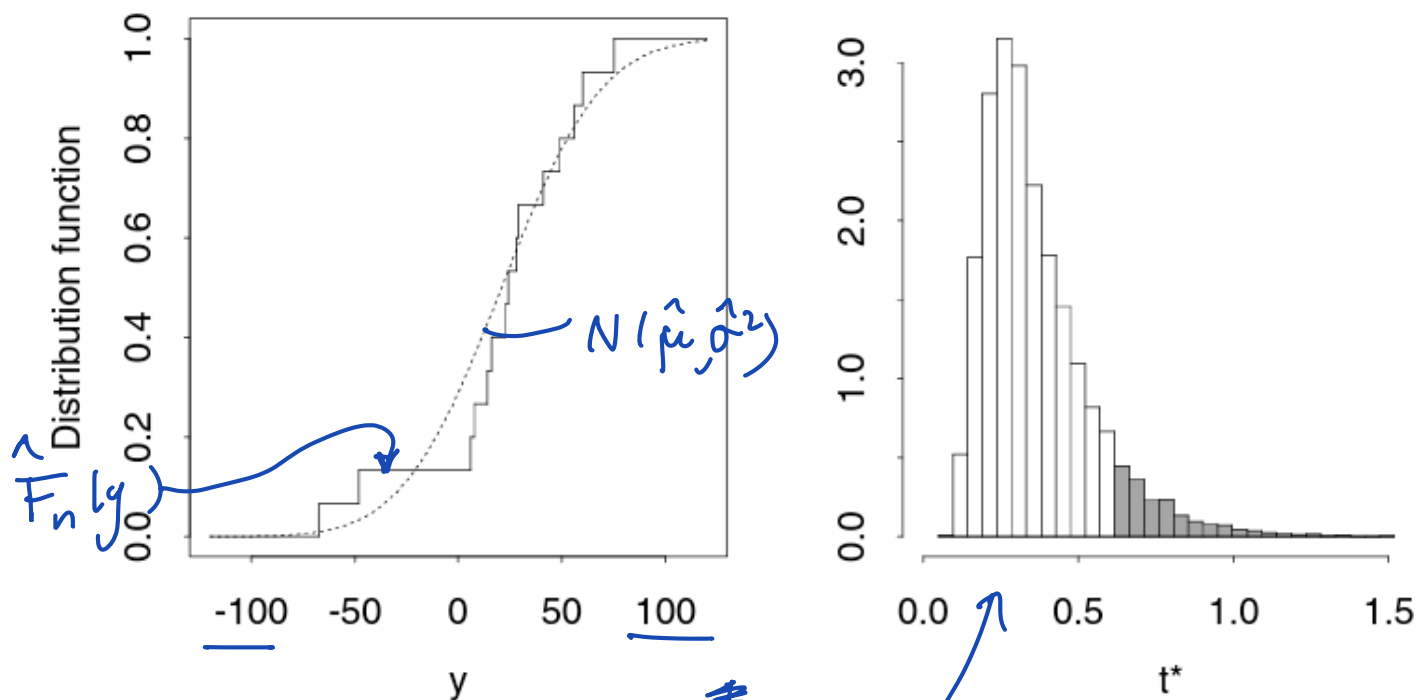
Exp<sub>N</sub> values of  $X_{(i)}$

```
library(SMPracticals)
data(darwin)
cross <- seq(1,30,by=2)
self <- cross+1
diffs <- darwin[self,4]-darwin[cross,4]
qqnorm(diffs)
```

e.g. test st. ?

## Example: SM 7.24

**Figure 7.5** Analysis of maize data. Left: empirical distribution function for height differences, with fitted normal distribution (dots). Right: null density of Anderson–Darling statistic  $T$  for normal samples of size  $n = 15$  with location and scale estimated. The shaded part of the histogram shows values of  $T^*$  in excess of the observed value  $t_{\text{obs}}$ .



SM Example 7.24 testing  $N(\mu, \sigma^2)$  distribution

$$\int_{-\infty}^{\infty} \frac{(\hat{F}_n(t) - \hat{F}_0(t))^2}{\hat{F}_0(t)(1 - \hat{F}_0(t))} d\hat{F}_0(t)$$

- Relatively simple case:  $\mathbf{X} \sim f(\mathbf{x}; \theta)$ ,  $H_0 : \theta = \theta_0$ ,  $H_1 : \theta \neq \theta_0$

$$\text{pr}(H_0 | \mathbf{x}) = \frac{f(\mathbf{x} | H_0) \text{pr}(H_0)}{f(\mathbf{x} | H_0) \text{pr}(H_0) + f(\mathbf{x} | H_1) \text{pr}(H_1)}$$

post. prob

$$\pi(\theta_0 | \underline{x}) = \frac{L(\theta_0; \underline{x}) \pi(\theta_0)}{L(\theta_0; \underline{x}) \pi(\theta_0) + \int f(\underline{x}; \theta) \pi(\theta) d\theta \cdot \pi(H_1)}$$

depends strongly  
on prior

some version of model  
selection

$f(\underline{x})$  marg'l

- Relatively simple case:  $\mathbf{X} \sim f(\mathbf{x}; \theta)$ ,  $H_0 : \theta = \theta_0$ ,  $H_1 : \theta \neq \theta_0$

•

$$\text{pr}(H_0 \mid \mathbf{x}) = \frac{f(\mathbf{x} \mid H_0)\text{pr}(H_0)}{f(\mathbf{x} \mid H_0)\text{pr}(H_0) + f(\mathbf{x} \mid H_1)\text{pr}(H_1)}$$

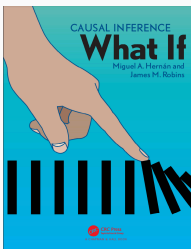
$$= \frac{f(\mathbf{x} \mid \theta_0)\text{pr}(H_0)}{f(\mathbf{x} \mid \theta_0)\text{pr}(H_0) + \int f(\mathbf{x} \mid \theta)\pi(\theta)d\theta\text{pr}(H_1)}$$

$$= \frac{L_n(\theta_0)}{L_n(\theta_0) + \int L_n(\theta)\pi(\theta)d\theta} \quad \left. \vphantom{\frac{L_n(\theta_0)}{L_n(\theta_0) + \int L_n(\theta)\pi(\theta)d\theta}} \right\} \mathcal{P}_n(H_0) = \mathcal{P}_n(H_1) = \frac{1}{2}$$

- can't use improper priors; result is sensitive to the prior for  $\theta$

- randomization; confounding; observational studies; experiments; “correlation is not causation”, Simpson’s ‘paradox’
- counterfactuals; average treatment effect; conditional average treatment effect; ...
- graphical models; directed acyclic graphs; causal graphs; Markov assumptions...

- The Book



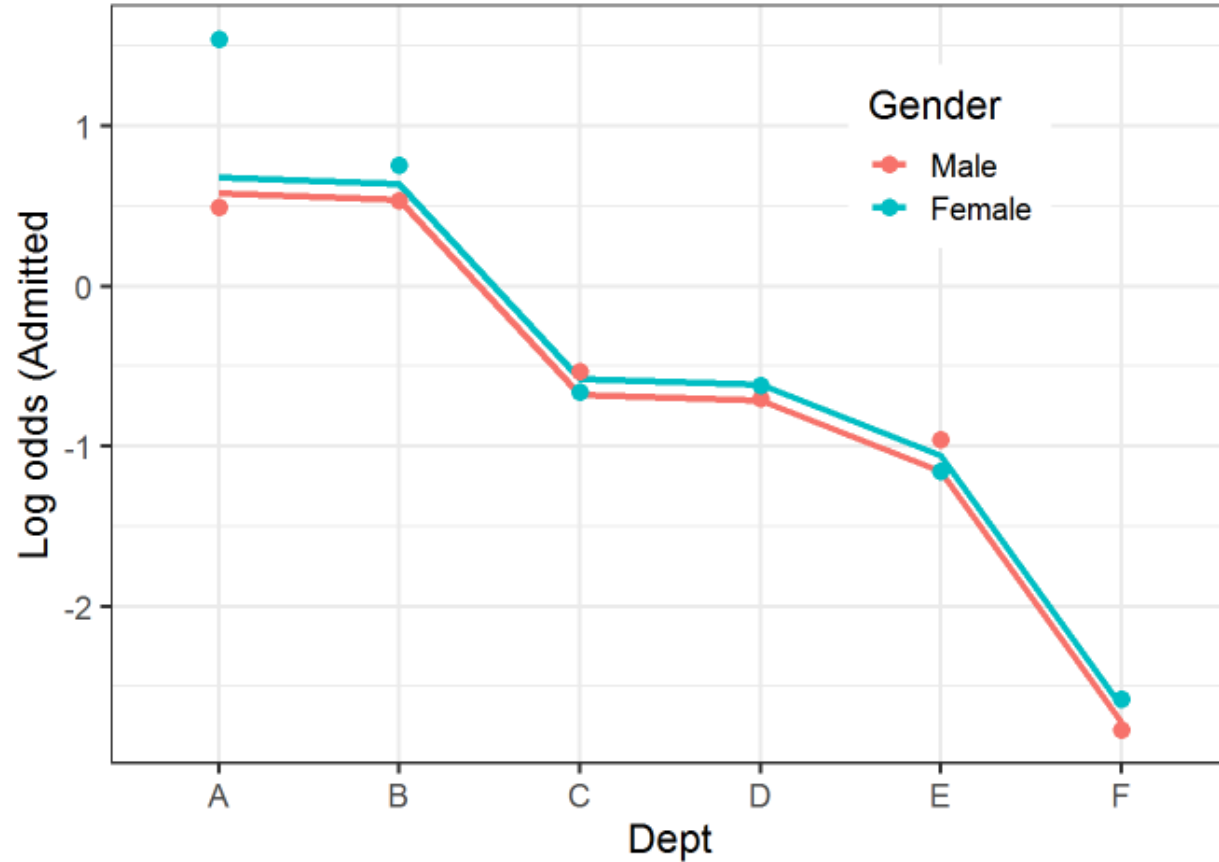


# Confounding variables

Major	Number of applicants	Men		Women		
		Number admitted	Percent admitted	Number of applicants	Number admitted	Percent admitted
A	825	512	62	108	89	82
B	560	353	63	25	17	68
C	325	120	37	593	202	34
D	417	138	33	375	131	35
E	191	53	28	393	94	24
F	373	22	6	341	24	7
Total	2691	1198	44	1835	557	30

data(UCBAdmissions)

## ... Confounding variables



race of defendant	death penalty imposed	death penalty not imposed	percentage
white	19	141	11.88%
black	17	149	10.24%

race of defendant	death penalty imposed	death penalty not imposed	percentage
white	19	141	11.88%
black	17	149	10.24%

white victim	race of defendant	death penalty imposed	death penalty not imposed	percentage
	white	19	132	12.58%
	black	11	52	17.46%

black victim	race of defendant	death penalty imposed	death penalty not imposed	percentage
	white	0	9	0%
	black	6	97	5.83%

258

6 · Stochastic Models

Age (years)	Smokers	Non-smokers
Overall	139/582 (24)	230/732 (31)
18–24	2/55 (4)	1/62 (2)
25–34	3/124 (2)	5/157 (3)
35–44	14/109 (13)	7/121 (6)
45–54	27/130 (21)	12/78 (15)
55–64	51/115 (44)	40/121 (33)
65–74	29/36 (81)	101/129 (78)
75+	13/13 (100)	64/64 (100)

**Table 6.8** Twenty-year survival and smoking status for 1314 women (Appleton *et al.*, 1996). The smoker and non-smoker columns contain number dead/total (% dead).

- $X$  – binary treatment indicator
- $Y$  – binary outcome
- “ $X$  **causes**  $Y$ ” to be distinguished from “ $X$  is associated with  $Y$ ”

“treatment”

could be continuous

- $X$  – binary treatment indicator
- $Y$  – binary outcome
- “ $X$  **causes**  $Y$ ” to be distinguished from “ $X$  is associated with  $Y$ ”

“treatment”  
could be continuous

- introduce **potential outcomes**  $C_0, C_1$

$$Y = \begin{cases} C_0 & \text{if } X = 0 \\ C_1 & \text{if } X = 1 \end{cases}$$

- equivalently  $Y = C_X$  or  $Y = C_0(1 - X) + C_1X$

consistency equation

- **causal treatment effect**       $\theta = E(C_1) - E(C_0)$

want to estimate this

- **association**

$$\alpha = E(Y \mid X = 1) - E(Y \mid X = 0)$$

have data to estimate  $\alpha$

Potential outcomes  $C_0, C_1$

$X$	$Y$	$C_0$	$C_1$
0	4	4	*
0	7	7	*
0	2	2	*
0	8	8	*
1	3	*	3
1	5	*	5
1	8	*	8
1	9	*	9

treatment  $X$ , response  $Y$

Potential outcomes  $Y^0, Y^1$

Table 2.1

	$A$	$Y$	$Y^0$	$Y^1$
Rheia	0	0	0	?
Kronos	0	1	1	?
Demeter	0	0	0	?
Hades	0	0	0	?
Hestia	1	0	?	0
Poseidon	1	0	?	0
Hera	1	0	?	0
Zeus	1	1	?	1
Artemis	0	1	1	?
Apollo	0	1	1	?
Leto	0	0	0	?
Ares	1	1	?	1
Athena	1	1	?	1
Hephaestus	1	1	?	1
Aphrodite	1	1	?	1
Cyclope	1	1	?	1
Persephone	1	1	?	1
Hermes	1	0	?	0
Hebe	1	0	?	0
Dionysus	1	0	?	0



## Potential outcomes

Table 1.1

	$Y^{a=0}$	$Y^{a=1}$
Rheia	0	1
Kronos	1	0
Demeter	0	0
Hades	0	0
Hestia	0	0
Poseidon	1	0
Hera	0	0
Zeus	0	1
Artemis	1	1
Apollo	1	0
Leto	0	1
Ares	1	1
Athena	1	1
Hephaestus	0	1
Aphrodite	0	1
Cyclope	0	1
Persephone	1	1
Hermes	1	0
Hebe	1	0
Dionysus	1	0

## Observed outcomes

Table 1.2

	$A$	$Y$
Rheia	0	0
Kronos	0	1
Demeter	0	0
Hades	0	0
Hestia	1	0
Poseidon	1	0
Hera	1	0
Zeus	1	1
Artemis	0	1
Apollo	0	1
Leto	0	0
Ares	1	1
Athena	1	1
Hephaestus	1	1
Aphrodite	1	1
Cyclope	1	1
Persephone	1	1
Hermes	1	0
Hebe	1	0
Dionysus	1	0

$$\theta = E(C_1) - E(C_0)$$

risk difference; ratio; odds

$$\alpha = E(Y \mid X = 1) - E(Y \mid X = 0)$$

If  $X$  is independent of  $(C_0, C_1)$ ,  $\theta = \alpha$

If  $X$  is randomly assigned, then  $X \perp (C_0, C_1)$

## Example 16.2

$X$	$Y$	$C_0$	$C_1$
0	0	0	0*
0	0	0	0*
0	0	0	0*
0	0	0	0*
1	1	1*	1
1	1	1*	1
1	1	1*	1
1	1	1*	1

$$\theta = 0; \quad \alpha = 1$$

$(C_0, C_1)$  not independent of  $X$

$X$	$Y$	$C_0$	$C_1$
0	0	0	0*
1	0	0	0*
1	0	0	0*
1	0	0	0*
1	1	1*	1
1	1	1*	1
1	1	1*	1
1	1	1*	1

$$\theta = 0, \quad \alpha = 4/7 < 1$$

thought experiment

## Potential outcomes

Table 1.1

	$Y^{a=0}$	$Y^{a=1}$
Rheia	0	1
Kronos	1	0
Demeter	0	0
Hades	0	0
Hestia	0	0
Poseidon	1	0
Hera	0	0
Zeus	0	1
Artemis	1	1
Apollo	1	0
Leto	0	1
Ares	1	1
Athena	1	1
Hephaestus	0	1
Aphrodite	0	1
Cyclope	0	1
Persephone	1	1
Hermes	1	0
Hebe	1	0
Dionysus	1	0

## Observed outcomes

Table 1.2

	$A$	$Y$
Rheia	0	0
Kronos	0	1
Demeter	0	0
Hades	0	0
Hestia	1	0
Poseidon	1	0
Hera	1	0
Zeus	1	1
Artemis	0	1
Apollo	0	1
Leto	0	0
Ares	1	1
Athena	1	1
Hephaestus	1	1
Aphrodite	1	1
Cyclope	1	1
Persephone	1	1
Hermes	1	0
Hebe	1	0
Dionysus	1	0

1. A well-understood evidence-based mechanism, or set of mechanisms, that links a cause to its effect
2. two phenomena are linked by a stable association, whose direction is established and which cannot be explained by mutual dependence on some other allowable variable
3. observed association may be linked to causal effect via counterfactuals if  
 $(C_o, C_o) \perp X$  not usually testable

- typically have additional explanatory variables (covariates)  $Z$
- causal effect of treatment when  $Z = z$

$$\theta_z = E(C_1 \mid Z = z) - E(C_0 \mid Z = z)$$

- marginal causal effect

$$\theta = E_Z\{E(C_1 \mid Z) - E(C_0 \mid Z)\}$$

Table 2.2

	$L$	$A$	$Y$
Rheia	0	0	0
Kronos	0	0	1
Demeter	0	0	0
Hades	0	0	0
Hestia	0	1	0
Poseidon	0	1	0
Hera	0	1	0
Zeus	0	1	1
Artemis	1	0	1
Apollo	1	0	1
Leto	1	0	0
Ares	1	1	1
Athena	1	1	1
Hephaestus	1	1	1
Aphrodite	1	1	1
Cyclope	1	1	1
Persephone	1	1	1
Hermes	1	1	0
Hebe	1	1	0
Dionysus	1	1	0

$$\theta_{L=0}$$

$$\theta_{L=1}$$

$L = 1$  critical condition

$L = 0$  stable condition  
conditional randomization

- continuous “treatment” variable  $X \in \mathbb{R}$
- counterfactual outcome  $(C_0, C_1) \rightarrow$  counterfactual function  $C(x)$
- observed response  $Y = C(X)$  consistency
- causal regression function  
 $\theta(x) = E\{C(x)\}$
- association regression function  
 $r(x) = E(Y | X)$

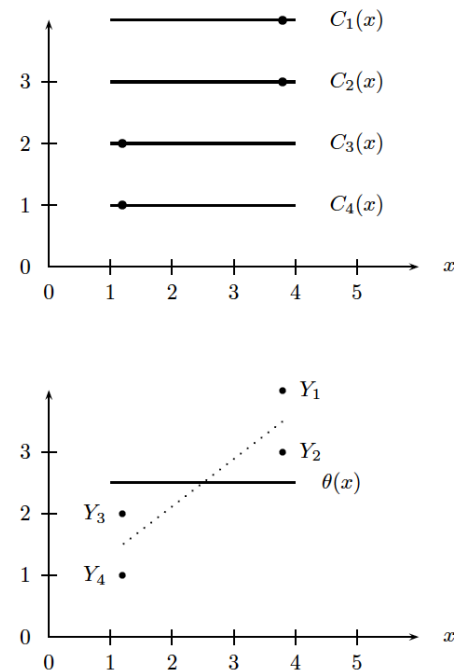


FIGURE 16.2. The top plot shows the counterfactual function  $C(x)$  for four subjects. The dots represent their  $X$  values. Since  $C_i(x)$  is constant over  $x$  for all  $i$ , there is no causal effect. Changing the dose will not change anyone's outcome. The lower plot shows the causal regression function  $\theta(x) = (C_1(x) + C_2(x) + C_3(x) + C_4(x))/4$ . The four dots represent the observed data points  $Y_1 = C_1(X_1)$ ,  $Y_2 = C_2(X_2)$ ,  $Y_3 = C_3(X_3)$ ,  $Y_4 = C_4(X_4)$ . The dotted line represents the regression  $r(x) = E(Y|X = x)$ . There is no causal effect since  $C_i(x)$  is constant for all  $i$ . But there is an association since the regression curve  $r(x)$  is not constant.



- in observational studies treatment is not randomly assigned  $\implies \theta(x) \neq r(x)$
- group subjects based on additional **confounding** variables
- **No unmeasured confounding:**

$$\{C(x); x \in \mathcal{X}\} \perp X \mid Z$$

- under the assumption of no unmeasured confounding,  
the causal regression function

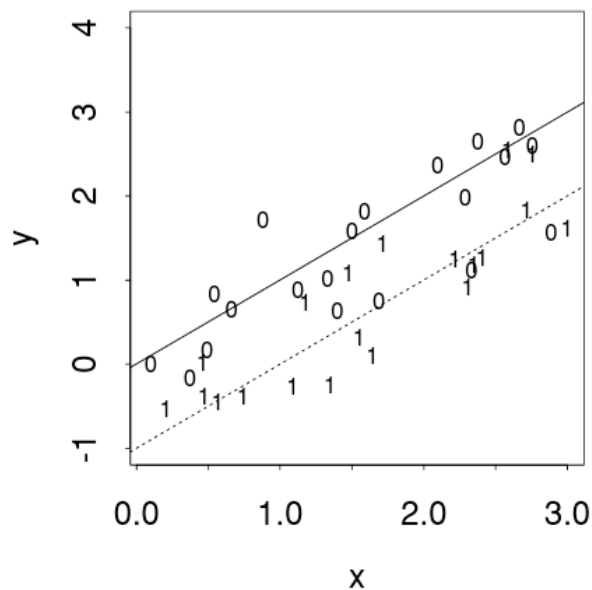
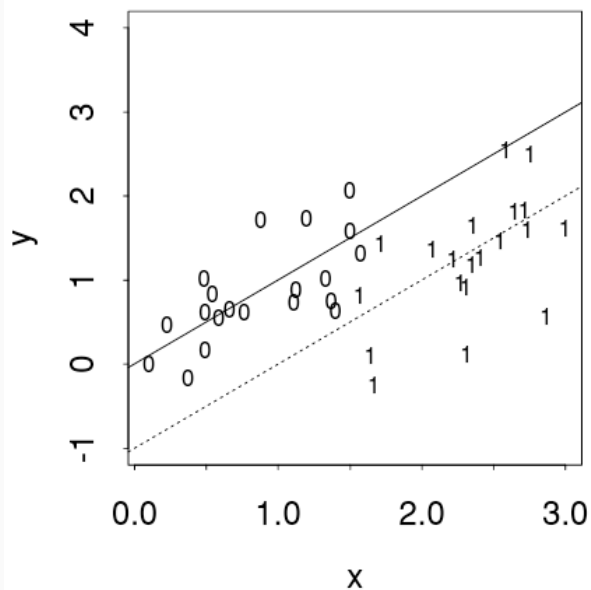
typo in (16.7)

$$\theta(x) = \int \mathbb{E}(Y \mid X = x, Z = z) dF_Z(z)$$

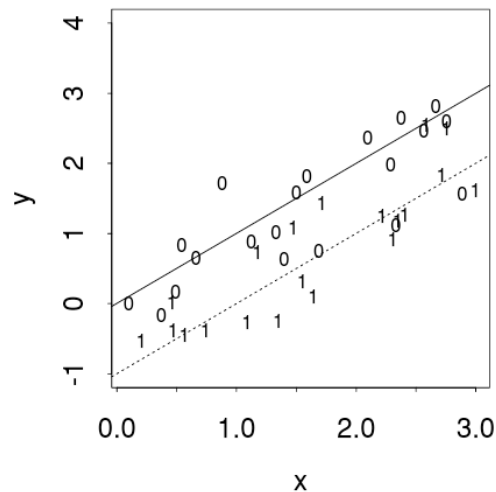
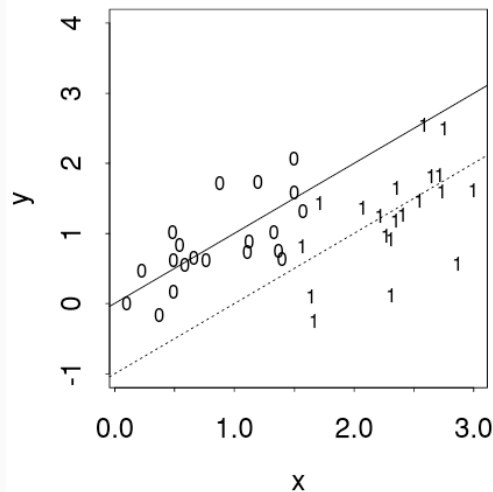
can be estimated by the association function

$$\hat{\theta}(x) = \frac{1}{n} \sum_{i=1}^n \hat{r}(x, Z_i) = \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 \bar{Z}_n$$

causal reg function  $\equiv$  adjusted treatment effect



**Figure 9.2** Simulated results from experiments to compare the effect of a treatment  $T$  on a response  $Y$  that varies with a covariate  $X$ . The lines show the mean response for  $T = 0$  (solid) and  $T = 1$  (dotted). Left: the effect of  $T$  is confounded with dependence on  $X$ . Right: the experiment is balanced, with random allocation of  $T$  dependent on  $X$ .



**Figure 9.2** Simulated results from experiments to compare the effect of a treatment  $T$  on a response  $Y$  that varies with a covariate  $X$ . The lines show the mean response for  $T = 0$  (solid) and  $T = 1$  (dotted). Left: the effect of  $T$  is confounded with dependence on  $X$ . Right: the experiment is balanced, with random allocation of  $T$  dependent on  $X$ .

$$C_1(x) - C_0(x) \equiv 1$$

Left:  $\bar{y}_1 - \bar{y}_0 = 0.2 \pm 0.3$

Right:  $\bar{y}_1 - \bar{y}_0 = -1.2 \pm 0.3$

adjust for covariate:  $y = \beta_0 + \beta_1 x + \delta t + \epsilon$

Left:  $\hat{\delta} = -0.7 \pm 0.3$  Right:  $\hat{\delta} = -1.25 \pm 0.16$

right randomized within pairs; matched on  $x$

**“Bradford-Hill guidelines”** Evidence that an observed association is causal is strengthened if:

- the association is strong
- the association is found consistently over a number of independent studies
- the association is specific to the outcome studied
- the observation of a potential cause occurs earlier in time than the outcome
- there is a dose-response relationship
- there is subject-matter theory that makes a causal effect plausible
- the association is based on a suitable natural experiment

see also AoS §16.3

260 16. Causal Inference

	$Y = 1$	$Y = 0$	$Y = 1$	$Y = 0$
$X = 1$	.1500	.2250	.1000	.0250
$X = 0$	.0375	.0875	.2625	.1125
	$Z = 1$ (men)		$Z = 0$ (women)	

The marginal distribution for  $(X, Y)$  is

	$Y = 1$	$Y = 0$	
$X = 1$	.25	.25	.50
$X = 0$	.30	.20	.50
	.55	.45	1

confusion of causal effect  
with association

From these tables we find that,

$$\begin{aligned}\mathbb{P}(Y = 1|X = 1) - \mathbb{P}(Y = 1|X = 0) &= -0.1 \\ \mathbb{P}(Y = 1|X = 1, Z = 1) - \mathbb{P}(Y = 1|X = 0, Z = 1) &= 0.1 \\ \mathbb{P}(Y = 1|X = 1, Z = 0) - \mathbb{P}(Y = 1|X = 0, Z = 0) &= 0.1.\end{aligned}$$

To summarize, we *seem* to have the following information:

- graphs can be useful for clarifying dependence relations among random variables  
SM Markov random fields
- a **Directed Acyclic Graph** has random variables on the vertices and edges joining random variables

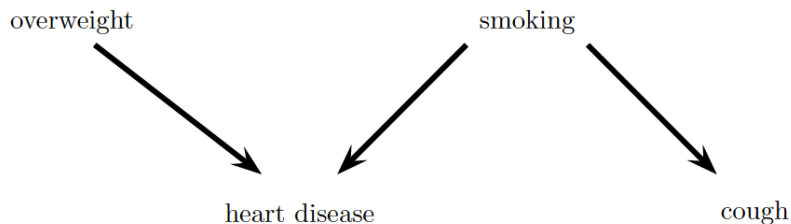


FIGURE 17.2. DAG for Example 17.4.

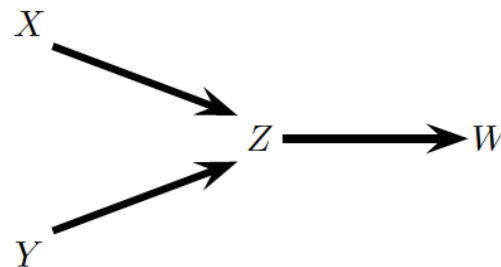


FIGURE 17.3. Another DAG.

- variables at parent nodes are potential causes for responses at child nodes
- directed graphs often helpful adjunct to modelling with baseline variables, intermediate responses, and outcome variables of interest
- much harder to study the full joint distribution than the usual supervised learning approaches
- DAGs can be used to represent confounders

276 17. Directed Graphs and Conditional Independence

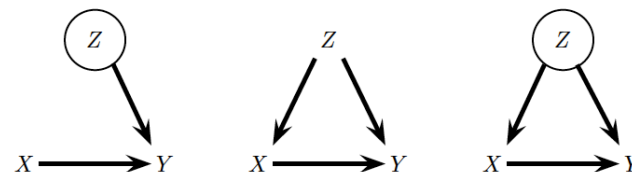


FIGURE 17.11. Randomized study; Observational study with measured confounders; Observational study with unmeasured confounders. The circled variables are unobserved.

276 17. Directed Graphs and Conditional Independence

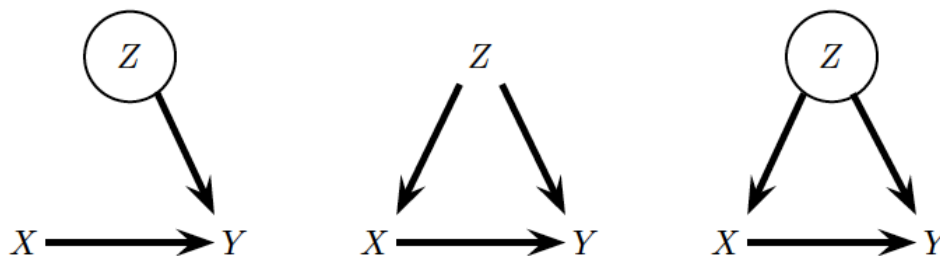


FIGURE 17.11. Randomized study; Observational study with measured confounders; Observational study with unmeasured confounders. The circled variables are unobserved.

randomized study

observational study  $E(Y | x) = \int E(Y | X, Z = z) dF_Z(z)$

unobserved confounder:  $\theta \neq \alpha$