

B555 2018

Probability

Read Wasserman Ch I

Probability deals with expts. having random outcomes

Ex

- ① Flip a coin
- ② Roll a die
- ③ Sample a voter

Defn

The sample space,  $\Omega$ , is set of all possible outcomes

Ex

Flip 2 coins :  $\Omega = \{ HH, HT, TH, TT \}$

Choose card from deck  $\Omega = \{ A, 2, \dots, 10, J, Q, K \} \times \{ \spadesuit, \heartsuit, \clubsuit, \diamondsuit \}$

Defn

An event is a subset of ~~prob~~  $\Omega$ .

Ex

$A =$  get a head in 2 coin flips

$$A = \{ HT, TH, HH \} \subseteq \Omega$$

Since events are subsets usual set operations are possible

$A^c = A$  did not occur

$A \cup B$ : Either A, B, or both occurred.

$A \cap B$ : Both A and B occurred

[ $\emptyset$  and  $\Omega$  also events]

(2)

Defn

A probability is a set function  $P(A)$ , for all possible subsets satisfying

$$\textcircled{1} \quad P(A) \geq 0$$

$$\textcircled{2} \quad P(\Omega) = 1$$

$$\textcircled{3} \quad \text{For disjoint } A_1, \dots, A_n \quad (A_i \cap A_j = \emptyset \text{ for } i \neq j) \quad P\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n P(A_i)$$

and

$$\text{For disjoint } A_1, A_2, \dots \quad P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i)$$

Simple relations follow

Ex

$$P(A^c) = 1 - P(A)$$

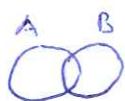
$$[1 = P(\Omega) = P(A \cup A^c) = P(A) + P(A^c)]$$

Ex

$$P(A \cup B) = P(A) + P(B) - P(A, B)$$

[ $P(A \cap B)$ ]

PF



$$\begin{aligned} P(A \cup B) &= (P(A \cap B^c) + P(A \cap B)) + (P(A^c \cap B) + P(A \cap B)) - P(A \cap B) \\ &= P(A) + P(B) - P(A \cap B) \end{aligned}$$

Equally Likely Outcomes

Often reasonable to assume all outcomes are equally likely

$$\Omega = \{HH, HT, TH, TT\} \quad \text{Then } P(A) = \frac{|A|}{|\Omega|}$$

Ex

Roll 2 dice

$\Omega = \{(1, 1), (1, 2), \dots, (6, 6)\} \times \{(1, 2, \dots, 6\}$

$P(S) = ?$

$|S| = 7$

$|\Omega| = 36$

(2)

36 outcomes non-eq prob  
each outcome has prob.  $1/36$ .

(3)

## Independence

### Defn

$A, B$  independent if  $P(A, B) = P(A)P(B)$

$A_1 \dots A_n$  mutually independent if for every subset  $J \subseteq \{1, \dots, n\}$

$$P\left(\bigcap_{j \in J} A_j\right) = \prod_{j \in J} P(A_j)$$

Intuitively if  $A$  has no effect on  $B$  and vice versa

$A$  occurs  $P(A)$  ~~times~~ of time

of times when  $A$  occurs,  $B$  occurs  $P(B)$  ~~times~~ of time

So  $A$  and  $B$  occur  $P(A)P(B)$  of time.

That is  $P(A, B) = P(A)P(B)$ .

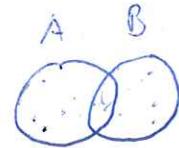
### Ex

Toss coin 10 times. What is  $P(\text{at least 1 H})$ ?

$$\begin{aligned} P(\text{at least 1 H}) &= 1 - P(\text{no heads}) \\ &= 1 - P(\text{1st fl., T}, \dots, 10^{\text{th}} \text{ fl., T}) \\ &= 1 - \left(\frac{1}{2}\right)^{10} \end{aligned}$$

## Conditional Probability

$P(A|B)$  is prob of  $A$  occurring when we know  $B$  occurs.



$P(A|B)$  = "prob of  $A$  given  $B$ "

Defn If  $P(B) > 0$   $P(A|B) = \frac{P(A, B)}{P(B)}$

Defn

$$P(A|B) = P(A \cap B) / P(B)$$

Note from

(4)

~~Ex~~ Have 2 white + 2 black balls in urn. Draw 2 balls w/o replacement. What is prob both black?

$$P(\text{Both black}) = P(B_1 \text{ black}, B_2 \text{ black}) = P(B_1 \text{ black} | B_2 \text{ black}) P(B_2 \text{ black})$$

$$= \left(\frac{1}{3}\right) \left(\frac{1}{2}\right) = \frac{1}{6}$$

Ex Law of Rare disease

Have rare disease existing in  $\frac{1}{1000}$  of population

Have test for disease s.t.

Test gives + result when person has disease w/ prob. .99  
 " " - " " " doesn't have " " "

Randomly chosen person tests positively for disease. What is prob. person has disease?

$$\text{Let } D = \text{person has disease} \quad P(D) = .001$$

$$P(D^c) = .999$$

$$\text{Have } P(+ | D) = P(- | D^c) = .99$$

$$P(D | +) = \frac{P(D, +)}{P(+)} = \frac{P(D) P(+ | D)}{P(D) P(+ | D) + P(D^c) P(+ | D^c)}$$

$$P(+) = P(+, D) \cup (+, D^c) = \frac{(.001)(.99)}{(.001)(.99) + (.999)(.01)} \approx \frac{1}{11}$$

(5)

## Law of Total Prob

Let  $A_1 \dots A_n$  be partition of  $\Omega$



$$(A_i \cap A_j = \emptyset \text{ for } i \neq j; \bigcup_{i=1}^n A_i = \Omega)$$

Then  $P(B) = \sum_{i=1}^n P(A_i) P(B|A_i)$

PF

$B = \bigcup_{i=1}^n A_i \cap B$  and  $A_1 \cap B, A_2 \cap B, \dots, A_n \cap B$  are disjoint so

$$P(B) = \sum_{i=1}^n P(A_i \cap B) = \sum_{i=1}^n P(A_i) P(B|A_i)$$

(Ex)

Have 3 black, 4 white, 2 red socks in drawer.

What is prob of choosing a matching pair?

$A_1 = 1^{st}$  sock black

[note  $A_1, A_2, A_3$  partition  $\Omega$ ]

$A_2 = \dots$  white

$B = \text{Matching pair}$

$A_3 = \dots$  red

$$P(B) = P(A_1) P(B|A_1) + P(A_2) P(B|A_2) + P(A_3) P(B|A_3)$$

$$= \frac{3}{9} \cdot \frac{2}{8} + \frac{4}{9} \cdot \frac{3}{8} + \frac{2}{9} \cdot \frac{1}{8} = \frac{15}{72}$$

## Bayes' Rule

Let  $A_1 \dots A_n$  partition  $\Omega$ . Then

Suppose we know  $P(B|A_i)$   $i = 1 \dots n$ .

Then  $P(A_i | B) = \frac{P(A_i, B)}{P(B)} = \frac{P(A_i) P(B|A_i)}{\sum_{j=1}^n P(A_j) P(B|A_j)}$

# Random Variables      Ross Wissneran Ch2

## Defn

A Random variable (rv) is function from  $\Omega$  to  $\mathbb{R}$  (Real #'s)

## Ex

Flip coin 3 times so  $\Omega = \{\text{HHH}, \text{HHT}, \text{HTH}, \text{HTT}, \text{THH}, \text{THT}, \text{THT}, \text{TTT}\}$

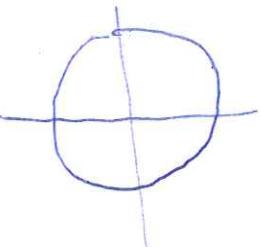
$X = \# \text{ H's}$        $\begin{matrix} \downarrow & \downarrow & \downarrow \\ 3 & 2 & 2 \end{matrix} \dots$

## Ex

Choose a point randomly from unit disc =  $\{(x, y) \in \mathbb{R}^2 : x^2 + y^2 \leq 1\}$

Let  $W = \text{dist. from origin}$

$$W = (x^2 + y^2)^{1/2}$$



## Defns

A discrete r.v. has enumerable outcomes  $\Omega = \{w_1, w_2, \dots\}$

$$\text{or } \Omega = \{w_1, w_2, \dots\}$$

A continuous r.v. has ~~one or several~~ intervals of  $\mathbb{R}$  as possible outcomes

$$[0, 1], \mathbb{R}^n, [0, \infty), [0, \eta] \cup [1, \eta]$$

## Important Discrete Distributions

### Bernoulli:

If  $X \in \{0, 1\}$  with  $P(X=1) = p$   
 $P(X=0) = q = 1-p$

$X$  is said to be Bernoulli. Often written

$$f(x) = p^x (1-p)^{1-x} \quad x=0, 1$$

Bernoulli is # of success from single S-F expt.  $P(S) = p$ .

Binomial

### Geometric

Have S-F expt. Perform repeated S-F expts until 1st S.

Let  $X = \# \text{ trials until } \cancel{S}$ .

$$P(X=1) = P(S) = p$$

$$P(X=2) = P(FS) = (1-p)p$$

$$P(X=3) = P(FFS) = (1-p)^2 p$$

$X$  has Geometric dist

$X \sim \text{Geometric}(p)$

$$P(X=x) = P(\underbrace{FF\dots F}_x S) = (1-p)^{x-1} p$$

### Poisson

$$X \sim \text{Poisson}(\lambda) \quad \lambda > 0 \quad \text{if} \quad f(x) = P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!} \quad x=0, 1, \dots$$

~~Poisson~~ Poisson models ~~unrelated~~ ~~events~~ <sup>in minute</sup> eg. # of radioactive quantities emitted by radioactive source where  $\lambda$  is average emissions per minute.

## Random Variables

Defn A Random Variable (rv) is a function,  $X$ ,

$$X: \Omega \rightarrow \mathbb{R}$$

↓  
sample space      ↓  
reals

Informally: an observed random quantity

Defn

If r.v.  $X$  has finite outcomes  $\{x_1, \dots, x_n\}$  or countably many outcomes it is discrete.

Write  $f(x) = P(X=x)$  or  $f_X(x) = P(X=x)$  for probabilities of outcomes.  $f(x)$  is probability mass function

Ex Flip a coin  $n$  times and let  $X$  be # H's.

Suppose coin is biased s.t.  $P(H) = p$ ;  $P(T) = q = 1-p$ .

Then  $\Omega = \underbrace{\{H,T\} \times \{H,T\} \times \dots \times \{H,T\}}_{n \text{ times}}$

$$X(HHTTHHTH) = 4$$

$\in \Omega$

Possible outcomes:  $x = 0, 1, \dots, n$

$$f(x) = P(X=x) = \binom{n}{x} p^x q^{n-x} \quad x=0, \dots, n$$

This is binomial dist. Write  $X \sim \text{Binomial}(n, p)$   
as shorthand.

## Poisson cont

In ~~discrete~~ <sup>continuous</sup> this situation divide interval (minute) into  $n$  gives  
 and suppose each subinterval has emission with prob.  $\frac{\lambda}{n}$

$\Rightarrow$  Total # of S's (emissions) would be

$$X_n \sim \text{Binomial}(n, \frac{\lambda}{n})$$

(Note "acc" # of S's is  $\lambda$  as it should be)

Can show  $\lim_{n \rightarrow \infty} \text{Binomial}(n, \frac{\lambda}{n}) = \text{Poisson}(\lambda)$

$$\lim_{n \rightarrow \infty} P(X_n = x) = \lim_{n \rightarrow \infty} \binom{n}{x} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x} \stackrel{x \rightarrow \infty}{\approx} \frac{e^{-\lambda} \lambda^x}{x!}, \dots$$

~~Continuous Random Variables~~  
~~X is continuous if it takes values~~

For any discrete dist. must have  $\sum_x f(x) = 1$

## Continuous Random Variables

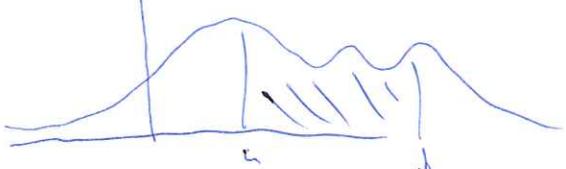
Continuous r.v.'s have a continuum of possible outcomes.

Eg  $[0, 1]$ ,  $[0, \infty]$ ,  $\mathbb{R}$ , etc.

A continuous r.v.  $X$  has prob. density function,  $f(x)$ , s.t.

$$P(a < X < b) = \int_a^b f(x) dx,$$

pdf  $f$  must have  $\int_{-\infty}^{+\infty} f(x) dx = 1$



NB IF  $X$  has pdf  $f(x)$   
 then  $P(X=x) = \int_x^\infty f(y) dy = 0$ .

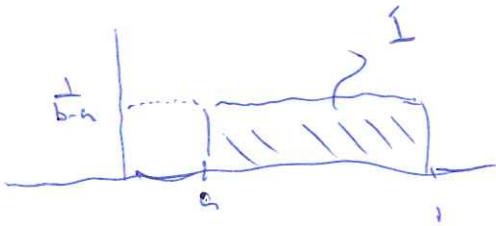
# Important Continuous Distributions

(9)

## Uniform

If  $X \in [a, b]$  and no outcome preferred over any other

$$X \sim \text{Unif}(a, b) \text{ with pdf } f(x) = \begin{cases} \frac{1}{b-a} & x \in [a, b] \\ 0 & \text{o.w.} \end{cases}$$



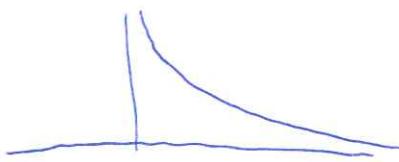
## Exponential

Recall radioactive particles where  $\lambda$  was ave # emissions per unit time.

How long do we wait for 1<sup>st</sup> emission? Let  $X$  be time we wait

$$X \sim \text{Exponential}(\lambda) \quad f(x) = \lambda e^{-\lambda x} \quad x \geq 0$$

0 o.w.



## Normal

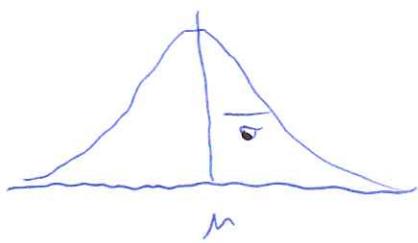
$X$  has normal dist. w/ parameters  $\mu, \sigma^2$  if

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

- (1) Sums of indep. rv's approx Normal.  
(Central Limit Thm)

- (2) Normals are easy to compute with

Eg  $X_1 \sim N(\mu_1, \sigma_1^2)$   $X_2 \sim N(\mu_2, \sigma_2^2)$   
 $X_1, X_2$  indep.  
 $X_1 + X_2 \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$



## **B555: Homework 1, Due Jan. 30 in class**

From Wasserman work problems

Ch 1: 5,6,9,11,12,15,17,19,20

Ch 2: 7,11,14,15,21



## B555 HW 1

## ( Chapter 1

5)  $\mathcal{A} = \{x_1 \dots x_n : x_1 = H, x_i = H \text{ some } i < n, x_j = T \text{ } j \neq i, j \neq n\}$

Let  $X$  be the flip giving the 2nd  $H$ .

$$P(X=x) = (x-1) \left(\frac{1}{2}\right)^x \quad x = 2, 3, \dots$$

6) The uniform dist must satisfy  $f(x) = c$  for  $x = 0, 1, \dots$

$$\text{If } c = 0 \text{ then } \sum f(x) = 0$$

$$\text{If } c > 0 \text{ then } \sum f(x) = \infty$$

So can't satisfy  $\sum f(x) = 1$ .

7) (i)  $P(A|B) = \frac{P(A,B)}{P(B)} > 0$

(ii)  $P(\neg A|B) = \frac{P(\neg A, B)}{P(B)} = \frac{P(B)}{P(B)} = 1$

8) Suppose  $A_1, A_2, \dots$  are disjoint events. Then

$$P\left(\bigcup_{i=1}^{\infty} A_i | B\right) = \frac{P\left(\left(\bigcup_{i=1}^{\infty} A_i\right) \cap B\right)}{P(B)} = \frac{P\left(\bigcup_{i=1}^{\infty} (A_i \cap B)\right)}{P(B)}$$

$$= \sum_{i=1}^{\infty} \frac{P(A_i \cap B)}{P(B)} = \sum_{i=1}^{\infty} P(A_i | B)$$

(2)

ii) Let  $A, B$  be independent. Then since  $A^c \cap B^c = (A \cup B)^c$

$$\begin{aligned} P(A^c \cap B^c) &= 1 - [P(A) + P(B) - P(A)P(B)] = (1 - P(A))(1 - P(B)) \\ &= P(A^c)P(B^c) \end{aligned}$$

iii) Let  $C \in \{1, 2, 3\}$  be chosen card and  $S \in \{\text{Red}, \text{Green}\}$  be the color we see.  $P(C=i) = \frac{1}{3} \quad i=1, 2, 3$ .

$$P(S = \text{Green} | C = i) = \begin{cases} 1 & i = 1 \\ 0 & i = 2 \\ \frac{1}{2} & i = 3 \end{cases}$$

$$P(C=i | S = \text{Green}) = \frac{P(C=i, S = \text{Green})}{P(S = \text{Green})} = \begin{cases} \frac{1/3}{1/2} & i = 1 \\ 0 & i = 2 \\ \frac{1/6}{1/2} & i = 3 \end{cases}$$

$$P(\text{other side Green} | S = \text{Green}) = P(C=1 | S = \text{Green}) = \frac{2}{3}.$$

15) a) Let  $N = \# \text{ children with blue eyes}$   $N \sim \text{Binomial}\left(\frac{5}{4}, \frac{1}{4}\right)$

$$P(N \geq 2 | N \geq 1) = 1 - P(N < 2 | N \geq 1) = \frac{P(N=1)}{P(N \geq 1)} = \frac{\frac{5}{4} \left(\frac{1}{4}\right)\left(\frac{3}{4}\right)^4}{1 - \left(\frac{3}{4}\right)^3}$$

b) The number of blue-eyed children from the remaining 2 children is  $\text{Binomial}(2, \frac{1}{4})$ .

$$P(\text{desired event}) = 1 - \left(\frac{3}{4}\right)^2$$

HW 1. Cont.

③

17)

Write  $P_C(\cdot)$  for the prob. given C.

$$\begin{aligned} P(A, B, C) &= P(A, B|C) P(C) = P_C(A|B) P_C(B) P(C) \\ &= P(A|B, C) P(B|C) P(C) \end{aligned}$$

18)  $P(W) = .5$ ;  $P(M) = .3$ ;  $P(L) = .2$

$$P(V|W) = .82; P(V|M) = .65; P(V|L) = .5$$

$$P(W|V) = \frac{P(W) P(V|W)}{P(W) P(V|W) + P(M) P(V|M) + P(L) P(V|L)}$$

$$= \frac{.41}{.41 + .195 + .1} = \frac{.41}{.705}$$

20) a)  $P(H) = \sum_{i=1}^5 P(C_i) P(H|C_i) = \frac{1}{5} \left( 0 + \frac{1}{4} + \frac{1}{2} + \frac{3}{4} + 1 \right) = \frac{1}{2}$

$$P(C_i|H) = \frac{P(C_i) P(H|C_i)}{P(H)} = \begin{cases} 0 & i=1 \\ \frac{1}{10} & i=2 \\ \frac{1}{5} & i=3 \\ \frac{3}{10} & i=4 \\ \frac{4}{10} & i=5 \end{cases}$$

b)  $P(H_2|H_1) = \cancel{P(H_2|H_1) \text{ or } P(H_1|H_2)} \quad 0 \cdot 0 + \frac{1}{10} \frac{1}{4} + \frac{2}{10} \frac{1}{2} + \frac{3}{10} \frac{3}{4} + \frac{4}{10} 1 = \frac{3}{4}$

c)  $P(C_i|B_4) = \frac{P(C_i, B_4)}{P(B_4)} \quad P(C_i|B_4) = \begin{cases} 0 & i=1 \\ \frac{\frac{1}{5} \left( \frac{3}{4} \right)^3 \frac{1}{4}}{41/80} & i=2 \\ \frac{\frac{1}{2} \left( \frac{1}{2} \right)^3 \left( \frac{1}{2} \right)}{41/80} & i=3 \end{cases}$

$$P(B_4) = \frac{1}{5} 0 + \frac{1}{5} \left( \frac{3}{4} \right)^3 \frac{1}{4} + \frac{1}{5} \left( \frac{1}{2} \right)^3 \left( \frac{1}{2} \right) + \frac{1}{5} \left( \frac{1}{4} \right)^3 \frac{3}{4} + \frac{1}{5} 0 = \frac{41}{80}$$

$$7) P(Z > z) = P(X > z, Y > z) = (1-z)^2$$

$$f_Z(z) = \frac{dP(Z \leq z)}{dz} = \frac{d(1-(1-z))^2}{dz} = 2(1-z) \quad 0 < z < 1$$

$$11) a) P(X=0, Y=0) = 0 \quad \text{yet} \quad P(X=0) = P(Y=0) = 1-p.$$

b) Note that  $X+Y = N$  so

$$\begin{aligned} P(X=x, Y=y) &= P(X=x, Y=y, N=n) = P(X=x | N=n) P(N=n)^y \\ &= \frac{e^{-\lambda} \lambda^{x+y}}{(x+y)!} \binom{x+y}{x} p^x (1-p)^y \\ &= \frac{e^{-\lambda} \lambda^x p^x}{x!} \frac{(1-p)^y}{y!} \end{aligned}$$

The joint pmf factors into  $g(x)h(y)$  so must be independent.

However, can show more directly by computing marginals.

$$\begin{aligned} P(X=x) &= \sum_{n=0}^{\infty} P(X=x | N=n) P(N=n) = \sum_{n=0}^{\infty} \binom{n}{x} p^x (1-p)^{n-x} \frac{e^{-\lambda} \lambda^n}{n!} \\ &= \sum_{n=x}^{\infty} \frac{p^x (1-p)^{n-x} e^{-\lambda} \lambda^n}{x! (n-x)!} \\ &= \sum_{n=0}^{\infty} \frac{p^x (1-p)^n e^{-\lambda} \lambda^{x+n}}{x! n!} \xrightarrow{1} \\ &= \frac{e^{-\lambda} (\lambda p)^x}{x!} \sum_{n=0}^{\infty} \frac{(-\lambda p)^n ((1-p)\lambda)^n}{n!} = \frac{e^{-\lambda} (\lambda p)^x}{x!} \end{aligned}$$

$X \sim \text{Poisson}(\lambda p)$  and by symmetry  $Y \sim \text{Poisson}((1-p)\lambda)$

Independence follows by direct comparison of  $P(X=x, Y=y)$  and  $P(X=x)P(Y=y)$

## (5)

### Chapter 2 cont

(14)  $F_R(r) = P(R \leq r) = r^2 \implies f_R(r) = 2r \quad 0 < r < 1$

(15)  $Y = F(X) \quad \text{note } 0 \leq Y \leq 1$

$$P(Y \leq y) = P(F(X) \leq y) = P(X \leq F^{-1}(y)) = y$$

(By defn  $F^{-1}(y)$  is # s.t.  $P(X \leq F^{-1}(y)) = y$ )

$$\implies Y = F(X) \sim \text{Unif}(0, 1)$$

Calculation shows if  $X \sim F$  then  $F(X) \sim \text{Unif}(0, 1)$

Reverse calc shows  ~~$U \sim \text{Unif}(0, 1) \implies F^{-1}(U) \sim F$~~

For  $\text{Ex}_\beta(B)$  dist.

$$f(x) = \frac{1}{\beta} e^{-x/\beta} \implies F(x) = 1 - e^{-x/\beta}$$

$$\implies F^{-1}(y) = \beta \log(1-y)$$

So if  $U \sim \text{Unif}(0, 1)$  then  $-\beta \log(1-U) \sim \text{Ex}_\beta(\beta)$

(21) Let  $Y = \max(X_1, \dots, X_n)$  where  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Ex}_\beta(\beta)$

$$F_Y(y) = P(Y \leq y) = P(X_1 \leq y, \dots, X_n \leq y) = (1 - e^{-y/\beta})^n$$

$$\implies f_Y(y) = \frac{n}{\beta} (1 - e^{-y/\beta})^{n-1} e^{-y/\beta} \quad y \geq 0$$



## Mean + Variance

For discrete r.v.  $X$  with pmf  $f(x)$  the expectation of  $X$  is

$$E(X) = \mu = \sum_x x f(x)$$

IF  $X$  is continuous with pdf  $f(x)$  " " " "

$$E(X) = \mu = \int_{-\infty}^{+\infty} x f(x) dx$$

$E(X) = \mu$  is the average value or "mean" of dist.

For discrete  $X$  with pmf  $f(x)$  the variance of  $X$  is

$$V(X) = \sigma^2 = \sum_x (x - \mu)^2 f(x) = \text{ave sq dist. from } \mu$$

so  $V(X)$  measures spread.

For continuous r.v.  $X$  with pdf  $f(x)$

$$V(X) = \sigma^2 = \int_{-\infty}^{+\infty} (x - \mu)^2 f(x) dx$$

## Back to Normal

Recall  $X \sim N(\mu, \sigma^2)$  if  $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$

The "parameters" of normal,  $\mu, \sigma^2$ , are in fact mean + variance

$$\int_{-\infty}^{+\infty} x \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx = \mu$$

$$\int_{-\infty}^{+\infty} (x - \mu)^2 \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx = \sigma^2$$

## Computational Ex

As stated before sums of ~~several~~ indeg. r.v.'s are approx normal.

~~More precisely:~~

Ex

$$X_1, \dots, X_5 \stackrel{\text{indep.}}{\sim} \text{Unif}(0, 1)$$

$$Y = \sum_{i=1}^5 X_i$$

$$\textcircled{1} \quad E(X_i) = \int_0^1 x \cdot 1 dx = \frac{1}{2}$$

$$\textcircled{2} \quad V(X_i) = \int_0^1 (x - \frac{1}{2})^2 \cdot 1 dx = \frac{1}{12}$$

$$\textcircled{3} \quad E(Y) = E\left(\sum_{i=1}^5 X_i\right) = \sum_{i=1}^5 E(X_i) = \frac{5}{2}$$

$$\textcircled{4} \quad V(Y) = V\left(\sum_{i=1}^5 X_i\right) = \sum_{i=1}^5 V(X_i) = \frac{5}{12}$$

$$\textcircled{5} \quad Y \stackrel{\text{approx}}{\sim} N\left(\frac{5}{2}, \frac{5}{12}\right)$$

Let's simulate  $Y$  many times and compare

a)  $\hat{P}(Y \leq y) = \frac{\# \text{Events} \leq y}{n}$   
with

b)  $P(N(\frac{5}{2}, \frac{5}{12}) \leq y)$  [can compute this]

## CLT

$X, X_1, \dots$  sequence of indeg. r.v.'s with same dist. with mean  $\mu$  + var  $\sigma^2$ . Then for any  $[a, b]$

$$P\left(\frac{\sum_{i=1}^n X_i - n\mu}{\sqrt{n}\sigma} \in [a, b]\right) \xrightarrow{n \rightarrow \infty} P(Z \in [a, b])$$

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

## An Introduction to R

### To get R:

1. Download R (it's free) from the website <http://cran.r-project.org> There are versions for Linux, Windows and Mac.
2. R will be installed in the public computing cluster in room 109 in the Informatics building.

### R can be used as a calculator:

Try typing the following expressions at the command line (followed by return): (> is the command prompt).

```
> 5+3  
> 10*10  
> log(9.4)  
> exp(exp(exp(20))) # R is only human! (anything following '#' is a comment)
```

R has most any mathematical function you can think of such as sqrt(), sin() ... mostly with easily guessable names. Expressions using the logical operators ==, !=, <, > give Boolean values (T,F)

```
> 4 > 3      # this evaluates to T (true)  
> 1 == exp(0) # so does this  
> 1 != exp(0) # this evaluates to F (false)
```

It is possible to have variables that hold values in your program. Most strings beginning with an alphabet character will be treated as variables. The assignment operator is <- (a "less than" followed by "minus"). Try typing the following lines in succession

```
> x <- 3  
> y <- x*x+x  
> y      # print the value of y
```

### Vectors

One of the nicest aspects of R is the way it handles vectors. Here are a several ways to create vectors:

```
> x <- 1:100          # x is now the vector (1,2,...,100)  
> y <- seq(-pi,pi,length=100) # y consists of 100 evenly spaced values from -pi to pi  
> z <- c(1,4,8,20)      # z is the vector (1,4,8,20)  
> a <- x+y            # vectors of same length can be added, multiplied, etc.  
> b <- 4*x            # this is interpreted correctly too
```

### Random Number Generation

R has lots of built-in functions for doing things with random numbers. For instance

```
> x <- runif(100)    # creates a vector of 100 (uniformly distributed) random numbers between 0 and 1.  
> punif(v)           # is the probability that a Unif(0,1) rand number is less than v  
> qunif(u)           # gives the u'th quantile of a Unif(0,1). More on this later.
```

There are similar functions for a variety of other distributions including the normal(0,1) (rnorm,pnorm,qnorm) Cauchy (rcauchy, pcauchy, qcauchy), Exponential, Binomial, Poisson, and others.

### Subsets

```
> x <- runif(100)    # creates a vector of 100 Unif(0,1) random numbers  
> x[1]                # the first element of x  
> x[c(1,3,5)]        # a vector containing 1st, 3rd and 5th elements of x  
> y <- x > .5         # a 100-long vector of Boolean values y[i] is T iff x[i] > .5  
> z <- x[x>.5]        # the "x's" that are greater than 5
```

## Plotting Try the following

```
> x <- seq(0,1,length=100)
> y <- x^2                      # y = x squared
> plot(x,y)                     # plot with (x[1],y[1]) \ldots, (x[100],y[100])
> plot(y,x)
> plot(y)      # same as plot(1:length(y),y)
```

**Source Files** You will want to write simple programs in R and this always requires some trial, error and iteration. I recommend the following procedure: Create a “source” file in any text editor containing your R commands. This could be emacs or the Windows “Notepad” or whatever you are comfortable using. Suppose you create the following file named “myprog.R” in your editor:

```
len <- 100
x <- runif(len,-.5,.4)
y <- cumsum(x)    # y[1] = x[1], y[2] = x[1]+x[2], etc.
plot(exp(y))
title("my stock price")
print("history is: ")
print(y)
```

This technique allows you to write a program in the usual incremental way. If you want to get a hard copy of the printout and the plot (for example, to submit as your homework), do the following

```
> postscript("myplot.ps")  # write plot in the postscript file "myplot.ps"
> sink("myout.txt")        # write text output to "myout.txt"
> source("myprog.R")       # run the program you created
> dev.off()                # redirect plots to screen. Don't forget this!
> sink()                   # redirect output to screen. ditto.
```

**A fun example** Suppose two decks of cards are shuffled. The first deck is arranged in a line, face up, and the second deck is arranged in a similar line right below. Count the number of places where the two decks have the same card. What is the probability that there are no matches between the two decks? (This is a hard calculation to do). One way to estimate this probability would be to perform the experiment many times and observe the proportion of times the event occurs.

```
> trials <- 1000           # number of trials
> zeros <- 0               # counts the number of times the "no match" event occurs
> for (i in 1:trials) {     # all statements in the "for" loop are executed ---
                           # once for each value of i in 1:trials
  >   x <- 1:52;
  >   deck1 <- sample(x,52,replace=F)  # a random permutation of the "cards"
  >   deck2 <- sample(x,52,replace=F)  # another random permutation
  >   matches <- sum(deck1==deck2)    # number of matches
  >   if (matches == 0) zeros <- zeros+1 # if (matches == 0) we had a no matches so count it
  > }
> print("my estimate is:")
> print(zeros/trials);     # the proportion of times the "no match" event occurred
```

## Quitting and help

```
> help("rnorm")  # gives information about the function rnorm. Of course this works
>                  # for other functions too.
> q()  # quitting the program. Hope you had fun.
```

## Joint Distributions

Several rvs  $X_1, \dots, X_n$ , perhaps interdependent, have a joint distribution,  $f(x_1, \dots, x_n)$

For  $X_1, \dots, X_n$  discrete  $f(x_1, \dots, x_n)$  is the joint pmf defined by

$$f(x_1, \dots, x_n) = P(X_1=x_1, X_2=x_2, \dots, X_n=x_n)$$

For  $X_1, \dots, X_n$  continuous  $f(x_1, \dots, x_n)$  is the joint pdf

$$P(a_1 < X_1 < b_1, \dots, a_n < X_n < b_n) = \int_{a_1}^{b_1} \dots \int_{a_n}^{b_n} f(x_1, \dots, x_n) dx_n \dots dx_1$$

Must have  $\sum_{x_1, \dots, x_n} f(x_1, \dots, x_n) = 1$  for discrete case

$$\int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} f(x_1, \dots, x_n) dx_n \dots dx_1 \quad \text{for continuous case.}$$

## Marginal Distributions

For bivariate discrete  $X, Y$  can represent pmf as table

		$X$		$P$
		0	1	
$Y$	0	.1	.2	.3
	1	.3	.4	.7
		.4	.6	

Summing cols/rows give "marginal" distributions  
(distributions of  $X$  or  $Y$  with no knowledge of other vble.)

More generally, for discrete  $X_1, \dots, X_n$  with pmf  $f(x_1, \dots, x_n)$   
the marginal distributions  $f_{X_i}(x_i), \dots, f_{X_n}(x_n)$  are

$$f_{X_i}(x_i) = \sum_{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n} f(x_1, \dots, x_n)$$

# Independence

Random variables  $X_1 \dots X_n$  are independent if ~~if~~ their joint is product of margins.

$$f(x_1 \dots x_n) = f_{x_1}(x_1) \dots f_{x_n}(x_n)$$

This holds for both continuous + discrete r.v's.

Important Examples  $\{X_i \in A_i\}, \{X_1 \in A_1\}, \dots, \{X_n \in A_n\}$  are independent.

## Multinomial

Discrete  $X_1, X_2, \dots$  have a multinomial dist.

$$x_1 \dots x_n \sim \text{Multinomial}(p_1, \dots, p_n, n)$$

Have  $n$  independent trials with  $K$  possible outcomes w/ prob  
 $p_1, \dots, p_K$ . Let  $X_k = \#$  times outcome  $k$  occurs.

$X_1 - X_k$  said to have multinomial dist  $X_1 - X_n \sim \text{Multinomial}(p_1 - p_k, n)$

$$f(x_1, \dots, x_n) = P(X_1=x_1, \dots, X_n=x_n) = \frac{N!}{x_1! x_2! \dots x_n!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k} \quad \begin{matrix} x_1, \dots, x_n = 0, 1, \dots, n \\ \text{s.t. } \sum x_i = n \end{matrix}$$

Q: What is  $f_{x_1}(x_1)$ ? A: Each trial either type I ( $\gamma_1$ ) or not ( $1 - \gamma_1$ )

$$\Rightarrow X_i \sim \text{Binomial}(n, p_i) \quad f_{X_i}(x_i) = \binom{n}{x_i} p_i^{x_i} (1-p_i)^{n-x_i}, \quad x_i = 0, \dots, n$$

~~3000 feet~~

## ~~Expectancy + Covariance~~

$$\left[ \sum_{\substack{x_1 \dots x_n \\ \sum x_i = n}} \frac{n!}{x_1! \dots x_n!} p_1^{x_1} \dots p_n^{x_n} = \frac{n!}{x_1! (n-x_1)!} p_1^{x_1} (1-p_1)^{n-x_1} \right]$$

Notation Often better to use vector notation and write

$$X = \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix} \quad \text{for collection of random variables.}$$

## (14)

### Expectation + Covariance

(3) For random vector  $X = \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix}$  the expectation of  $X$ ,  $E(X) = \begin{pmatrix} EX_1 \\ \vdots \\ EX_n \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix}$

where  $EX_i$  is expectation of  $i$ th vble.

That is

$$EX_i = \sum_{x_i} x_i f_{X_i}(x_i) = \sum_{x_1 \dots x_n} x_i f(x_1 \dots x_n) \quad (\text{discrete})$$

$$EX_i = \int_{x_i} x_i f_{X_i}(x_i) dx_i = \int_{x_1} \dots \int_{x_n} x_i f(x_1 \dots x_n) dx_1 \dots dx_n \quad (\text{continuous})$$

Covariance ① For  $X, Y$  rvs  $\text{Cov}(X, Y) = E(X - \mu_X)(Y - \mu_Y)$  [the covariance] Do correlation

For random vectors variance generalizes to covariance matrix. ③

(4) For random vector  $X = \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix}$  the covariance matrix,  $\Sigma$ , is given by

$$\Sigma_{ij} = E(X_i - \mu_i)(X_j - \mu_j)$$

Observe that  $\Sigma_{ii} = E(X_i - \mu_i)^2 = V(X_i)$  so diagonal elements of  $\Sigma$  are variances.

② Intuitively, if  $X_i, X_j$  vary together (one high (low) when other high (low)) then  $(X_i - \mu_i)$  and  $(X_j - \mu_j)$  tend to have same sign, so

$$\Sigma_{ij} = \text{Cov}(X_i, X_j) > 0$$

$$\begin{pmatrix} X_1 & X_2 & \dots & X_n \end{pmatrix} \begin{pmatrix} X_1 & X_2 & \dots & X_n \end{pmatrix}^T$$

③ Alternatively, if  $X_i, X_j$  vary inversely (one high when other low)  $(X_i - \mu_i)$  and  $(X_j - \mu_j)$  tend to have opposite sign so

$$\Sigma_{ij} = \text{Cov}(X_i, X_j) < 0$$

Often write  $X = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{pmatrix}$ ,  $\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix}$

$$\Sigma = \text{Cov}(X) = E(X - \mu)(X - \mu)^T = E$$

## Important Facts about Expectations

① Expectation is linear. That is,

Expectation of sum is  
sum of expectations

**[HW]** If  $X_1, \dots, X_n$  have  $\mathbb{E}X_1 = \mu_1, \dots, \mathbb{E}X_n = \mu_n$

$$\mathbb{E}(\lambda_1 X_1 + \lambda_2 X_2 + \dots + \lambda_n X_n) = \lambda_1 \mu_1 + \lambda_2 \mu_2 + \dots + \lambda_n \mu_n.$$

In vector notation, if  $\lambda^t = (\lambda_1 \dots \lambda_n)$  and  $x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$

$$E(\lambda^t X) = \lambda^t E(X)$$

or, if  $A_{k \times n}$  is matrix

$$\left( \begin{array}{ccc} c_{11} & \dots & c_{1n} \\ c_{21} & \dots & c_{2n} \\ \vdots & & \\ c_{k1} & \dots & c_{kn} \end{array} \right) = A$$

$$E(AX) = A E(X)$$

② If  $X_1, X_2$  independent,  $E(X_1 X_2) = E(X_1) E(X_2)$

## Joint Normal

$Z$  is standard normal if  $Z \sim \mathcal{N}(0, 1)$ .  $[f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2}]$

Suppose  $Z_1, \dots, Z_n$  are standard normal + independent? If

$$Z \in \begin{pmatrix} z_1 \\ \vdots \\ z_n \end{pmatrix} \quad E(Z) = O = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} \quad \text{and}$$

$$\text{Cov}(Z_i, Z_j) = E Z_i Z_j = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases} \quad (EZ_i Z_j = EZ_i EZ_j)$$

So the covariance matrix of  $Z$  is  $I_{n \times n}$  (identity matrix)

Correlation

$$\text{Recall } |\sum_{i=1}^n x_i y_i| \leq (\sum_{i=1}^n x_i^2)^{1/2} (\sum_{i=1}^n y_i^2)^{1/2} \quad \text{Cauchy-Schwarz}$$

Also by C.S.

$$\begin{aligned} |\text{Cov}(X, Y)| &= \left| \sum_{x,y} (x - \mu_x)(y - \mu_y) f(x, y) \right| \leq \left( \sum_{x,y} (x - \mu_x)^2 f(x, y) \right)^{1/2} \left( \sum_{x,y} (y - \mu_y)^2 f(x, y) \right)^{1/2} \\ &= \text{Var}(X)^{1/2} \text{Var}(Y)^{1/2} = \sigma_x \sigma_y \end{aligned}$$

$\Rightarrow$

$$-1 \leq \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y} \leq 1$$

$\text{Define correlation of } X, Y \text{ to be}$

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y} = \rho_{XY}$$

said to be.

Have  $-1 \leq \rho_{XY} \leq 1$ . If  $\rho_{XY} = 0$   $X, Y$  uncorrelated

Correlation + Independence

$$\text{If } X, Y \text{ indep. } \text{Cov}(X, Y) = E((X - \mu_X)(Y - \mu_Y)) = E(X - \mu_X) E(Y - \mu_Y) = 0$$

Indep  $\Rightarrow$  Uncorrelated

HW Shows examp.



## Multivariate Normal

(16)

$X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$  has multivariate normal dist.  $\wedge X \sim N(\mu, \Sigma)$  if w mean  $\mu$  and cov  $\Sigma$

$$f(x) = f\left(\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}\right) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

① By direct calc. can show  $E(x) = \int x f(x) dx = \mu$

$$\text{cov}(x) = \int (x-\mu)(x-\mu)^T f(x) dx = \sum$$

② Consider  $z_1, \dots, z_n \stackrel{iid}{\sim} N(0, 1)$   $Z = \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{pmatrix}$  has

$$f(z) = \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2} z^T z} \Rightarrow Z \sim N(0, I)$$

③ For invertible  $A_{k \times n}$   $Az$  has

a)  $E(Az) = AE(z) = A0 = 0$

b)  ~~$\Sigma =$~~   $\text{cov}(Az) = E(Az(Az)^T) = E(Azz^T A^T) = A Ezz^T A^T$

$$= AIA^T = AAT$$

$$(\Sigma = AA^T \Rightarrow \Sigma^{-1} = A^{-T}A^{-1})$$

$$\begin{array}{c} z \\ + \\ + \\ x = Az \end{array}$$

④ Let  $X = AZ$  What is pdf of  $X$ ?

$$f_x(x) = \frac{1}{|A|} f_z(A^T x) = \frac{1}{(2\pi)^{n/2} |A|} e^{-\frac{1}{2} x^T A^T A^{-1} A^{-T} x} = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} x^T \Sigma^{-1} x}$$

$$\Sigma = AA^T$$

$$\Rightarrow X \sim N(0, \Sigma = AA^T)$$

$$|\Sigma| = |A||A| \Rightarrow |\Sigma|^n = |A|^n$$

(5) Any time have  $X \sim \mathcal{N}(\mu, \Sigma)$  can think of  $X = AZ + \mu$

where  $A$  is chosen s.t.  $AA^T = \Sigma$ .

Do pairs. To demonstrate.

### Conditional Distributions

The conditional dist. of  $X$  given  $y$ ,  $f_{x|y}(x|y) = f(x|y)$

is pmf or pdf for  $X$  when known that  $y=y$ .

For  $X, Y$  discrete

$$f(x|y) = \frac{P(X=x, Y=y)}{P(Y=y)} = p(x|y) = \frac{f(x,y)}{f_y(y)}$$

For  $X, Y$  continuous still have

$$f(x|y) = \frac{f(x,y)}{f_y(y)}$$

where  $P(X \in A | Y=y) = \int_{x \in A} f(x|y) dx$

Similarly, for  $X, Y, Z$   $f(x,y|z) = \frac{f(x,y,z)}{f_z(z)}$

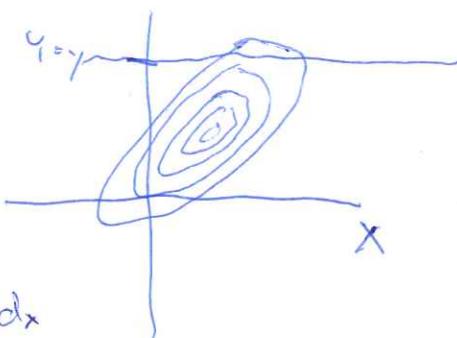
is conditional pmf, pdf of  $X, Y$  given  $Z=z$ .

### Conditional Independence

For  $X, Y, Z$ , if  $X, Y$  indep. in conditional dist. given  $Z=z$  then  $f(x,y|z) = f(x|z)f(y|z)$

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

Do region.  $\rightarrow$  could have all different giving M/F, yet overall having different acceptance rates to M/F



Indep  
X  
Y  
N

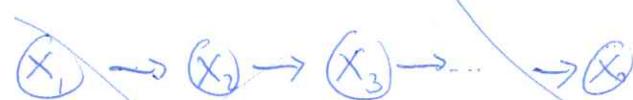
Indep  
X  
Y  
Z

NB

## (R)

### Important Graphical Models

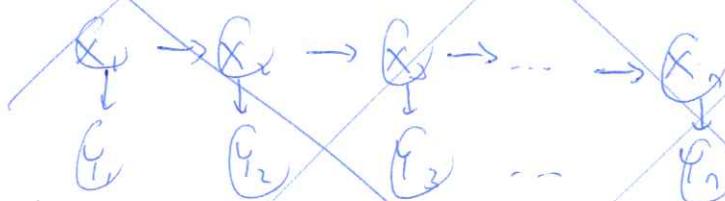
①  $X_1 \dots X_n$  respect



$$p(x_1 \dots x_n) = p(x_1) p(x_2|x_1) p(x_3|x_1, x_2) \dots p(x_n|x_{n-1})$$

$X_1 \dots X_n$  form a Markov Chain.

②  $X_1 \dots X_n, Y_1 \dots Y_n$  respect



$$p(x_1 \dots x_n, y_1 \dots y_n) = p(x_1) p(x_2|x_1) \dots p(x_n|x_{n-1}) p(y_1|x_1) p(y_2|x_1, x_2) \dots p(y_n|x_1, \dots, x_{n-1})$$

③  $X, Y$  are a hidden Markov Model.

### Mixture Distributions Read Ch 3 Wasserman (not moment generating fns)

Have  $M$  different prob. models  $p_1(x) \dots p_M(x)$ . To generate  $X$ :

choose model 1 with prob  $q_1$

" " 2 " "  $q_2$

" " " M " "  $q_M$

$$\left( \sum_{m=1}^M q_m = 1 \right)$$

and sample from chosen model.

What is dist. for  $X$  obtained this way? By LoTP

$$p(x) = \cancel{\sum_{m=1}^M p(m) p(x|m)} \sum_{m=1}^M q_m p(x|m) = \sum_m q_m p_m(x)$$

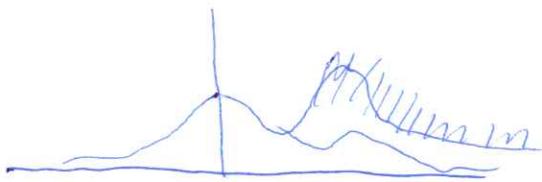
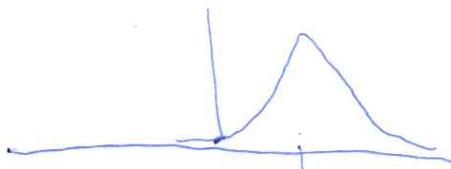
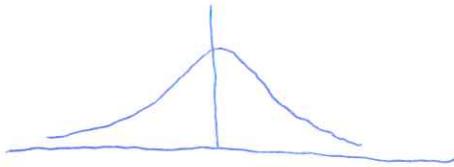
$p(x)$  is mixture of  $p_1 \dots p_M$  with mixing weight  $q_1 \dots q_M$

Ex

$$\text{Have } q_1(x) = \mathcal{N}(x; 0, 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}$$

$$q_2(x) = \mathcal{N}(x; 1, 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-1)^2}$$

Choose model 1 with prob  $\frac{3}{4}$   
 " 2 " "  $\frac{1}{4}$



Many (maybe most) real world distributions are mixtures.

Population contains several types of individuals, each with characteristic dist.

The mixing weights are the proportions of each type.

In sampling characteristic from pop., choose type according to ~~characteristic representation~~. Then sample from ~~type dist.~~ ~~type dist.~~.

① Ex

② Ex Poisson-binom. r

Find mixture dists. in NYT Game of Thrones

Likely Mixtures: Joffrey Baratheon, Arya Stark, Bran Stark, Tyrion Lannister.

Are judgements of Goodness + Beauty independent?

$Z \sim \text{Poisson}(\lambda)$   $X, Y = H(S(T))$  from  $Z$  flips

$P(H) = p$   $P(T) = q$   $X \sim \text{Binom}(Z, p)$   $Y = Z - X$

$$P(X=x) = \sum_{z=0}^{\infty} \frac{\lambda^z e^{-\lambda}}{z!} p(z=x) P(X=x|Z=z)$$

$$\text{Data.} = \sum_{z=0}^{\infty} \frac{\lambda^z e^{-\lambda}}{z!} \binom{z}{x} p^x (1-p)^{z-x}$$

Observe  $X, Y$  indep.

## Conditional Independence

(20)

Write  $g(x)$ ,  $g(y|x)$  etc. for pmf, pdf, conditional pmf/pdf etc.

① For r.v's  $X, Y, Z$ ,  $X, Y$  conditionally independent given  $Z$

if  $X, Y$  independent under conditional dist given  $Z=z$ . That is

$$g(x, y|z) = g(x|z)g(y|z)$$

② Ex

Suggest roll die ( $Z$ ). Then take 2 fair coins and let

$X = \# H's$  in  $Z$  flips of coin 1

$Y = \# H's$  in  $Z$  flips of coin 2

When  $X$  is high  $Y$  will also tend to be high so (informally)  $X, Y$  not independent.

But once we know  $Z$  (conditioned on  $Z$ )  $X, Y$  independent.

Given  $Y, X, Z$  indep.  
What is a C.I. relationship?  
reasonable

Given  $Y, X, Z$  indep.

person owns smokes gets lung cancer  
person owns smokes gets lung cancer  
X Y Z

④ Note

If  $X, Y$  independent then conditioning on  $X$  adds no info on  $Y$ . That is

$$g(y|x) = g(y) \quad [ \underbrace{g(y|x)g(x)}_{=g(x,y)} = g(x)g(y) ]$$

Similarly, if  $X, Y$  conditionally independent given  $Z$ , once we know  $Z$  additional knowledge of  $X$  gives no new info about  $Y$ . I.e.

$$\begin{aligned} g(x|z) &= g(x) \\ g(x|y,z) &= g(x|z) \end{aligned}$$

So factorization simplifies

$$g(x,y,z) = g(z)g(y|z)g(x|y,z) = g(z)g(y|z)g(x|z).$$

⑤

NB
$X, Y$ conditionally independent given $Z$

$\cancel{\Rightarrow}$   $X, Y$  indep

Do Simpson's paradox  
(Simpson's Paradox)

# Bayesian Network

61

As shown in HW can always factor joint dist through "chain rule": ( )

$$\gamma(x, y, z) = \gamma(x) \gamma(y|x) \gamma(z|y,x)$$

Factorization not unique since other possible orderings:

$$\gamma(x,y,z) = \gamma(z)\gamma(y|z)\gamma(x|y,z)$$

There are  $n!$  possible orderings of  $n$  variables so  $n!$  possible factorizations

But they are all equal. Eg

$$g(x,y) = g(x) \underbrace{g(y|x)}_{\rightarrow} = g(x) \frac{g(x,y)}{g(x)} = \underbrace{\frac{g(y)}{g(y)}}_{\rightarrow} g(x,y) = g(y) g(x|y)$$

(Bcx)

End. Ind. makes factorization simpler.

For instance, if  $X_n$  cond. ind. of  $X_1 \dots X_{n-1}$  given  $X_{n-1}$  then

$$g(x_n | x_1, x_2, \dots, x_{n-1}) = g(x_n | x_{n-1})$$

IF the  $\{x_i\}$  are binary r.v's.

LHS:  $2^{n-1}$  Bernoulli distributions      RHS: 2 Bernoulli dists.

Roughly speaking <sup>training</sup> LHS requires amount of data prop. to  $2^{n-1}$

LTS may well be intractable for large  $n$ .

Bayesian Network uses cond. indeg. to simplify factorization.

Eg, suppose for  $X, Y, Z, W$

$y$  cond. ind. of  $x$  given  $w$        $p(y|x,w) = p(y|w)$

$$\text{Second inde of } W \text{ given } X, Y \quad \gamma(z|w, x, y) = \gamma(z|x, y)$$

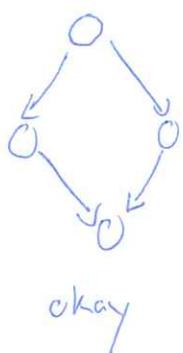
## Graphical Representation

Factorization of joint dist can be represented by directed acyclic graph (DAG)

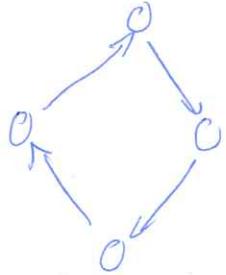
G: Graph is collection of nodes/vertices (random variables) and edges (dependency)

D: Edges are directed

A: no directed cycles



okay



not okay.

Ex

$$\gamma(x, y, z) = \gamma(x)\gamma(y|x)\gamma(z|x, y)$$



$$\gamma(x, y, z) = \gamma(x)\gamma(y|x)\gamma(z|y)$$



[ Given Y, X, Z cond. indep ]

Parents of variable in graph are variables we condition on.

More generally

IF  $A \subseteq \{1, \dots, n\}$  write  $X_A = \{X_i : i \in A\}$

still write  $X_i$  for  $X_{\{i\}}$

Suppose  $X_1, \dots, X_n$  r.v.'s w/~~XXXXXX~~ described by DAG.

The following statements are equivalent

(1)  $X_1, \dots, X_n$  respects a given DAG

(2)  $q(x_1, \dots, x_n) = \prod_{i=1}^n q(x_i | X_{P(i)})$  where  $P(i)$  are "parents" of  $i$  in DAG (all nodes connected to  $i$ )

(3) Let  $A(i)$  be the ancestors of  $i$ . That is

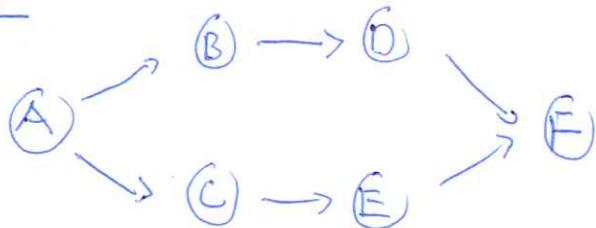
$j \in A(i)$  iff there is directed path from  $j$  to  $i$

$$j \rightarrow k_1 \rightarrow k_2 \rightarrow \dots \rightarrow i$$

set subtraction.

Given  $P(i)$   $X_i$  cond indeg. of  $A(i) \setminus P(i)$ .

Ex



$q(a, b, c, d, e, f)$  respects graph

$$q(a, b, c, d, e, f) = q(a) q(b|a) q(c|a) q(d|b) q(e|c) q(f|d, e)$$

C.I's

$P(F)$

$A(F) \setminus P(F)$

For ex. Given  $\{D, E\}$ ,  $F$  cond indeg. of  $\{A, B, C\}$   
with similar statements for other vbls.

- ① For each node with no parents simulate from marginal dist.  
② Find node  $i$  s.t. all of  $i$ 's parents have been simulated  
Draw  $X_i$  from  $q(x_i | X_{P(i)})$

Simulation

Can think of  $q(x_1, \dots, x_n) = \prod_{i=1}^n q(x_i | X_{P(i)})$  as algorithm for simulating from  $\text{joint dist.}$

## B555: Homework 2

1. Suppose that  $X$  is a discrete  $n$ -dimensional random vector

$$X = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{pmatrix}$$

and that  $a \in \mathbb{R}^n$  ( $a$  is a constant  $n$ -dimensional vector). Show  $E(a^t X) = a^t E(X)$  where  $a^t$  denotes the transpose of  $a$ . (The result also holds for continuous random vectors).

2. Let  $X, Y$  be independent continuous random variables with density  $f(x, y)$ , where  $f(x, y) = g(x)h(y)$  for some functions  $g, h$ . Show that  $X$  and  $Y$  are independent.
3. For this problem, recall that random variables  $X, Y$  are conditionally independent given  $Z$  if  $X$  and  $Y$  are independent under the conditional distribution  $f(x, y|z)$  (the distribution of  $X, Y$  having observed  $Z = z$ ).

Suppose that we consider random minutes from randomly chosen drivers to determine three binary random variables:

- $H$  = the driver is within 1 mile of home
- $M$  = the driver is on a main road
- $A$  = the driver has had an accident in the sampled minute

Suppose that  $P(H = 1) = .8$  along with  $P(M|H)$  and  $P(A = 1|H, M)$  in the following two tables.

$P(M H)$	$M = 0$	$M = 1$
$H = 0$	.1	.9
$H = 1$	.6	.4

$P(A = 1 H, M)$	$M = 0$	$M = 1$
$H = 0$	.0003	.0008
$H = 1$	.0003	.0008

- (a) Compute  $P(A = 1|H = 0)$  and  $P(A = 1|H = 1)$ . Are accidents more likely near home?
- (b) Compute  $P(A = 1|M = 0)$  and  $P(A = 1|M = 1)$ . What conditional independence relation can you conclude?
- (c) Give a cause and effect hypothesis that is consistent with your conditional independence statement.
- (d) Suppose that a driver is driving on a stretch of highway (main road) on her way home. As she crosses the 1-mile-to-home road marker does she encounter an elevated probability of accident?

### 4. Simulation of Binomial

- (a) A random variable has  $X \sim \text{Unif}(0, 1)$  when its density function is

$$f(x) = \begin{cases} 1 & 0 < x < 1 \\ 0 & \text{otherwise} \end{cases}$$

The R command “runif(n)” creates a vector of  $n$  independent  $\text{Unif}(0, 1)$  outcomes. Use this function to simulate a sequence of 10,000 independent  $\text{Binomial}(500, 1/2)$  random variables and plot the resulting histogram. You can plot a histogram of a data vector  $x$  with “hist(x)”.

- (b) Letting your sequence of Binomial variables be  $X_1, \dots, X_n$ , “standardize” these variables by taking  $Y_n = (X_n - \mu)/\sigma$  where  $\mu$  and  $\sigma$  are the mean and standard deviation of the variables. What are the mean and variance of the  $Y_n$  variables?

- (c) Now simulate 10,000 *standardized* Binomials and plot the empirical cumulative distribution function  $\hat{F}(x)$  where

$$\hat{F}(x) = \hat{P}(X \leq x) = \frac{|\{i : X_i \leq x\}|}{n}$$

Use “plot(x,type='l’)” to plot with a solid line

- (d) On the same plot show the cumulative distribution function of the standard normal. “pnorm(x)” gives the probability that a  $N(0, 1)$  variable is less than  $x$ . You can add to an existing plot with “lines(x).” Your two plots should look very similar. If  $X$  is a Binomial with large  $n$ , the fact that

$$\frac{X - np}{\sqrt{np(1-p)}} \stackrel{\text{approx}}{\sim} N(0, 1)$$

is known as the normal approximation to the Binomial.

5. It is well-known that independent random variables are also uncorrelated.

6. Consider the variables  $X = (X_1, X_2)$  having probability density function

$$p(x) = \begin{cases} \frac{1}{\pi} & x_1^2 + x_2^2 < 1 \\ 0 & \text{otherwise} \end{cases}$$

Show that  $X_1$  and  $X_2$  are uncorrelated but not independent.

7. Suppose that  $X = (X_1, X_2) \sim N(\mu, \Sigma)$  where

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix}$$

- (a) Show that in this case  $X_1, X_2$  uncorrelated implies they are also independent. It might help to know that the marginal distribution of  $X_1$  is normal with mean  $\mu_1$  and variance  $\sigma_1^2 = \sigma_{11}$  with a similar statement for  $X_2$ .
- (b) Suppose,  $X = (X_1, X_2, \dots, X_p) \sim N(\mu, \Sigma)$ . Show that if  $\Sigma_{ij}^{-1} = 0$  then  $X_i$  and  $X_j$  are conditionally independent given *all* the other variables. (One can also show the stated conditional independence assumption implies  $\Sigma_{ij}^{-1} = 0$  for when dealing with normal random variables).
- (c) Suppose  $X_1, X_2, \dots, X_n$  defined by  $X_1 = Z_1$  and  $X_i = X_{i-1} + Z_i$  for  $i = 2, \dots, n$  where  $Z_1, \dots, Z_n \stackrel{\text{iid}}{\sim} N(0, 1)$ . This is known as a random walk.
- Show that  $(X_1, \dots, X_n) \sim N(0, \Sigma)$  for some covariance matrix  $\Sigma$ .
  - Without computing argue which elements of  $\Sigma^{-1}$  are 0.

8. Suppose that  $X \sim \text{Unif}(0, 1)$ .

- (a) Show that the cdf of  $Y = \frac{-\log(X)}{\lambda} = 1 - e^{-\lambda y}$  and conclude that  $Y \sim \text{Exponential}(\lambda)$ .
- (b) Suppose that  $X_1, X_2, \dots \stackrel{\text{iid}}{\sim} \text{Exponential}(\lambda)$ . Interpret  $X_1$  as the time we need to wait for the first event,  $X_2$  as the *additional* time we wait for the 2nd event,  $X_3$  as the additional time for the 3rd event and so on. In R take  $\lambda = 3$  and simulate 100,000 times the number of events that occur in the interval  $[0, 1]$ . Compute the empirical pmf,  $\hat{p}(n)$  for  $n = 0, 1, \dots$  by

$$\hat{p}(n) = \frac{\# \text{ times experiments gives } n}{100000}$$

On the same plot compare these probabilities with the  $\text{Poisson}(\lambda)$  distribution. The process  $N(t)$  that counts the number of events that occur before time  $t$  is known as a Poisson process.

9. Let  $M$  be a discrete variable with probability mass function  $p(m)$ . Suppose  $X$  is a continuous random vector that *depends* on  $M$ . Specifically, assume  $X|M = m$  is  $N(\mu_m, \Sigma_m)$ .

- (a) Even though  $M$  is discrete and  $X$  is continuous we can still get the joint distribution of  $M, X$  as  $p(m, x) = p(m)p(x|m)$ . In this problem we'll use the term density for discrete, continuous and mixed cases. Compute
- The joint density  $p(m, x)$
  - The marginal density  $p(x)$ .
  - The conditional density  $p(m|x)$
- (b) Suppose that a population is composed of  $C$  classes  $1, \dots, C$ . Let  $p(i)$  be the proportion of the population in class  $i$  so  $p(1) + \dots + p(C) = 1$ . Suppose we take several measurements,  $X$ , on the population and that these measurements behave differently under the classes. We model the “class conditional” distribution of  $X$  as multivariate normal with mean  $\mu_i$  and covariance  $\Sigma_i$ .
- If we observe  $X = x$  give the conditional probabilities  $P(M = m|X = x)$ .
  - Suppose that you observe  $X = x$  and want to choose the value of  $m$ ,  $\hat{M}$ , that has the best chance of being right. What should  $\hat{M}$  be?



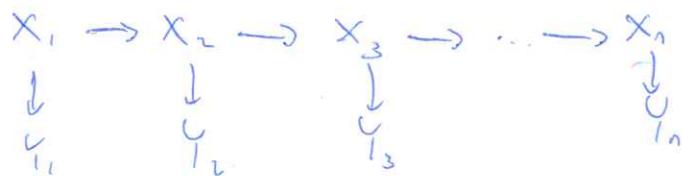
Important Graphical Models

①  $X_1 \dots X_n$  respect:  $(X_1) \rightarrow (X_2) \rightarrow (X_3) \rightarrow \dots \rightarrow (X_n)$

$$p(x_1 \dots x_n) = p(x_1) p(x_2 | x_1) p(x_3 | x_2) \dots p(x_n | x_{n-1})$$

$X_1 \dots X_n$  form Markov Chain

②  $X_1 \dots X_n, Y_1 \dots Y_n$  respect



③  $p(x_1 \dots x_n, y_1 \dots y_n) = p(x_1) p(x_2 | x_1) \dots p(x_n | x_{n-1}) p(y_1 | x_1) \dots p(y_n | x_n)$

$X, Y$  are hidden Markov Model.

## Classification

(25)

Have set of possible classes  $1 \dots K$

Eg: Regub/Dem/Indep or Spam/Not Spam or ...

Have data vector,  $X$ , of observations whose distrib. depends on class

Observe  $X=x$  and try to predict correct class.

## Modeling

Let  $c \in \{1, \dots, k\}$  be class, modeled as r.v. with dist  $\pi(c)$ .

Let  $\pi(x|c) = \pi_c(x)$  be class conditional distribution

(dist on  $X$  when class,  $c$ , known)

$\pi(c)$  is "prior" dist. on class [prior to seeing  $X$ ]

$\pi(c|x)$  is "posterior" " " " " [after " " " ]

By Bayes' Rule

$$\pi(c|x) = \frac{\pi(c)\pi(x|c)}{\pi(x)} = \frac{\pi(c)\pi(x|c)}{\sum_{c'} \pi(c')\pi(x|c')} \quad (\text{LorP})$$

Note:  $\pi(x) = \sum_c \pi(c)\pi(x|c)$  is mixture dist. with mixing weights  $\pi(c)$

## (26)

# Optimal Classification (Bayes Classifier)

Notation Will write  $\hat{c}$  for our estimate of true class  $c$ .

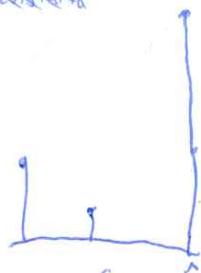
Suppose we take

$$\hat{c} = \arg \max_c q(c|x) = \arg \max_c \frac{q(c)q(x|c)}{\cancel{q(x)}} = \arg \max_c q(c)q(x|c)$$

The  $c$  that maximizes  $q(c|x)$

Fix  $x$  and consider errors

$$q(\text{error } | x) = \sum_{c \neq \hat{c}} q(c|x)$$



By choosing  $\hat{c}$  as most likely class we minimize sum

$\Rightarrow \hat{c} = \arg \max_c q(c|x)$  is optimal classifier.

Do: gaussian-classifier.r  
coin-flip-class-error.r

Comment: Hard to implement Bayes Class.

## Parameter Estimation

Bayes Classifier assumes we know class-conditional dists ( $q_c(x) = q(x|c)$ )  
though we don't. How to get these?

### Maximum Likelihood Estimation (MLE)

Have family of distributions indexed by unknown parameter  $G$ .

Ex:  $N(\theta, \sigma^2)$ , Exponential( $\theta$ ), Binomial( $n, \theta$ ), Poisson( $\theta$ )  
 $\sim N(\theta_1, \sigma_1^2)$ , ...

Observe random sample  $x_1, \dots, x_n \stackrel{iid}{\sim} p_G(x)$ . MLE,  $\hat{\theta} = \hat{\theta}_{MLE}(x_1, \dots, x_n)$

## Ex (Poisson)

Suppose  $x_1, \dots, x_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$ .

Fact

$$X \sim \text{Poisson}(\lambda) \Rightarrow EX = \lambda \quad [E(X) = \sum_{x=0}^{\infty} x \frac{e^{-\lambda} \lambda^x}{x!} = \sum_{x=1}^{\infty} e^{-\lambda} \frac{\lambda^x}{(x-1)!}$$

$$= \lambda \left( \sum_{x=0}^{\infty} \frac{e^{-\lambda} \lambda^x}{x!} \right) = \lambda 1 = \lambda]$$

$$\hat{\lambda}_{MLE}(x_1, \dots, x_n) = \arg \max_{\lambda} p_{\lambda}(x_1, \dots, x_n)$$

$$= \arg \max_{\lambda} \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{(x_i!)}$$

const wrt  $\lambda$

$$= \arg \max_{\lambda} e^{-n\lambda} \lambda^{\sum x_i}$$

Since  $\log$  is increasing fn  $\arg \max_{\lambda} p_{\lambda}(x_1, \dots, x_n) = \arg \max_{\lambda} \log p_{\lambda}(x_1, \dots, x_n)$

$$= \arg \max_{\lambda} -n\lambda + \sum x_i \log \lambda$$

To maximize set  $\frac{\partial}{\partial \lambda} = 0$  giving

the "sample mean"

$$0 = -n + \frac{\sum x_i}{\lambda} \Rightarrow \hat{\lambda} = \frac{\sum x_i}{n} \doteq \bar{x}$$

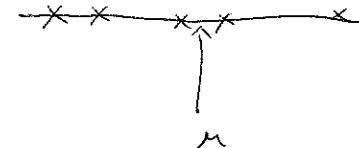
Comment

When dist parametrized by  $\theta = \text{mean}$  (as with Poisson) often happens

$$\text{that } \hat{\theta}_{MLE} = \bar{x}$$

Ex (Normal with known  $\sigma^2$ ) (28)

Have  $x_1, \dots, x_n \sim i.i.d N(\mu, \sigma^2)$  where  $\sigma^2$  known,  $\mu$  unknown.

$$\begin{aligned}
 \hat{\mu}_{MLE}(x_1, \dots, x_n) &= \arg \max_{\mu} q_{\mu}(x_1, \dots, x_n) = \arg \max_{\mu} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2} \\
 &= \arg \max_{\mu} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2} \\
 &= \arg \max_{\mu} \log e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2} \\
 &= \arg \max_{\mu} -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \\
 &= \arg \min_{\mu} \sum_{i=1}^n (x_i - \mu)^2
 \end{aligned}$$


$$=\bar{x}$$

Q: What is  $\#$  that minimizes sum sq diff from set of  $\#$ 's  $x_1, \dots, x_n$ ?

A:  $\bar{x}$

$$\begin{aligned}
 \sum_{i=1}^n (x_i - \mu)^2 &= \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - \mu)^2 \\
 &= \underbrace{\sum_{i=1}^n (x_i - \bar{x})^2}_{\text{const wrt } \mu} + \sum (\bar{x} - \mu)^2 + 2 \sum (x_i - \bar{x})(\bar{x} - \mu)
 \end{aligned}$$

To minimize  $\sum_{i=1}^n (x_i - \mu)^2$  take  $\mu = \bar{x}$ .

(29)

Ex  $\left(\text{Normal with known } \mu \text{ unknown } \sigma^2\right)$

$$\begin{aligned}\hat{\sigma}_{MLE}^2(x_1 - \bar{x}_n) &= \arg \max_{\sigma^2} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_i - \bar{x})^2} \\ &= \arg \max_{\sigma^2} \left(2\pi\sigma^2\right)^{-n/2} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \arg \max_{\sigma^2} -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum (x_i - \bar{x})^2\end{aligned}$$

Setting  $\frac{d}{d\sigma^2} = 0$  gives

$$0 = -\frac{n/2}{\sigma^2} + \frac{1}{2} \sum (x_i - \bar{x})^2 \cdot \frac{1}{(\sigma^2)^2}$$

$$\Rightarrow \hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum (x_i - \bar{x})^2 \quad \left( = \text{the "sample" avg sq dist. from } \mu \right)$$

## Bayesian Estimation

As before: have  $\gamma_{\theta}(x)$  for unknown  $\theta$ . Observe  $x$  ( $= x_1, \dots, x_n$ ) and want to estimate  $\theta$ .

Bayesian View Assume prior dist on  $\theta$ ,  $\gamma(\theta)$ . Then

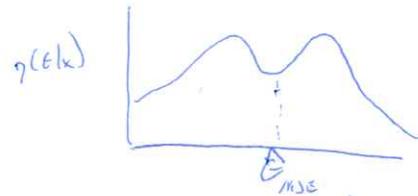
$$\gamma(\theta|x) = \frac{\gamma(\theta)\gamma(x|\theta)}{\gamma(x)}$$

Will base estimate for  $\theta$  on  $\gamma(\theta|x)$ . Eg

$$\hat{\theta}_{MAP} = \arg \max_{\theta} \gamma(\theta|x) = \arg \max_{\theta} \gamma(\theta) \gamma(x|\theta)$$

Maximum A Posteriori

$$\hat{\theta}_{MSE} = E(\theta|x)$$



(Recall that mean = expectation minimizes ave sq. error)  
 $\left[ \arg \min_c \int (x-c)^2 \gamma(x) dx = \int x \gamma(x) dx = E[X] \right]$

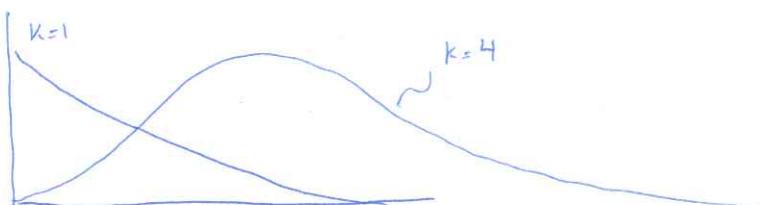
## Conjugate Priors

### Gamma Dist

$$X \sim \text{Gamma}(\theta, k) \quad \gamma(x) = \frac{e^{-x/\theta} x^{k-1}}{P(k) \theta^k}$$

$x > 0$

$$P(k) = \int_0^\infty e^{-x} x^{k-1} dx$$



Interesting interaction between Gamma and Poisson.

(31)

Observe  $X \sim \text{Poisson}(\lambda)$  and want to estimate  $\lambda$ .

Assume prior on  $\lambda$ ,  $\lambda \sim \text{Gamma}(\theta, k)$

$$\begin{aligned} p(\lambda|x) &= \frac{p(\lambda)p(x|\lambda)}{p(x)} = \underbrace{\frac{e^{-\lambda/\theta}\lambda^{k-1}}{r(k)\theta^k} \frac{e^{-\lambda}\lambda^x}{x!} \frac{1}{p(x)}}_{\text{denom is const. wrt } \lambda} \\ &= c e^{-\lambda/\theta} \lambda^{x+k-1} \end{aligned}$$

$$\Rightarrow \lambda|x \sim \text{Gamma}\left(\frac{1}{1/\theta}, x+k\right)$$

The parameter update after observation is

$$\theta_{\text{new}} = \frac{1}{1 + 1/\theta_{\text{old}}}$$

$$k_{\text{new}} = x + k_{\text{old}}$$

Note: We ended up where we started with a Gamma Dist on  $\lambda$

Could observe another  $x, x_2 \sim \text{Poisson}(\lambda)$ , then

$$\lambda|x_1, x_2 \sim \text{Gamma}\left(\frac{1}{2+1/\theta}, x_1+x_2+k\right)$$

$$\lambda|x_1, x_2, x_3 \sim \text{Gamma}\left(\frac{1}{3+1/\theta}, x_1+x_2+x_3+k\right)$$

$$\lambda|x_1, \dots, x_n \sim \text{Gamma}\left(\frac{1}{n+1/\theta}, \sum_{i=1}^n x_i + k\right)$$

Do earthquakes.

## Conjugate Prior - Cont

Have  $x \sim N(\mu, \sigma^2)$  and want to estimate  $\mu$ .

Have prior dist. on  $\mu$ ,  $\mu \sim N(r, \rho^2)$

What is posterior on  $\mu$ ?

$$g(\mu | x) = \frac{g(\mu) g(x|\mu)}{g(x)} = c e^{-\frac{1}{2\rho^2}(\mu-r)^2} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

$$= c' e^{q(\mu)}$$

$$\text{where } q(\mu) = a\mu^2 + b\mu + c = \frac{1}{\rho^2}(\mu-d)^2 + r$$

By completing the square get

$$\therefore \mu | x \sim N\left(\frac{\frac{r}{\rho^2} + \frac{x}{\sigma^2}}{\frac{1}{\rho^2} + \frac{1}{\sigma^2}}, \frac{1}{\frac{1}{\rho^2} + \frac{1}{\sigma^2}}\right)$$

Said another way, let

$\rho_{\text{old}}$   $\rho_{\text{new}}$  be original (prior) params for  $\mu$

$\rho_{\text{new}}$  " new (posterior) " " "

$$\rho_{\text{new}} = \frac{\frac{\rho_{\text{old}}}{\rho_{\text{old}}^2} + \frac{x}{\sigma^2}}{\frac{1}{\rho_{\text{old}}^2} + \frac{1}{\sigma^2}}$$

As before can use this update eqn. to revise our knowledge of  $\mu$  with sequence  $x_1, \dots, x_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ .

$$\rho_{\text{new}} = \frac{1}{\frac{1}{\rho_{\text{old}}^2} + \frac{1}{\sigma^2}}$$

Continued in HW.

# Generative vs. Discriminative Classification

Bishop Ch. 4

(33)

In generative model have  $p(c)$ ,  $p(x|c)$  and compute

$$p(c|x) = \frac{p(c)p(x|c)}{\sum_c p(c)p(x|c)}$$

posterior is everything we need to know for inferences about  $c$ .

Typically involves

- ① Estimating  $p(c)$   $\sim \text{Eg } N(\mu_c, \Sigma_c)$
- ② Estimating parameters for  $p(x|c)$

Discriminative classification models  $p(c|x)$  directly (and parametrically)

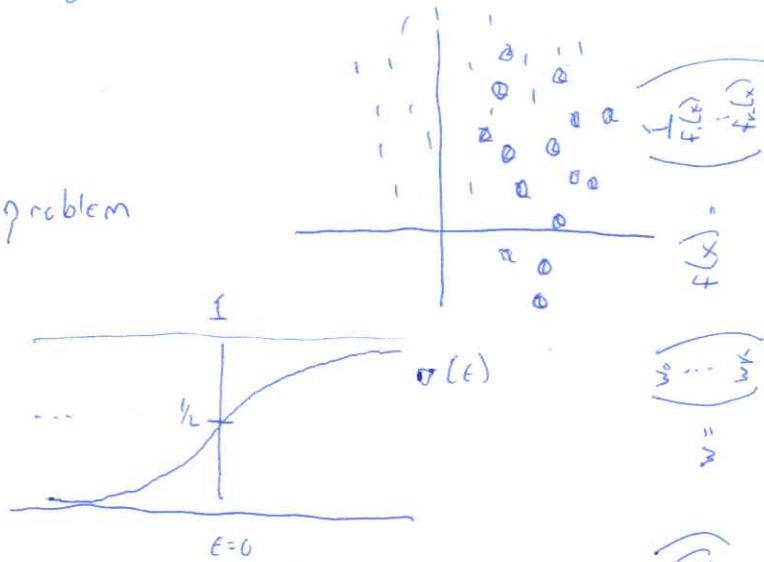
Most famous example: Logistic regression.

## Logistic Regression

Consider 2-class classification problem

Define logistic sigmoid  $f_n$

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$



Important properties

- ①  $\sigma(t) + \sigma(-t) = 1$
- ②  $\sigma'(t) = \sigma(t)(1 - \sigma(t))$

$f(x)$  = "feature vector" contains whatever observables, or functions of observables we think are useful,

## Model

Have observation  $x = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$

Compute features:  $f_1(x), \dots, f_k(x)$

## Model

$$p(c=1|x) = \sigma(w_0 1 + w_1 f_1(x) + \dots + w_k f_k(x))$$

Ex

2 classes  $c_1 = 1$  has disease

$c_2 = 0$  no disease

Measure 3 important biomarkers  $x_1, x_2, x_3$ .

Believe interactions between vbles also important so let

$$f(x) = (x_1, x_2, x_3, x_1x_2, x_2x_3, x_1x_3, x_1x_2x_3, 1)$$

$$p(c=1 | x) = \sigma(w^t f(x)) = \sigma(w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_1 x_2 + \dots)$$

b/w 2 sols

Notes,

Practice problems

Change due date

Logistic assumes  $p(c=1|x) = \sigma(w^t f(x))$

but how to find  $w = (w_0 \dots w_k)^t$ ?

$$\textcircled{1} \quad \text{Since } p(c=1|x) = \sigma(w^t f(x))$$

$$p(c=0|x) = 1 - \sigma(w^t f(x))$$

\textcircled{2} If  $c \in \{0, 1\}$  could write as

$$p(c|x) = \sigma(w^t f(x))^c (1 - \sigma(w^t f(x)))^{1-c}$$

\textcircled{3} Suppose have sample  $(c_1, x_1), (c_2, x_2) \dots (c_n, x_n)$

Assume the ~~samples~~<sup>observations</sup> independent.

$$p(c_1 \dots c_n | x_1 \dots x_n) = \prod_{i=1}^n p(c_i | x_i) = \prod_{i=1}^n \sigma(w^t f(x_i))^c_i (1 - \sigma(w^t f(x_i)))^{1-c_i}$$

choose  $w$  by max like. That is

$$\hat{w} = \arg \max_w \prod_{i=1}^n p(c_i | x_i) = \arg \max_w \prod_{i=1}^n \sigma(w^t f(x_i))^c_i (1 - \sigma(w^t f(x_i)))^{1-c_i}$$

Can't optimize in closed form as in previous examples, but can optimize numerically.

Recall

Want to maximize/minimize

$$\nabla g(x)$$



$$\nabla g(x)$$

$$g(x) : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\begin{pmatrix} \frac{\partial g(x)}{\partial x_1} \\ \vdots \\ \frac{\partial g(x)}{\partial x_n} \end{pmatrix}$$

is gradient.

4)  $\nabla g(x) = 0$  at cusp or saddle

5)  $\nabla g(x)$  "points" in

direction of maximal (local)

2)  $-\nabla g(x)$  " " " decrease

3)  $\nabla g(x)$  orthogonal to level sets

2)  $\{x : g(x) = c\}$  = "Level sets of g"

## Back to Logistic Regression

$$g(c_1 \dots c_n | x_1 \dots x_n) = \prod_{i=1}^n \sigma(w^T f(x_i))^{c_i} (1 - \sigma(w^T f(x_i)))^{1-c_i}$$

$$\log g(c_1 \dots c_n | x_1 \dots x_n) = \sum_{i=1}^n c_i \log \sigma(w^T f(x_i)) + (1-c_i) \log (1 - \sigma(w^T f(x_i)))$$

$$\frac{\partial \log g(c_1 \dots c_n | x_1 \dots x_n)}{\partial w_k} = \sum_{i=1}^n \frac{c_i \sigma'(w^T f(x_i)) f_k(x_i)}{\sigma(w^T f(x_i))} - \frac{(1-c_i) \sigma'(w^T f(x_i)) f_k(x_i)}{1 - \sigma(w^T f(x_i))}$$

$$\sigma'(t) = \sigma(t)(1-\sigma(t))$$

$$= \sum c_i (1 - \sigma(w^T f(x_i)) f_k(x_i)) - (-c_i) \sigma(w^T f(x_i)) f_k(x_i)$$

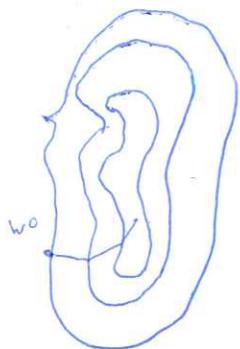
$$= \sum_{i=1}^n (c_i - \sigma(w^T f(x_i))) f_k(x_i)$$

$$\Rightarrow \nabla \log g(c_1 \dots c_n | x_1 \dots x_n) = \left( \sum_{i=1}^n (c_i - \sigma(w^T f(x_i))) \begin{pmatrix} f(x_i) \\ \vdots \end{pmatrix} \right)$$

Can optimize  $\log g(c_1 \dots c_n | x_1 \dots x_n)$  by gradient descent

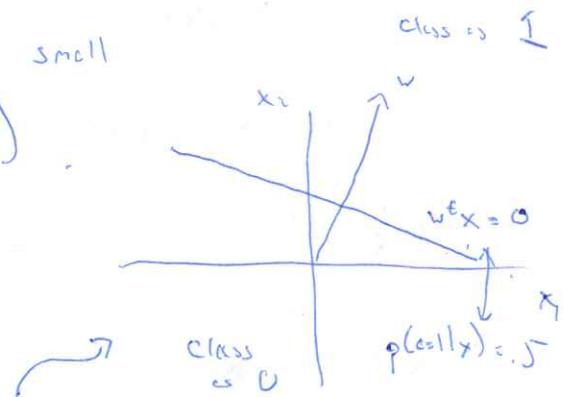
Choose initial guess for  $w, w^0$

$$w^{\text{new}} = w^{\text{old}} + \epsilon \nabla \log g(c_1 \dots c_n | x_1 \dots x_n) \quad \text{for some stepsize } \epsilon$$



Continue until  $\nabla$  sufficiently small  
(at optimal value  $\nabla = 0$ )

- ① Do logistic-model.r
- ② Logistic-grad.descent.r



## Newton - Raphson Method

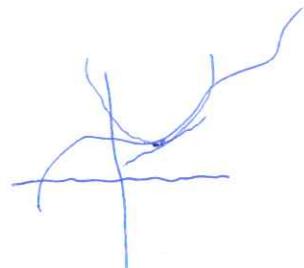
(37)

NR improves on simple-minded gradient descent.  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

$$x^{new} = x^{old} + \varepsilon \nabla f(x^{old})$$

Recall Taylor series in 1-d  $f: \mathbb{R} \rightarrow \mathbb{R}$

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2} f''(x_0)(x - x_0)^2$$



The multivariate T.S. is  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

$$f(x) \approx f(x_0) + \underbrace{\nabla f(x_0)}_{\text{dot grad.}} \cdot (x - x_0) + \frac{1}{2} (x - x_0)^T H(x - x_0)$$

where  $H$  is "Hessian" matrix  $H_{ij} = \frac{\partial^2 f(x_0)}{\partial x_i \partial x_j}$

Interesting fact on quadratic forms

$$q(x) = \frac{1}{2} (x - x_0)^T H(x - x_0) \quad \nabla q(x) = H(x - x_0)$$

Take gradient of approx above to get

$$\nabla f(x_0) \approx \nabla f(x_0) + H(x - x_0)$$

To find <sup>maximal</sup> ~~optimal~~ value of  $f(x)$  set  ~~$\nabla f(x) = 0$~~  to get

$$0 = \nabla f(x) \approx \nabla f(x_0) + H(x - x_0)$$

Solving for  $x$  gives

$$-H \nabla f(x_0) = x - x_0 \implies x^{opt} = x_0 - H \nabla f(x_0)$$

$$NR: \quad \cancel{x^{old}} \quad x^{new} = x^{old} - H \nabla f(x^{old})$$

# Logistic Regression with Newton - Raphson

38

## Review

$$\sigma(t) = \frac{1}{1+e^{-t}} \quad \text{Model} \quad \gamma(c=1|x) = \sigma(w^t x)$$

Have data vectors  $x_1, x_2, \dots, x_n$  We know

$$\frac{\partial \log \gamma(c_i, \dots, c_n | x_1, \dots, x_n)}{\partial w_j} = \sum_{i=1}^n (c_i - \sigma(w^t x_i)) x_{ij} \Leftrightarrow \nabla \log \gamma(c_i, \dots, c_n | x_1, \dots, x_n) = \sum (c_i - \sigma(w^t x_i)) x_i$$

Writing

$$X = \begin{pmatrix} x_1 & \cdots & x_K \\ x_{11} & \cdots & x_{1K} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{nK} \end{pmatrix} \quad c = \begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix} \quad \gamma = \begin{pmatrix} \sigma(w^t x_1) \\ \vdots \\ \sigma(w^t x_n) \end{pmatrix}$$

$$\nabla \log \gamma(c_i, \dots, c_n | x_1, \dots, x_n) = X^t (c - \gamma)$$

$$\frac{\partial^2 \log \gamma(c_i, \dots, c_n | x_1, \dots, x_n)}{\partial w_j \partial w_l} = - \sum_{i=1}^n \sigma(w^t x_i) (1 - \sigma(w^t x_i)) x_j x_{il} \\ = - (X^t C X)_{jl}$$

$$C = \begin{pmatrix} \sigma(w^t x_1) (1 - \sigma(w^t x_1)) & & \\ & \ddots & \\ & & \sigma(w^t x_n) (1 - \sigma(w^t x_n)) \end{pmatrix} \quad \text{i.e. } -X^t C X \text{ is Hessian}$$

$$NR: \quad w^{new} = w^{old} + H^{-1} \nabla f(w^{old}) = w^{old} + (X^t C X)^{-1} X^t (c - \gamma)$$

NB,  $C$  and  $\gamma$  depend on  $w^{old}$  so must be recomputed each iter

