

# Introduction to Data Science

(AOS)

## Chapter 1: Probability (mathematical language to quantify uncertainty)

sample space  $(\Omega)$  is the set of possible outcomes of an experiment. Subsets of  $\Omega$  are called events. Points of  $\omega$  in  $\Omega$  are called sample outcomes / elements

ex 1. If we toss a coin forever, then the sample space is the infinite set

$$\Omega = \{ \omega = (\omega_1, \omega_2, \omega_3, \dots) : \omega_i \in \{H, T\} \}$$

Let  $E$  be the event that the first head appears on the third toss. Then

$$E = \{ (\omega_1, \omega_2, \omega_3, \dots) : \omega_1 = T, \omega_2 = T, \omega_3 = H, \omega_i \in \{H, T\} \text{ for } i \geq 3 \}$$

Given an event  $A$  let  $A^c = \{ \omega \in \Omega : \omega \notin A \}$  : complement of  $A$  (not  $A$ )

true event  
(always true)

$$\Omega^c = \emptyset \text{ (empty set)} \rightarrow \text{null event - always false}$$

union of events  $A$  and  $B$   $A \cup B = \{ \omega \in \Omega : \omega \in A \text{ or } \omega \in B \text{ or } \omega \in \text{both} \}$

"A or B" If  $A$  is a sequence of sets then  $\bigcup_{i=1}^{\infty} A_i = \{ \omega \in \Omega : \omega \in A_i \text{ for at least one } i \}$

Intersection "A and B"  $A \cap B = \{ \omega \in \Omega : \omega \in A \text{ and } \omega \in B \}$

if  $A_1, A_2, \dots$  is a sequence of sets  $\bigcap_{i=1}^{\infty} A_i = \{ \omega \in \Omega : \omega \in A_i \text{ for all } i \}$

The set difference is defined by  $A \setminus B = \{ \omega : \omega \in A, \omega \notin B \}$

If every element of  $A$  contained in  $B$   $A \subset B$  or  $B \supset A$

$A_1, A_2, \dots$  are disjoint / mutually exclusive if  $A_i \cap A_j = \emptyset$  whenever  $i \neq j$

eg.  $A_1 = [0, 1)$   $A_2 = [1, 2)$   $A_3 = [2, 3)$

A partition of  $\Omega$  is a sequence of disjoint sets  $A_1, A_2, \dots$  such that  $\bigcup_{i=1}^{\infty} A_i = \Omega$

Event  $A$  the indicator function of  $A$

$$I_A(\omega) = I(\omega \in A) = \begin{cases} 1, & \text{if } \omega \in A \\ 0, & \text{if } \omega \notin A \end{cases}$$

A sequence sets  $A_1, A_2, \dots$  is monotone increasing if  $A_1 \subset A_2 \subset \dots$  and we

define  $\lim_{n \rightarrow \infty} A_n = \bigcup_{i=1}^{\infty} A_i$ , if ——— decreasing if  $A_1 \supset A_2 \supset \dots$  and we

define  $\lim_{n \rightarrow \infty} A_n = \bigcap_{i=1}^{\infty} A_i$

A function  $P$  that assigns a real number  $P(A)$  to each event  $A$  is a probability distribution or a probability measure if it satisfies the following axioms.

Axiom 1:  $P(A) \geq 0$  for every  $A$

Axiom 2:  $P(\Omega) = 1$

Axiom 3: If  $A_1, A_2, \dots$  are disjoint then  $P(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$

Derived Properties from Axioms:

$$P(\emptyset) = 0 \text{ (use Axiom 1)}$$

$$P(A^c) = 1 - P(A) \text{ (use Axiom 2)}$$

$$A \subset B \Rightarrow P(A) \leq P(B)$$

$$A \cap B = \emptyset \Rightarrow P(A \cup B) = P(A) + P(B)$$

$$0 \leq P(A) \leq 1$$



$\sum_{j \in A \cap B} x_j = \sum_{j \in A} x_j + \sum_{j \in B} x_j - \sum_{j \in A \cap B} x_j \Leftrightarrow$  implications  $\Rightarrow$  individuals needs

Proof  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ . We can start from the following events:  
 $= P((A \cap B^c) \cup (A \cap B) \cup (A^c \cap B)) = P(A \cap B^c) + P(A \cap B) + P(A^c \cap B)$   
 $= P(A \cap B^c) + P(A \cap B) + P(A^c \cap B) + P(A \cap B) - P(A \cap B)$   
 $= P((A \cap B^c) \cup (A \cap B)) + P((A^c \cap B) \cup (A \cap B)) - P(A \cap B)$   
 $= P(A) + P(B) - P(A \cap B)$

Theorem (Continuity of Probabilities) If  $A_n \rightarrow A$  then  $P(A_n) \rightarrow P(A)$  as  $n \rightarrow \infty$

Proof Suppose  $A_n$  is monotone increasing so that  $A_1 \subset A_2 \subset \dots$ . Let  $A = \lim_{n \rightarrow \infty} A_n = \bigcup_{i=1}^{\infty} A_i$ . Define  $B_1 = A_1$ ,  $B_2 = \{w \in \Omega : w \in A_2, w \notin A_1\}$ ,  $B_3 = \{w \in \Omega : w \in A_3, w \notin A_2, w \notin A_1\}$ ,  $B_1, B_2, B_3, \dots$  disjoint.  $A_n = \bigcup_{i=1}^n A_i = \bigcup_{i=1}^n B_i$  for each  $n$  and  $\bigcup_{i=1}^{\infty} B_i = \bigcup_{i=1}^{\infty} A_i$  from Axiom 3.  $P(A_n) = P(\bigcup_{i=1}^n B_i) = \sum_{i=1}^n P(B_i)$

$$\Rightarrow \lim_{n \rightarrow \infty} P(A_n) = \lim_{n \rightarrow \infty} \sum_{i=1}^n P(B_i) = \sum_{i=1}^{\infty} P(B_i) = P(\bigcup_{i=1}^{\infty} B_i) = P(A)$$

$\rightarrow$  If  $\Omega$  is finite and if each outcome is equally likely then  $P(A) = \frac{|A|}{|\Omega|}$  and it's called uniform probability distribution.

Given  $n$  objects, the number of ways of ordering these is  $n! = (n-1)(n-2) \dots 3 \cdot 2 \cdot 1$   
 $\binom{n}{k} = \frac{n!}{k!(n-k)!}$  "n choose k"  $\binom{n}{0} = \binom{n}{n} = 1$

Independent Events Two events  $A$  and  $B$  are independent if  $P(AB) = P(A)P(B)$ .  $A \perp B$ , when are not we write  $A \not\perp B$ .

Suppose  $A$  and  $B$  are ~~independent~~ disjoint each with a +ve probabilities. Can they be independent?  
 No.  $P(A) \cdot P(B) > 0$  yet  $P(AB) = P(\emptyset) = 0$

Conditional Probability: If  $P(B) > 0$  then  $P(A|B) = \frac{P(AB)}{P(B)}$

\* If  $A$  and  $B$  independent events  $\Rightarrow P(A|B) = P(A)$

$$P(AB) = P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$

1) Bayes' Theorem: Let  $A_1, \dots, A_k$  be a partition of  $\Omega$  such that  $P(A_i) > 0$  for each  $i$ . If  $P(B) > 0$  then for each  $i = 1, \dots, k$ :

$$P(A_i|B) = \frac{P(B|A_i) \cdot P(A_i)}{\sum_{j=1}^k P(B|A_j) \cdot P(A_j)}$$

2) The Law of Total Probability. Let  $A_1, \dots, A_k$  be a partition of  $\Omega$  (mutually exclusive & exhaustive events). Then for any event  $B$ :  
 $P(B) = \sum_{i=1}^k P(B|A_i) P(A_i)$

Proof 2:  $C_j = B \cap A_j$  and  $C_1, \dots, C_k$  disjoint and  $B = \bigcup_{j=1}^k C_j$   
 $P(B) = \sum P(C_j) = \sum P(B \cap A_j) = \sum P(B|A_j) P(A_j)$

$$\text{Proof 1: } P(A_i|B) = \frac{P(A_i \cap B)}{P(B)} = \frac{P(B|A_i) \cdot P(A_i)}{P(B)} = \frac{P(B|A_i) \cdot P(A_i)}{\sum_{j=1}^k P(B|A_j) \cdot P(A_j)}$$



$$P(A) = P(A \cap B) + P(A \cap B^c) = P(A|B)P(B) + P(A|B^c)P(B^c)$$



A RV is a mapping  $X: \Omega \rightarrow \mathbb{R}$  that assigns a real number  $X(\omega)$  to each outcome  $\omega$ .

Let  $\Omega = \{(x, y); x^2 + y^2 \leq 1\}$ , w. typical outcome  $= (x, y)$

RV:  $X(\omega) = x$ ,  $Y(\omega) = y$ ,  $Z(\omega) = x + y$ ,  $W(\omega) = \sqrt{x^2 + y^2}$

$X$  is a subset of  $A$  define  $X^{-1}(A) = \{\omega \in \Omega : X(\omega) \in A\}$  and

$$P(X \in A) = P(X^{-1}(A)) = P(\{\omega \in \Omega; X(\omega) \in A\})$$


$$P(X=x) = P(X^{-1}(x)) = P(\{\omega \in \Omega; X(\omega) = x\})$$

12V  $\rightarrow$  particular valve

Ex. Flip a coin twice and let  $X$  be the number of heads. Then  $P(X=0) = P(\{TT\}) = 1/4$ ,  $P(X=1) = P(\{HT, TH\}) = 1/2$  and  $P(X=2) = P(\{HH\}) = 1/4$ .

$\omega$	$IP(\omega)$	$X(\omega)$	}	$x$	$P(X=x)$	Discrete	Probable Mass Function (PMF)		
TT	1/4	0				0	1/4	Continuous	Probability density function (PDF)
TH	1/4	1				1	1/2		
HT	1/4	1				2	1/4		
HH	1/4	2							

cumulative distribution function (CDF)

\* The cumulative distribution function (CDF) is the function  $F_X: \mathbb{R} \rightarrow [0, 1]$  defined by  $F_X(x) = \mathbb{P}(X \leq x)$   {right continuous, non-decreasing}

~~Ex. flip a coin twice and X~~

$$F_X(x) = \begin{cases} 0 & , x < 0 \\ 1/4 & , 0 \leq x < 1 \\ 3/4 & , 1 \leq x < 2 \\ 1 & , x \geq 2 \end{cases}$$

Let  $X$  have CDF  $F$  and  $Y$  have GDF  $G$ . If  $F(x) = G(x)$  for all  $x$  then  $P(X \in A) = P(Y \in A)$  for all  $A$ .

A function  $F$  mapping to  $[0, 1]$  is a CDF for some probability  $P$  if and only if  $F$  satisfies:

- $F$  is non-decreasing  $x_1 < x_2$  implies  $F(x_1) \leq F(x_2)$

4)  $f$  is normalised  $\lim_{x \rightarrow -\infty} f(x) = 0$  and  $\lim_{x \rightarrow \infty} f(x) = 1$

iii)  $F$  is right-continuous  $F(x) = F(x^+)$  for all  $x$  where  $F(x^+) = \lim_{\substack{y \rightarrow x \\ y > x}} F(y)$

$X$  is discrete if it takes countably many values  $\{x_1, x_2, \dots\}$ . We define the probability function or probability mass function for  $X$  by  $f_X(x) = P(X=x)$ .

Thus,  $f_X(x) \geq 0$  for all  $x \in \mathbb{R}$  and  $\sum_i f(x_i) = 1$ . The CDF of  $X$ :

$$F_X(x) = P(X \leq x) = \sum_{x_i \leq x} f_X(x_i)$$



$$f_X(x) = \begin{cases} 1/4, & x=0 \\ 1/2, & x=1 \\ 1/4, & x=2 \\ 0, & \text{otherwise} \end{cases}$$



A RV  $X$  is continuous if there exists a function  $f_X$  such that  $f_X(x) \geq 0$  for all  $x$   $\int_{-\infty}^{\infty} f_X(x) dx = 1$  and for every  $a \leq b$

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

The function  $f_X$  is called the probability density function (PDF):

$$F_X(x) = \int_{-\infty}^x f_X(t) dt \quad \text{and} \quad f_X(x) = F'_X(x) \text{ at all points } x \text{ at}$$

which  $F_X$  is differentiable

Ex. 1 Let say  $X$  has PDF:

$$f_X(x) = \begin{cases} 1 & \text{for } 0 \leq x \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

$f_X(x) \geq 0$  and  $\int_{-\infty}^{\infty} f_X(x) dx = 1$ . A RV with that density is said to have a Uniform(0, 1) distribution. The CDF is given by

$$F_X(x) = \begin{cases} 0, & x < 0 \\ x, & 0 \leq x \leq 1 \\ 1, & x > 1 \end{cases}$$



Ex. 2:  $X$  has PDF:

$$f_X(x) = \begin{cases} 0 & \text{for } x < 0 \\ \frac{1}{(1+x)^2} & \text{otherwise} \end{cases}$$

since  $\int_{-\infty}^{\infty} f_X(x) dx = 1 \rightarrow$  well-defined PDF

\* If  $X$  continuous  $\rightarrow P(X=x) = 0$  for every  $x \neq P(x) = P(X=x)$

\* A PDF can be bigger than 1  $\neq$  mass function. A PDF can be unbounded

Ex. 3: Let  $f_X(x) = \begin{cases} 0, & \text{for } x < 0 \\ \frac{1}{1+x}, & \text{otherwise} \end{cases}$  this is not a PDF since  $\int_0^{\infty} \frac{1}{1+x} dx = \int_1^{\infty} \frac{1}{u} du = \log(\infty) = \infty$

Lemma 2.15 Let  $F$  be the CDF for a RV  $X$ . Then:

$$1. P(X=x) = F(x) - F(x^-) \text{ where } F(x^-) = \lim_{y \uparrow x} F(y)$$

$$2. P(x < X \leq y) = F(y) - F(x)$$

$$3. P(X > x) = 1 - F(x)$$

$$4. \text{If } X \text{ continuous: } F(b) - F(a) = P(a < X < b) = P(a \leq X < b) = P(a < X \leq b) = P(a \leq X \leq b)$$

Let  $X$  be a RV with CDF  $F$ . The inverse CDF <sup>quantile function</sup> is denoted by:

$F^{-1}(q) = \inf \{x : F(x) \geq q\}$  for  $q \in [0, 1]$ . If  $F$  is strictly increasing & continuous then  $F^{-1}(q)$  is the unique real number  $x$  such that  $F(x) = q$

$$1^{st} \text{ quantile: } F^{-1}(1/4)$$

$$3^{rd} \text{ quantile: } F^{-1}(3/4)$$

$$2^{nd} = 1 - \text{median: } F^{-1}(1/2)$$



- \* Two variables are equal in distribution  $X \stackrel{d}{=} Y$  if  $F_X(x) = F_Y(x)$  for all  $x$  but  $X$  and  $Y$  not equal eg  $P(Y=1) = P(Y=-1) = 1/2$ ...
- $\hookrightarrow P(X=Y) = 0$

Some Important Discrete RV.

1. Point Mass Distribution:  $X \sim \delta_a$  if  $P(X=a) = 1$   $F_X(x) = \begin{cases} 0, & x < a \\ 1, & x \geq a \end{cases}$   
PMF:  $f(x) = 1$  for  $x=a$  and 0 otherwise.
2. Discrete Uniform Distribution:  $k \geq 1$  (given integer), then  $X$  has PDF:  
 $f(x) = \begin{cases} 1/k, & \text{for } x=1, \dots, k \\ 0, & \text{otherwise} \end{cases}$
3. Bernoulli Distribution:  $P(X=1)=p$  and  $P(X=0)=1-p$ ,  $p \in [0,1]$   
 $X \sim \text{Bernoulli}(p)$  Probability function:  $f(x) = p^x (1-p)^{1-x}$ ,  $x \in \{0,1\}$
4. Binomial Distribution:  $P(\text{heads}) = p$  flip it  $n$  times,  $X = \#$  of heads  
 $f(x) = P(X=x)$  mass function  
 $f(x) = \begin{cases} \binom{n}{x} p^x (1-p)^{n-x}, & \text{for } x=0, \dots, n \\ 0, & \text{otherwise} \end{cases}$   $X \sim \text{Binomial}(n, p)$   
and you can add binomials
5. Geometric Distribution:  $X \sim \text{Geom}(p)$ ,  $p \in (0,1)$   
 $P(X=k) = p(1-p)^{k-1}$ ,  $k \geq 1$   
 $\sum_{k=1}^{\infty} P(X=k) = p \sum_{k=1}^{\infty} (1-p)^{k-1} = \frac{p}{1-(1-p)} = 1$   
" $X = \#$  of flips until the first success (heads)"
6. Poisson Distribution  $X \sim \text{Poisson}(\lambda)$   
 $f(x) = e^{-\lambda} \cdot \frac{\lambda^x}{x!}$ ,  $x \geq 0$   $\sum_{x=0}^{\infty} f(x) = e^{-\lambda} \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} = e^{-\lambda} e^{\lambda} = 1$   
eg. use for rare events like radioactive decay or traffic accidents  $\rightarrow$  you can add

Some important ~~Discrete~~ Continuous RV:

1. Uniform Distribution  $X \sim \text{uniform}(a,b)$  if:  
 $f(x) = \begin{cases} \frac{1}{b-a} & \text{for } x \in [a,b] \\ 0 & \text{otherwise} \end{cases}$  where  $a < b$  Distribution function is:  
 $f(x) = \begin{cases} 0, & x < a \\ \frac{x-a}{b-a}, & x \in [a,b] \\ 1, & x > b \end{cases}$
2. Normal (Gaussian)  $X \sim N(\mu, \sigma^2)$   $f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2} (x-\mu)^2\right\}$ ,  $x \in \mathbb{R}$   
standard normal distribution  $\sim N(0, 1^2)$  PDF  $\rightarrow \phi(z)$  CDF  $\rightarrow \Phi(z)$   
i) If  $X \sim N(\mu, \sigma^2)$  then  $Z = (X-\mu)/\sigma \sim N(0,1)$   
ii) If  $Z \sim N(0,1)$  then  $X = \mu + \sigma Z \sim N(\mu, \sigma^2)$   
iii) if  $X_i \sim N(\mu_i, \sigma_i^2)$ ,  $i=1, \dots, n$  independent then  $\sum_{i=1}^n X_i \sim N\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2\right)$



$$P(a < X < b) = P\left(\frac{a-t}{\sigma} < Z < \frac{b-t}{\sigma}\right) = \Phi\left(\frac{b-t}{\sigma}\right) - \Phi\left(\frac{a-t}{\sigma}\right)$$

Ex.  $X \sim N(3, 5)$  find  $P(X > 1)$

$$1 - P(X < 1) = 1 - P\left(Z < \frac{1-3}{\sqrt{5}}\right) = 1 - \Phi(-0.8944) = 0.81$$

$$\text{now find } q = \Phi^{-1}(0.2) \quad P(X < q) = 0.2$$

$$= P\left(Z < \frac{q-t}{\sigma}\right) = \Phi\left(\frac{q-t}{\sigma}\right) = \Phi(-0.8416) = 0.2$$

$$= -0.8416 = \frac{q-t}{\sigma} = \frac{q-3}{\sqrt{5}} \quad \Rightarrow q = 3 - 0.8416\sqrt{5} = 1.1181$$

3. Exponential Distribution.  $X \sim \text{Exp}(\theta)$   $f(x) = \frac{1}{\theta} e^{-x/\theta}, x > 0, \theta > 0$

Used for the lifetimes of electronic components and the waiting times between rare events

4. Gamma Distribution:  $a > 0 \quad \Gamma(a) = \int_0^\infty y^{a-1} e^{-y} dy$

$$X \sim \text{Gamma}(a, \theta) \quad \text{if} \quad f(x) = \frac{1}{\theta^a \Gamma(a)} x^{a-1} e^{-x/\theta}, x > 0, a, \theta > 0$$

Exponential distribution is Gamma(1,  $\theta$ )

5. Beta Distribution:  $a > 0, b > 0 \quad X \sim \text{Beta}(a, b)$  if  $f(x) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}, 0 < x < 1$

6. t and Cauchy Distribution  $X \sim t_\nu$  (similar to normal but with heavier tails)

normal corresponds to t with  $\nu = \infty$

$$\text{special case when } \nu = 1 \quad f(x) = \frac{1}{\pi(1+x^2)}$$

7.  $\chi^2$  Distribution:  $X \sim \chi_p^2$

if  $Z_1, \dots, Z_p$  independent standard Normal RV then  $\sum_{i=1}^p Z_i^2 \sim \chi_p^2$

Bivariate Distribution

$X, Y$  discrete - Joint Mass function  $f(x, y) = P(X=x, Y=y)$

$\rightarrow X, Y$  continuous we call  $f(x, y)$  a PDF if:

$$i) f(x, y) \geq 0 \text{ for all } (x, y)$$

$$ii) \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy = 1$$

$$iii) \text{ for any set } A \subset \mathbb{R} \times \mathbb{R}, P((X, Y) \in A) = \int_A f(x, y) dx dy$$

In both discrete & continuous the Joint CDF  $F_{XY}(x, y) = P(X \leq x, Y \leq y)$

Ex. Let  $(X, Y)$  uniform

$$f(x, y) = \begin{cases} 1 & \text{if } 0 \leq x \leq 1, 0 \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

$$P(X < 1/2, Y < 1/2) \text{ (here } A = \{X < 1/2, Y < 1/2\})$$

integrating  $f$  over this subset (compute area of  $A$ ) =  $1/4, P(X < 1/2, Y < 1/2) = 1/4$

Ex. Let  $(X, Y)$  have a density

$$f(x, y) = \begin{cases} x+y & \text{if } 0 \leq x \leq 1, 0 \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases} \quad \text{then}$$

$$\int_0^1 \int_0^1 (x+y) dx dy = \int_0^1 \left[ \int_0^1 x dx \right] dy + \int_0^1 \left[ \int_0^1 y dx \right] dy$$

$$= \int_0^1 \frac{1}{2} dy + \int_0^1 y dy = \frac{1}{2} + \frac{1}{2} = 1 \Rightarrow \text{verifies it is a PDF}$$



If  $X_1, X_2, \dots, X_n$  continuous RV then the marginal pdf of  $X_1$  is found by integrating  $X_2, X_3, \dots, X_n$  out of the joint pdf.

## Marginal Distributions $\rightarrow$

If  $(X, Y)$  have joint distribution with mass function  $f_{X,Y}$  then marginal mass function for  $X$  is -

$$f_X(x) = P(X=x) = \sum_y P(X=x, Y=y) = \sum_y f(x,y) \text{ and the same for } y$$

Ex. The marginal distributions for  $X$  corresponds to the row totals and for  $Y$  the column totals

	$Y=0$	$Y=1$	
$X=0$	1/10	2/10	3/10
$X=1$	3/10	4/10	7/10
	4/10	6/10	1

$$f_X(\cdot) = 3/10, f_X(1) = 7/10$$

For continuous RV the marginal densities are (same x,y) -

$$f_X(x) = \int f(x,y) dy, f_Y(y) = \int f(x,y) dx$$

Ex. ①  $f(x,y) = \begin{cases} x+y & \text{if } 0 \leq x \leq 1, 0 \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$

then

$$f_Y(y) = \int_0^1 (x+y) dx = \int_0^1 x dx + \int_0^1 y dx = \frac{1}{2} + y$$

②  $f(x,y) = \begin{cases} \frac{21}{4} x^2 y & \text{if } x^2 \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$

then  $f_X(x) = \int f(x,y) dy = \frac{21}{4} x^2 \int_{x^2}^1 y dy = \frac{21}{8} x^2 (1-x^2)$  for  $-1 \leq x \leq 1$   
 $\neq$  and  $f_X(x) = 0$  otherwise

Independent RV. Two RVs  $X$  and  $Y$  are independent if for every  $A$  and  $B$   $P(X \in A, Y \in B) = P(X \in A) P(Y \in B) \rightarrow X \perp Y$  otherwise  $X \text{ and } Y$

Ex. 

	$Y=0$	$Y=1$	
$X=0$	1/4	1/4	1/2
$X=1$	1/4	1/4	1/2
	1/2	1/2	1

 $f_X(0) = f_X(1) = 1/2, f_Y(0) = f_Y(1) = 1/2$   
 $f_X(0) f_Y(0) = f(0,0), f_X(0) f_Y(1) = f(0,1)$   
 $f_X(1) f_Y(0) = f(1,0), f_X(1) f_Y(1) = f(1,1)$   
 $X, Y$  independent

Ex. 2  $X, Y$  independent and have same density.

$$f(x) = \begin{cases} 2x & \text{if } 0 \leq x \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

$P(X+Y \leq 1)$  - using independence, the joint density is

$$f(x,y) = f_X(x) f_Y(y) = \begin{cases} 4xy & \text{if } 0 \leq x \leq 1, 0 \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned} P(X+Y \leq 1) &= \int_{0 \leq x+y \leq 1} f(x,y) dy dx \\ &= 4 \int_0^1 x \left[ \int_0^{1-x} y dy \right] dx \\ &= 4 \int_0^1 x \left( \frac{1-x^2}{2} \right) dx = 1/6 \end{aligned}$$



## Conditional Distributions

for discrete: The conditional probability mass function  $\rightarrow$  PDF  
 $f_{X|Y}(x|y) = P(X=x|Y=y) = \frac{P(X=x, Y=y)}{P(Y=y)} = \frac{f_{XY}(x, y)}{f_Y(y)}$   
 $P(Y=y)$  marginal PDF of  $y$

for continuous  $f_Y(y) > 0$ :  $P(X \in A | Y=y) = \int_A f_{X|Y}(x|y) dx$   
 $P(X < 1/4 | Y=1/3)$

Ex 1.  $f(x, y) = \begin{cases} x+y, & \text{if } 0 \leq x \leq 1, 0 \leq y \leq 1 \\ 0, & \text{otherwise} \end{cases}$

$$f_Y(y) = \int_0^1 x+y dx = \left[ \frac{x^2}{2} + yx \right]_0^1 = \frac{1}{2} + y$$

$$f_{X|Y}(x|y) = \frac{f_{XY}(x, y)}{f_Y(y)} = \frac{x+y}{y+1/2}$$

$$\Rightarrow \int_0^{1/4} f_{X|Y}(x|1/3) dx = \int_0^{1/4} \frac{x+1/3}{1/3+1/2} dx = \frac{\frac{1}{3}x + \frac{1}{2}x}{\frac{1}{3} + \frac{1}{2}} = \frac{11}{80}$$

Ex 2. Suppose that  $X \sim \text{Uniform}(0, 1)$ .  $Y|X=x \sim \text{Uniform}(x, 1)$  What is the marginal distribution of  $Y$

$$f_X(x) = \begin{cases} 1, & \text{if } 0 \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases} \quad \text{and} \quad f_{Y|X}(y|x) = \begin{cases} \frac{1}{1-x}, & \text{if } 0 \leq x < y < 1 \\ 0, & \text{otherwise} \end{cases}$$

$$f_{X,Y}(x, y) = f_{Y|X}(y|x) f_X(x) = \begin{cases} \frac{1}{1-x}, & \text{if } 0 \leq x < y < 1 \\ 0, & \text{otherwise} \end{cases}$$

the marginal for  $Y$  is  $f_Y(y) = \int_0^y f_{X,Y}(x, y) dx = \int_0^y \frac{dx}{1-x} = -\int_0^y \frac{du}{u} = -\log(1-y)$   
 $0 < y < 1$

## Multivariate Distributions and IID Samples

Let  $X = (X_1, \dots, X_n)$ , where  $X_1, \dots, X_n$  RVs.  $X$  random vector. Let  $f(x_1, \dots, x_n)$  denote the PDF

we say  $X_1, \dots, X_n$  independent for every  $A_1, \dots, A_n$ ,

$$P(X_1 \in A_1, \dots, X_n \in A_n) = \prod_{i=1}^n P(X_i \in A_i). \text{ It suffices to check } f(x_1, \dots, x_n) = \prod_{i=1}^n f_{X_i}(x_i)$$

$\Rightarrow$  If  $X_1, \dots, X_n$  independent and each has the same marginal distribution with PDF  $F$  we say that  $X_1, \dots, X_n$  are independent and identically distributed (iid).

$X_1, X_2, \dots, X_n \sim F$ . If  $F$  has density  $f$  we also write  $X_1, \dots, X_n \sim f$   
random sample size  $n$  from  $F$

**Multinomial** is the multivariate version of Binomial. eg. Draw a coloured ball from urns with  $r$  different colours. Let  $p = (p_1, \dots, p_r)$  where  $p_j \geq 0$  and  $\sum_{j=1}^r p_j = 1$ .  $p_j$  probability of drawing a ball of colour  $j$ . Draw  $n$  times (independent with replacement) and let  $X = (X_1, \dots, X_r)$  where  $X_j$  number of times that colour  $j$  appears.

$X \sim \text{Multinomial}(n, p)$  with prob function:  $f(x) = \frac{n!}{x_1! \dots x_r!} p_1^{x_1} \dots p_r^{x_r}$

⊗ Suppose that  $X \sim \text{Multinomial}(n, p)$  where  $X = (X_1, \dots, X_r)$  and  $p = (p_1, \dots, p_r)$ . The marginal distribution of  $X_j$  is Binomial  $(n, p_j)$



Multivariate Normal has two parameters  $\mu$  and  $\sigma$

$Z = \begin{pmatrix} z_1 \\ \vdots \\ z_n \end{pmatrix}$  where  $z_1, \dots, z_n \sim N(0,1)$  independent. The density  $f$  is

$$f(z) = \prod_{i=1}^n f(z_i) = \frac{1}{(2\pi)^{n/2}} \exp\left\{-\frac{1}{2} \sum_{i=1}^n z_i^2\right\} = \frac{1}{(2\pi)^{n/2}} \exp\left\{-\frac{1}{2} z^T \Sigma z\right\} \quad \Sigma: \text{variance-covariance matrix}$$

Transformations of RVs: Let  $X$ : RV with PDF  $f_X$  and CDF  $F_X$ . Let  $Y = r(X)$

eg  $Y = X^2$ ,  $Y = e^X$ .  $Y$  is a function of  $X$  and  $r(X)$  transformation of  $X$

Compute the PDF and CDF of  $Y$  In the discrete case is easy and more function

$$\text{is } f_Y(y) = P(Y=y) = P(r(X)=y) = P(\{x; r(x)=y\}) = P(X \in r^{-1}(y))$$

Ex. Suppose  $P(X=-1) = P(X=1) = 1/4$  and  $P(X=0) = 1/2$ . Let  $Y = X^2$

Then  $P(Y=0) = P(X=0) = 1/2$  and  $P(Y=1) = P(X=1) + P(X=-1) = 1/2$

$x$	$f_X(x)$	$y$	$f_Y(y)$	$Y$ takes fewer values than $X$ because the transformation is not one-to-one
-1	1/4	0	1/2	
0	1/2	1	1/2	
1	1/4			

Continuous (3 steps to find  $f_Y$ ) 1. For each  $y \in \mathbb{R}$ , find the set  $A_y = \{x; r(x) \leq y\}$

$$2. \text{ find CDF } F_Y(y) = P(Y \leq y) = P(r(X) \leq y) = P(\{x; r(x) \leq y\}) = \int_{A_y} f_X(x) dx$$

$$3. \text{ PDF is } f_Y(y) = F'_Y(y)$$

Ex. Let  $f_X(x) = e^{-x}$ ,  $x > 0$ . Hence  $F_X(x) = \int_0^x f_X(s) ds = 1 - e^{-x}$

Let  $Y = r(X) = \log(X)$ . Then  $A_y = \{x; x \leq e^y\}$

$$F_Y(y) = P(Y \leq y) = P(\log X \leq y) = P(X \leq e^y) = F_X(e^y) = 1 - e^{-e^y} \text{ for } y \in \mathbb{R}$$

Ex 2: Let  $X \sim \text{Uniform}(-1, 3)$  find the PDF of  $Y = X^2$ . The density of  $X$  is

$$f_X(x) = \begin{cases} 1/4 & \text{if } -1 \leq x \leq 3 \\ 0 & \text{otherwise} \end{cases} \quad Y \text{ can only take } (0, 9)$$

Two cases. i)  $0 < y < 1$  and ii)  $1 \leq y < 9$  for (i)  $A_y = [-\sqrt{y}, \sqrt{y}]$  and

$$F_Y(y) = \int_{A_y} f_X(x) dx = 1/2 \cdot \sqrt{y} \quad \text{for (ii) } A_y = [-1, \sqrt{y}] \text{ and } F_Y(y) =$$

$$\int_{A_y} f_X(x) dx = (1/4)(\sqrt{y} + 1) \quad \text{if we differentiate:}$$

$$f_Y(y) = \begin{cases} \frac{1}{4\sqrt{y}} & \text{if } 0 < y < 1 \\ \frac{1}{8\sqrt{y}} & \text{if } 1 \leq y < 9 \\ 0 & \text{otherwise} \end{cases}$$

\* when  $r$  strictly increasing or decreasing then  $r$  has an inverse

$$s = r^{-1} \quad f_Y(y) = f_X(s(y)) \left| \frac{ds(y)}{dy} \right|$$

Transformations of several RVs: if  $X, Y$  RVs we might want to know

the distribution of  $X/Y$ ,  $X+Y$ ,  $\max\{X, Y\}$ ,  $\min\{X, Y\}$ . Let  $Z = r(X, Y)$

the steps to find  $f_Z$ .



Independence: then:

1)  $P(A|B) = P(A)$

2)  $P(A \cap B) = P(A) \cdot P(B)$

1. For each  $z$ , find the set  $A_z = \{(x, y) : r(x, y) \leq z\}$

2. Find the CDF:  $F_Z(z) = P(Z \leq z) = P(r(X, Y) \leq z)$   
 $= P(\{(x, y) : r(x, y) \leq z\}) = \iint_{A_z} f_{X,Y}(x, y) dx dy$

3. Then  $f_Z(z) = F'_Z(z)$

Ex. Let  $X_1, X_2 \sim \text{Uniform}(0, 1)$  be independent. Find density  $Y = X_1 + X_2$ . the joint density of  $(X_1, X_2)$

$f(x_1, x_2) = \begin{cases} 1, & 0 \leq x_1 < 1, 0 \leq x_2 < 1 \\ 0, & \text{otherwise} \end{cases}$

Let  $r(x_1, x_2) = x_1 + x_2$

$F_Y(y) = P(Y \leq y) = P(r(X_1, X_2) \leq y)$   
 $= P(\{(x_1, x_2) : r(x_1, x_2) \leq y\}) = \iint_{A_y} f(x_1, x_2) dx_1 dx_2$

find  $A_y$ : Suppose  $0 < y \leq 1$ . Then  $A_y$  is the triangle with vertices  $(0, 0)$ ,  $(y, 0)$  and  $(0, y)$ . Then  $\iint_{A_y} f(x_1, x_2) dx_1 dx_2$  is the area of the triangle with vertices  $(0, 0)$ ,  $(y, 0)$  and  $(0, y)$ . This set has area  $y^2/2$ .

$F_Y(y) = \begin{cases} y^2/2, & 0 \leq y < 1 \\ 1 - [(1-y)^2]/2, & 1 \leq y < 2 \\ 1, & y \geq 2 \end{cases}$  if we differentiate the PDF:

$f_Y(y) = \begin{cases} y & 0 \leq y \leq 1 \\ 2-y & 1 \leq y \leq 2 \\ 0 & \text{otherwise} \end{cases}$  (xx) helps you find the probability calculate the expected value of a function of a RV without having to find the probab. distribution of a function itself.

### Chapter 3: Expectation

Expected value, mean. First moment of  $X$  is defined to be

$E(X) = \int x dF(x) = \begin{cases} \sum x_i \cdot f(x_i) & \text{if } X \text{ discrete} \\ \int x \cdot f(x) dx & \text{if } X \text{ continuous} \end{cases}$  assuming integral is well defined

$E(X)$  well defined if  $\int |x| df_X(x) < \infty$

$X \sim \text{Bernoulli}(p)$  Then  $E(X) = \sum_{x=0,1} x \cdot f(x) = 0(1-p) + (1)p = p$

$X \sim \text{Uniform}(-1, 3)$  Then  $E(X) = \int x dF_X(x) = \int x f_X(x) dx = \frac{1}{4} \int_{-1}^3 x dx = 1$

\* Cauchy distribution:  $f_X(x) = \{\pi(1+x^2)\}^{-1}$  using  $u=x$   $v=\tan^{-1}x$

$\int |x| dF(x) = \frac{2}{\pi} \int_0^\infty \frac{x dx}{1+x^2} = [x \tan^{-1}(x)]_0^\infty - \int_0^\infty \tan^{-1}x dx = \infty$

Mean does not exist. If you simulate many times the average never settles down  $\hookrightarrow$  because has thick tails  $\rightarrow$  extreme observations are common.

(\*) The Rule of the Lazy Statistician: Let  $Y = r(X)$ . Then  $\rightarrow$  function transformation

$E(Y) = E(r(X)) = \int r(x) dF_X(x) = \sum_{i=1}^n r(x_i) \cdot P(X=x_i)$



\* The  $k$ -th moment of  $X$  is defined  $E(X^k)$  assuming  $E(|X|^k) < \infty$

If the  $k$ -th moment exists and if  $j < k$  then the  $j$ -th moment exists.

Proof:  $E|X|^k = \int_{-\infty}^{\infty} |x|^k p_X(x) dx$

$k$ -th moment is defined to be  $E((X-\mu)^k)$

Properties of Expectations:

1. If  $X_1, \dots, X_n$  are RVs and  $a_1, \dots, a_n$  are constants:

$$E\left(\sum_{i=1}^n a_i X_i\right) = \sum_{i=1}^n a_i E(X_i)$$

Let  $X$  be RV with mean  $\mu$ . The variance of  $X$  ( $\sigma^2$  or  $\sigma_X^2$  or  $V(X)$ )

$$\sigma^2 = E(X-\mu)^2 = \int (x-\mu)^2 dF(x)$$

Properties: 1)  $V(X) = E(X^2) - \mu^2$

2. If  $a, b$  constants:  $V(aX+b) = a^2 V(X)$

3. If  $X_1, \dots, X_n$  independent and  $a_1, \dots, a_n$  constants then:

$$V\left(\sum_{i=1}^n a_i X_i\right) = \sum_{i=1}^n a_i^2 V(X_i)$$

Let  $X$  and  $Y$  be RVs with means  $\mu_X$  and  $\mu_Y$  and s.d.  $\sigma_X, \sigma_Y$

Covariance:  $\text{Cov}(X, Y) = E((X-\mu_X)(Y-\mu_Y))$

Correlation:  $\rho = \rho_{XY} = \rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$

Covariance satisfies:  $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$

Correlation:  $-1 \leq \rho(X, Y) \leq 1$

\* If  $X, Y$  independent  $\text{Cov}(X, Y) = \rho = 0$

$$V(X+Y) = V(X) + V(Y) + 2\text{Cov}(X, Y)$$

$$V(X-Y) = V(X) + V(Y) - 2\text{Cov}(X, Y)$$

Variance - Covariance matrix.

Conditional Expectation:

$X, Y$  RVs, the conditional expectation of  $X$  given  $Y=y$  is:

$$E(X|Y=y) = \begin{cases} \sum x f_{X|Y}(x|y) & \text{discrete} \\ \int x f_{X|Y}(x|y) dx & \text{continuous} \end{cases}$$

If  $r(x, y)$  is a function of  $x$  and  $y$  then:

$$E(r(X, Y)|Y=y) = \begin{cases} \sum r(x, y) f_{X|Y}(x|y) & \text{discrete} \\ \int r(x, y) f_{X|Y}(x|y) dx & \text{continuous} \end{cases}$$

Ex:  $X \sim \text{Unif}(0, 1)$ . After we observe  $X=x$  we draw  $Y|X=x \sim \text{Unif}(0, 1)$



intuitively we expect:  $E(Y|X=2) = (1+2)/2$  In fact:

$f_{Y|X}(y|x) = 1/(1-x)$  for  $x < y < 1$  and

$$E(Y|X=x) = \int_x^1 y \cdot f_{Y|X}(y|x) \cdot dy = \frac{1}{1-x} \int_x^1 y \cdot dy = \frac{1+x}{2} \text{ thus}$$

$E(Y|X) = (1+X)/2$  is a RV whose value is the number  $E(Y|X=x) = (1+x)/2$  once  $X=x$  is observed

→ The Rule of Iterated Expectations: for RV  $X$  and  $Y$  we have:

$$E[E(Y|X)] = E(Y) \quad \text{and} \quad E[E(X|Y)] = E(X)$$

for any function  $r(x,y)$   $E[E(r(X,Y)|X)] = E(r(X,Y))$

compute  $E(Y)$ : we know  $E(Y|X) = (1+X)/2$

$$\Rightarrow E(Y) = E[E(Y|X)] = E\left(\frac{1+X}{2}\right) = \frac{1+E(X)}{2} = \frac{1+1/2}{2} = 3/4$$

Conditional Variance

$$V(Y|X=x) = \int (y - \mu(x))^2 f(y|x) \cdot dy, \quad \mu(x) = E(Y|X=x)$$

Moment Generating Functions: (MGF) or Laplace transformation of  $X$  is:

$$\psi_X(t) = E(e^{tX}) = \int e^{tx} dF(x) \quad t: \text{real number}$$

Properties 1) If  $Y = aX + b \Rightarrow \psi_Y(t) = e^{bt} \cdot \psi_X(at)$

2) If  $X_1, \dots, X_n$  independent and  $Y = \sum_{i=1}^n X_i$  then  $\psi_Y(t) = \prod_{i=1}^n \psi_{X_i}(t)$ , where  $\psi_{X_i}$  is the MGF of  $X_i$

Moment Generating Function For Some Common Distributions:

Distribution	MGF $\psi(t)$
Bernoulli ( $p$ )	$p e^t + (1-p)$
Binomial ( $n, p$ )	$(p \cdot e^t + (1-p))^n$
Poisson ( $\lambda$ )	$e^{\lambda(e^t - 1)}$
Normal ( $\mu, \sigma$ )	$\exp\left\{t\mu + \frac{\sigma^2}{2} t^2\right\}$
Gamma ( $\alpha, \beta$ )	$\left(\frac{1}{1-t\beta}\right)^\alpha \text{ for } t < 1/\beta$



$E(X)$  - first moment

$\text{Var}(X)$  - second moment

Chapter 4 - AOS: Inequalities.  $\rightarrow$  useful for bounding quantities that are otherwise difficult to compute. For Bernoulli r.v.s:

1. Theorem 4.1: Markov's inequality: Let  $X$  be a non-negative and suppose that  $E(X)$  exists for any  $t > 0$ .  $P(X > t) \leq \frac{E(X)}{t}$  You need to know the first moment (you know by the expected value)

Proof: Since  $X \geq 0$ ,  $E(X) = \int_0^\infty x f(x) dx = \int_0^t x f(x) dx + \int_t^\infty x f(x) dx \geq \int_t^\infty x f(x) dx \geq t \int_t^\infty f(x) dx = t P(X > t)$

2. Theorem 4.2: Chebyshev's inequality: Let  $\mu = E(X)$  and  $\sigma^2 = \text{Var}(X)$ . Then

$P(|X - \mu| \geq t) \leq \frac{\sigma^2}{t^2}$  and  $P(|Z| \geq k) \leq \frac{1}{k^2}$ , where  $Z = (X - \mu)/\sigma$ . You need to know the second moment

In particular,  $P(|Z| \geq 2) \leq 1/4$  and  $P(|Z| \geq 3) \leq 1/9$ . You need to know the second moment

Proof: (use Markov's inequality)  $P(|X - \mu| \geq t) = P(|X - \mu|^2 \geq t^2) \leq \frac{E(X - \mu)^2}{t^2} = \frac{\sigma^2}{t^2}$ . to proof second part:  $t = k \cdot \sigma$ .

3. Theorem 4.4: Hoeffding's Inequality: Let  $Y_1, \dots, Y_n$  be independent observations such that  $E(Y_i) = 0$  and  $a_i \leq Y_i \leq b_i$ . Let  $t > 0$ . Then for any  $t > 0$ ,  $P(\sum_{i=1}^n Y_i \geq t) \leq e^{-t^2 / \sum_{i=1}^n (b_i - a_i)^2 / 4}$

Theorem 4.5: Let  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ . Then, for any  $\epsilon > 0$ ,  $P(|\bar{X}_n - p| > \epsilon) \leq 2e^{-2n\epsilon^2}$  where  $\bar{X} = n^{-1} \sum_{i=1}^n X_i$

eg. Let  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$  let  $n = 100$ ,  $\epsilon = 0.2$  using Chebyshev's inequality yielded  $P(|\bar{X}_n - p| > \epsilon) \leq 0.0625$ . According to Hoeffding's inequality  $P(|\bar{X}_n - p| > 0.2) \leq 2e^{-2 \cdot 100 \cdot 0.2^2} = 0.00067$  (much smaller)

Fix  $\alpha > 0$  and let  $\epsilon_n = \sqrt{\frac{1}{2n} \log(\frac{2}{\alpha})}$ . By Hoeffding's inequality:  $P(|\bar{X}_n - p| > \epsilon_n) \leq 2e^{-2n\epsilon_n^2} = \alpha$

Let  $C = (\bar{X}_n - \epsilon_n, \bar{X}_n + \epsilon_n)$ . Then  $P(p \notin C) = P(|\bar{X}_n - p| > \epsilon_n) \leq \alpha$ .  $P(p \in C) \geq 1 - \alpha \rightarrow$  this interval  $C$  traps the true parameter  $p$  with prob  $1 - \alpha$  so  $C$  is a  $1 - \alpha$  confidence interval.

for normal RV:

Theorem: Mill's inequality: Let  $Z \sim N(0, 1)$  then  $P(|Z| > t) \leq \frac{2}{\sqrt{t}} e^{-t^2/2}$

Inequalities for Expectations

1. Theorem 4.7: Cauchy-Schwarz inequality: If  $X$  and  $Y$  have finite variances then  $E|XY| \leq \sqrt{E(X^2) \cdot E(Y^2)}$ .  $|x^T y| \leq \|x\| \|y\|$  the only way  $|x^T y| = \|x\| \|y\|$  when  $x$  is a multiple of  $y$

\* A function  $g$  is convex if for each  $x, y$  and each  $\alpha \in [0, 1]$   $g(\alpha x + (1 - \alpha)y) \leq \alpha g(x) + (1 - \alpha)g(y)$



If  $g$  twice differentiable and  $g''(x) \geq 0$  for all  $x \Rightarrow g$  is convex.  
convex then  $g$  lies above any line that touches  $g$  at some point (tangent)

\* A function is concave if  $-g$  is convex.

Convex functions:  $x^2, e^x$ , concave functions:  $-x^2, \log(x)$

Theorem 4.9 Jensen's Inequality: If  $g$  is convex then  $Eg(X) \geq g(E(X))$

If  $g$  is concave then  $Eg(X) \leq g(E(X))$

Proof: Let  $L(x) = a + bx$  line tangent to  $g(x)$  at the point  $E(X)$ . Since  $g$  convex  $\rightarrow$  lies above the line  $L(x)$ .

$$Eg(X) \geq EL(X) = E(a + bX) = a + bE(X) = L(E(X)) = g(E(X))$$

\* From Jensen's inequality we see  $E(X^2) \geq (EX)^2$  if  $X$  positive, then  $E(1/X) \geq 1/EX$  since  $\log$  is concave  $E(\log X) \leq \log EX$

### Chapter 5: (AOS) Convergence of RVs

Behaviour of sequences of RVs: Large sample theory, limit theory and asymptotic theory. What happens if we gather more data?

A sequence of real numbers  $x_n$  converges to a limit  $x$  if for every  $\epsilon > 0$ ,

$|x_n - x| < \epsilon$  for all large  $n$ . Suppose  $x_n = x$  for all  $n$ . Then  $\lim_{n \rightarrow \infty} x_n = x$ .

Suppose  $X_1, X_2, \dots$  sequence of RV independent and each  $\sim N(0, 1)$ .

$X_n$  "converges" to  $X \sim N(0, 1)$  is not quite right since  $P(X_n = X) = 0$  for all  $n$  (two continuous RVs are equal with prob. zero)

1) The Law of Large Numbers: says that the sample average  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  converges in probability to the expectation  $\mu = E(X_i)$ . This means that  $\bar{X}_n$  is close to  $\mu$  with high prob.

2) The Central Limit Theorem says that  $\sqrt{n}(\bar{X}_n - \mu)$  converges in distribution to a Normal distribution. This means that the average has approximately a Normal dist. for large  $n$ .

Types of Convergence: (2 main types) Let  $X_1, X_2, \dots$  be a sequence of RV and let

$X$  be another random variable. Let  $F_n$  denote the CDF of  $X_n$  and let  $F$  denote the CDF of  $X$ .

1.  $X_n$  converges to  $X$  in probability written  $X_n \xrightarrow{p} X$  if for every  $\epsilon > 0$ :  
 $IP(|\bar{X}_n - X| > \epsilon) \rightarrow 0$  as  $n \rightarrow \infty$

2.  $X_n$  converges to  $X$  in distribution written  $X_n \xrightarrow{d} X$  if  
 $\lim_{n \rightarrow \infty} F_n(t) = F(t)$  at all  $t$  for which  $F$  is continuous.



$X_n$  converges to  $X$  in quadratic mean (convergence in  $L_2$ ), written  $X_n \xrightarrow{qm} X$  if  $E(X_n - X)^2 \rightarrow 0$  as  $n \rightarrow \infty$

Ex. Let  $X_n \sim N(0, 1/n)$ . Let  $F$  = distribution function for a point mass at 0. Note  $\sqrt{n} X_n \sim N(0, 1)$ .  $Z$  = standard normal RV. For  $t < 0$ :  $F_n(t) = P(X_n < t) = P(\sqrt{n} X_n < \sqrt{n} t) = P(Z < \sqrt{n} t) \rightarrow 0$  since  $\sqrt{n} t \rightarrow -\infty$ . For  $t > 0$ :  $F_n(t) = P(X_n < t) = P(\sqrt{n} X_n < \sqrt{n} t) = P(Z < \sqrt{n} t) \rightarrow 1$  since  $\sqrt{n} t \rightarrow \infty$ . Hence  $F_n(t) \rightarrow F(t)$  for all  $t \neq 0$  and so  $X_n \xrightarrow{qm} 0$ .  $F_n(0) = 1/2 \neq F(0) = 1 \Rightarrow$  convergence fails at  $t=0$ . That it doesn't matter because it's a continuity point of  $F$  and definition of convergence in distribution requires convergence at continuity points. Now consider convergence in probability.

For any  $\epsilon > 0$  using Markov's inequality,  $P(|X_n| > \epsilon) = P(|X_n|^2 > \epsilon^2) \leq \frac{E[X_n^2]}{\epsilon^2} = \frac{1/n}{\epsilon^2} \rightarrow 0$  as  $n \rightarrow \infty$ .  $X_n \xrightarrow{P} 0$

Theorem 5.4: The following relationships hold:

(a)  $X_n \xrightarrow{qm} X$  implies that  $X_n \xrightarrow{P} X$

(b)  $X_n \xrightarrow{P} X$  implies that  $X_n \xrightarrow{w} X$

(c) If  $X_n \xrightarrow{w} X$  and if  $P(X=c) = 1$  for some real number  $c$ , then  $X_n \xrightarrow{P} X$ . \* Inverse implications don't hold except the special case in c).

Theorem 5.5: Let  $X_n, X, Y_n, Y$  RVs. Let  $g$  continuous function.

(a) If  $X_n \xrightarrow{P} X$  and  $Y_n \xrightarrow{P} Y$  then  $X_n + Y_n \xrightarrow{P} X + Y$

(b) If  $X_n \xrightarrow{qm} X$  and  $Y_n \xrightarrow{qm} Y$  then  $X_n + Y_n \xrightarrow{qm} X + Y$

(c) If  $X_n \xrightarrow{w} X$  and  $Y_n \xrightarrow{w} c$  then  $X_n + Y_n \xrightarrow{w} X + c$

(d) If  $X_n \xrightarrow{P} X$  and  $Y_n \xrightarrow{P} Y$  then  $X_n Y_n \xrightarrow{P} XY$

(e) If  $X_n \xrightarrow{w} X$  and  $Y_n \xrightarrow{w} c$  then  $X_n Y_n \xrightarrow{w} cX$

(f) If  $X_n \xrightarrow{P} X$  then  $g(X_n) \xrightarrow{P} g(X)$

(g) If  $X_n \xrightarrow{w} X$  then  $g(X_n) \xrightarrow{w} g(X)$

Parts c & g are known as Slutsky's theorem. Note  $X_n \xrightarrow{w} X$  and  $Y_n \xrightarrow{w} Y$  doesn't in general imply  $X_n + Y_n \xrightarrow{w} X + Y$

The Law of Large Numbers: the mean of a large sample is close to the mean of the distribution. Let  $X_1, X_2, \dots$  be an IID sample. Let  $\mu = E(X_1)$  and  $\sigma^2 = V(X_1)$ .  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ .  $E(\bar{X}_n) = \mu$ ,  $V(\bar{X}_n) = \frac{\sigma^2}{n}$

The Weak Law of Large Numbers (WLLN): If  $X_1, \dots, X_n$  IID then  $\bar{X}_n \xrightarrow{P} \mu$ . The distribution of  $\bar{X}_n$  becomes more concentrated around  $\mu$  as  $n$  gets large.



Ex. Flip a coin. Let  $X_i$  = outcome of a single toss (0 or 1). Hence  $p = E(X_i)$ . After  $n$  times  $\bar{X}_n$ . According to the Law of Large Numbers  $\bar{X}_n \rightarrow p$  in probability. This doesn't mean that  $\bar{X}_n$  will be numerically equal to  $p$  means when  $n$  is large the distribution of  $\bar{X}_n$  is tightly concentrated around  $p$ .

Let  $p = 1/2$ . How large should be  $n$  so that  $P(0.4 \leq \bar{X}_n \leq 0.6) \geq 0.7$ .

First,  $E(\bar{X}_n) = p = 1/2$  and  $V(\bar{X}_n) = \sigma^2/n = p(1-p)/n = 1/4n$ .

From Chebyshev's inequality:  $P(0.4 \leq \bar{X}_n \leq 0.6) = P(|\bar{X}_n - p| \leq 0.1)$

$$= 1 - P(|\bar{X}_n - p| > 0.1) \geq 1 - \frac{1}{4n(0.1)^2} = 1 - \frac{25}{n}$$

$$1 - \frac{25}{n} \geq 0.7 \Rightarrow n = 84$$

The Central Limit Theorem (CLT): Let  $X_1, \dots, X_n$  be IID with mean  $\mu$  and variance  $\sigma^2$ . Let  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  then  $Z_n = \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \rightarrow Z$

where  $Z \sim N(0, 1)$ . In other words,  $\lim_{n \rightarrow \infty} P(Z_n \leq z) = \Phi(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$ .

Ex. Suppose the number of errors per computer program has a Poisson distribution with mean 5. ~~at  $Z_n$  converging to~~ we get 125 programs. Let  $X_1, \dots, X_{125}$  be number of errors in the programs. We want to approximate  $P(\bar{X}_n < 5.5)$ . Let  $\mu = E(X_i) = \lambda = 5$  and  $\sigma^2 V(X_i) = \lambda = 5$ .

$$P(\bar{X}_n < 5.5) = P\left(\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} < \frac{\sqrt{n}(5.5 - \mu)}{\sigma}\right) \approx P(Z < 2.5) = 0.9738$$

\* we rarely know  $\sigma \rightarrow$  we can estimate  $\sigma^2$  from  $X_1, \dots, X_n$  by  $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  if replace  $\sigma$  with  $S_n \rightarrow$  CLT still holds.

Theorem 5.10: Assume the same conditions as the CLT then:  $\frac{\sqrt{n}(\bar{X}_n - \mu)}{S_n} \rightarrow N(0, 1)$

Theorem 5.11: The Berry-Esseen Inequality: Suppose that  $E|X_i|^3 < \infty$ . Then  $\sup_z |P(Z_n \leq z) - \Phi(z)| \leq \frac{33}{4} \frac{E|X_i - \mu|^3}{\sqrt{n} \cdot \sigma^3}$

Theorem 5.12: Multivariate central limit theorem: Let  $X_1, \dots, X_n$  be IID random vectors where

$$X_i = \begin{pmatrix} X_{i1} \\ X_{i2} \\ \vdots \\ X_{ik} \end{pmatrix} \text{ with mean } \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_k \end{pmatrix} = \begin{pmatrix} E(X_{i1}) \\ E(X_{i2}) \\ \vdots \\ E(X_{ik}) \end{pmatrix} \text{ and variance matrix } \Sigma.$$



Let  $\bar{X} = \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \end{pmatrix}$  where  $\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ji}$ . Then  $\sqrt{n}(\bar{X} - \mu) \rightsquigarrow N(0, \Sigma)$   
 The Delta Method: If  $Y_n$  has a limiting Normal distribution then the delta method allows us to find the limiting distribution of  $g(Y_n)$  where  $g$  any smooth function. Suppose that  $\sqrt{n}(Y_n - \mu) \rightsquigarrow N(0, \Sigma)$  and that  $g$  is a differentiable function such that  $g'(\mu) \neq 0$ . Then  

$$\frac{\sqrt{n}(g(Y_n) - g(\mu))}{|g'(\mu)|\sigma} \rightsquigarrow N(0, 1)$$
 In other words.

$Y_n \approx N(\mu, \frac{\sigma^2}{n})$  implies that  $g(Y_n) \approx N(g(\mu), (g'(\mu))^2 \cdot \frac{\sigma^2}{n})$   
 Ex. Let  $X_1, \dots, X_n$  be IID with finite mean  $\mu$  and finite variance  $\sigma^2$ . By the central limit theorem  $\sqrt{n}(\bar{X}_n - \mu)/\sigma \rightsquigarrow N(0, 1)$ . Let  $W_n = e^{\bar{X}_n}$ .  
 Thus  $W_n = g(\bar{X}_n)$  where  $g(s) = e^s$ . Since  $g'(s) = e^s$  the delta method implies that  $W_n \approx N(e^\mu, e^{2\mu} \cdot \sigma^2/n)$ .

Theorem 5.15: The Multivariate Delta Method. Suppose that  $Y_n = (Y_{n1}, Y_{n2}, \dots, Y_{nk})$  is a sequence of Random vectors such that  $\sqrt{n}(Y_n - \mu) \rightsquigarrow N(0, \Sigma)$ .  
 Let  $g: \mathbb{R}^k \rightarrow \mathbb{R}$  and let  $\nabla g(y) = \begin{pmatrix} \frac{\partial g}{\partial y_1} \\ \vdots \\ \frac{\partial g}{\partial y_k} \end{pmatrix}$ . Let  $\nabla_\mu$  denote  $\nabla g(y)$  evaluated at  $y = \mu$  and assume that the elements of  $\nabla_\mu$  are non-zero.

Then  $\sqrt{n}(g(Y_n) - g(\mu)) \rightsquigarrow N(0, \nabla_\mu^T \Sigma \nabla_\mu)$

Ex. Let  $\begin{pmatrix} x_{1i} \\ x_{2i} \end{pmatrix}, \begin{pmatrix} x_{1i} \\ x_{2i} \end{pmatrix}, \dots, \begin{pmatrix} x_{1i} \\ x_{2i} \end{pmatrix}$  be IID random vectors with mean  $\mu = (\mu_1, \mu_2)^T$  and variance  $\Sigma$ . Let  $\bar{X}_1 = \frac{1}{n} \sum_{i=1}^n x_{1i}$ ,  $\bar{X}_2 = \frac{1}{n} \sum_{i=1}^n x_{2i}$  and define  $Y_n = (\bar{X}_1, \bar{X}_2)$ . Thus  $Y_n = g(\bar{X}_1, \bar{X}_2)$  where  $g(s_1, s_2) = s_1 s_2$ . By the CLT:  
 $\sqrt{n} \begin{pmatrix} \bar{X}_1 - \mu_1 \\ \bar{X}_2 - \mu_2 \end{pmatrix} \rightsquigarrow N(0, \Sigma)$  Now:  $\nabla g(s) = \begin{pmatrix} \frac{\partial g}{\partial s_1} \\ \frac{\partial g}{\partial s_2} \end{pmatrix} = \begin{pmatrix} s_2 \\ s_1 \end{pmatrix}$  and so

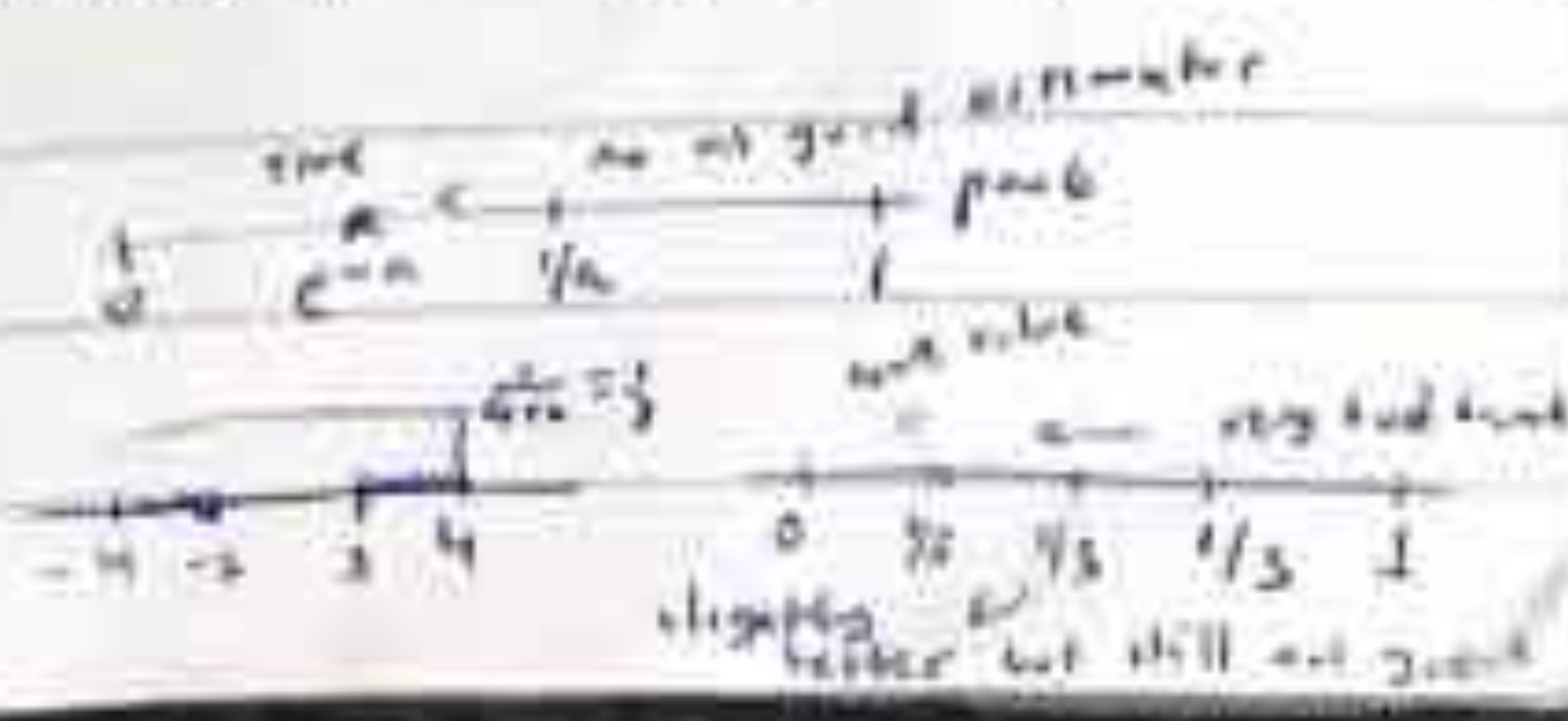
$$\nabla_\mu^T \Sigma \nabla_\mu = (\mu_2 \ \mu_1) \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{pmatrix} \begin{pmatrix} \mu_2 \\ \mu_1 \end{pmatrix} = \mu_2^2 \sigma_{11} + 2\mu_1 \mu_2 \sigma_{12} + \mu_1^2 \sigma_{22}$$

Therefore,  $\sqrt{n}(\bar{X}_1 \bar{X}_2 - \mu_1 \mu_2) \rightsquigarrow N(0, \mu_2^2 \sigma_{11} + 2\mu_1 \mu_2 \sigma_{12} + \mu_1^2 \sigma_{22})$

Markov's Inequality: Let's say a RV, is non-negative, <sup>the</sup> probability that the RV exceeds <sup>a</sup> that particular number <sup>(a)</sup> is bounded by the ratio  $\frac{E[X]}{a}$ . If  $f(x)$  very small  $\Rightarrow$  the probability of exceeding that value of  $a$  will also be small, if  $a$  very large  $\Rightarrow$  prob of exceeding that large value drops.

ex 1)  $X$  exponential ( $\lambda=1$ )  $P(X \geq a) \leq \frac{E[X]}{a} = \frac{1}{a}$

ex 2)  $X$  uniform  $[-4, 4]$   $P(X \geq 3) \leq \frac{E[X]}{3}$   
 (non-negative so we cannot apply it directly)  
 $P(X \geq 3) \leq P(|X| \geq 3) \leq \frac{E[|X|]}{3} = \frac{2}{3}$  (because of symmetry)





2) The Chebyshev's Inequality:  $RV \sim$  with finite mean  $\mu$  and variance  $\sigma^2$

"If the variance is small, then  $X$  is unlikely to be too far from the mean."

$$P(|X - \mu| \geq c) \leq \frac{\sigma^2}{c^2} \quad \text{if the variance is small the prob of falling far from the mean will be small}$$

if  $k=3$ : the probability that you fall 3 standard deviations from the mean or more will be less than or equal to  $1/9$  (this is true for most distributions)

Ex. 2:  $X \sim \exp(\lambda=1)$   $P(X \geq a) \leq \frac{1}{a} \sim$  Markov's, real value  $\sim e^{-a}$

$$P(X \geq a) = P(X-1 \geq a-1) \leq P(|X-1| \geq a-1)$$

using Chebyshev's:  $\leq \frac{1}{(a-1)^2} \rightarrow$  if  $a$ : large number this quantity will behave as  $\frac{1}{a^2}$   
for large  $a$ : the Chebyshev bound

3) Hoeffding's Inequality:  $P(X_1 + \dots + X_n \geq na)$   $X_i \sim iid$  the  $RV$  are equally likely to take the values  $-1$  and  $1$  with  $prob = \frac{1}{2}$

$$Y = X_1 + \dots + X_n \quad E[Y] = 0, \text{Var}(Y) = n$$

$E[X_i] = 0$  be symmetric  $-1$   $1$  What we know about the  $RV$   $Y$ ? Using the CLT  $\rightarrow Y$

$\text{Var}(X_i) = 1$  has approximately Normal distr ( $\mu=0$ ):  $P(\frac{Y}{\sqrt{n}} \geq a) \approx 1 - Q(a)$

using Chebyshev's:  $P(Y \geq na) \leq \frac{\text{Var}(Y)}{n^2 a^2} = \frac{1}{na^2}$   $\rightarrow$  this prob goes to zero as fast as  $\frac{1}{n}$  goes to 0 (very conservative)

Hoeffding's Inequality says that this prob falls exponentially with  $n$ :  $P(X_1 + \dots + X_n \geq na)$  let  $s > 0$ ;  $a > 0$ :  $P(e^{s(X_1 + \dots + X_n)} \geq e^{sna})$

using Markov's inequality:  $\leq E[e^{s(X_1 + \dots + X_n)}] / e^{sna} = E[e^{sX_1} \dots e^{sX_n}] / e^{sna}$  since  $X_i$  independent

$= E[e^{sX_1}] \dots E[e^{sX_n}] = [E[e^{sX_1}]]^n / e^{sna} = \left( \frac{E[e^{sX_1}]}{e^{sa}} \right)^n = p^n$  if  $p < 1$ : this bound falls exponentially with  $n$  and this prob will fall the same if you choose  $s'$

to be as informative for the bound:  $\left[ \frac{\frac{1}{2}(e^s + e^{-s})}{e^{sa}} \right]^n$

you need an exponentially decaying bound

The denominator has a positive derivative at 0

For "small"  $s$ :  $p < 1$  If  $s=a$  using the Hoeffding's:  $\leq e^{-na/2}$

we Taylor's reverse  $e^x = 1 + x + \frac{x^2}{2!} + \dots = \sum_{i=0}^{\infty} \frac{x^i}{i!}$

$$\frac{1}{2}(e^s + e^{-s}) = \frac{1}{2} \left( 1 + s + \frac{s^2}{2!} + \frac{s^3}{3!} + \dots \right) + \frac{1}{2} \left( 1 - s + \frac{s^2}{2!} - \frac{s^3}{3!} + \dots \right) = \sum_{i=0}^{\infty} \frac{s^{2i}}{(2i)!} \leq \sum_{i=0}^{\infty} \frac{s^{2i}}{i! 2^i}$$

$$(2i)! = 1 \cdot 2 \cdot 3 \cdot \dots \cdot i \cdot (i+1) \cdot \dots \cdot (2i) \geq i! 2^i$$

$$= \sum_{i=0}^{\infty} \frac{(s^2/2)^i}{i!} = e^{s^2/2} \quad \text{if } s=a \quad \leq \left( \frac{e^{s^2/2}}{e^{sa}} \right)^n = e^{-na/2}$$



Jensen's Inequality: Comparing  $E[g(X)]$  to  $g(E[X])$  if  $g$ : linear these two are equal  
 $g$ : convex if  $0 \leq p \leq 1$  then  $g(px + (1-p)y) \leq pg(x) + (1-p)g(y)$   
 or  $g''(x) \geq 0$

or  $g(x) \geq g(c) + g'(c)(x-c)$   $c$ : any number  
 number = expected value of number = number

$$g(x) \geq g(E[X]) + g'(E[X])(x - E[X])$$

$$Eg(x) \geq g(E[X]) + 0$$

$$\text{eg } g(x) = x^4$$

$$(E[X])^4 \leq E[X^4]$$

$$g(x) = -\log(x) : -\log(E[X]) \leq E[-\log(x)]$$

$$\log(E[X]) \geq E[\log(x)]$$

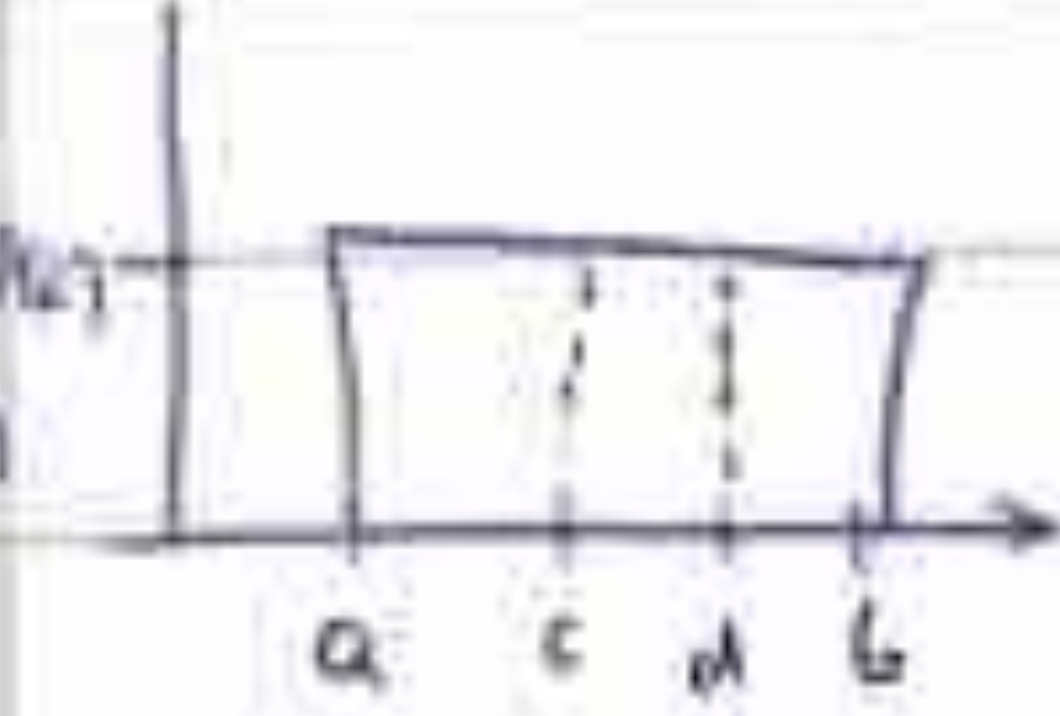


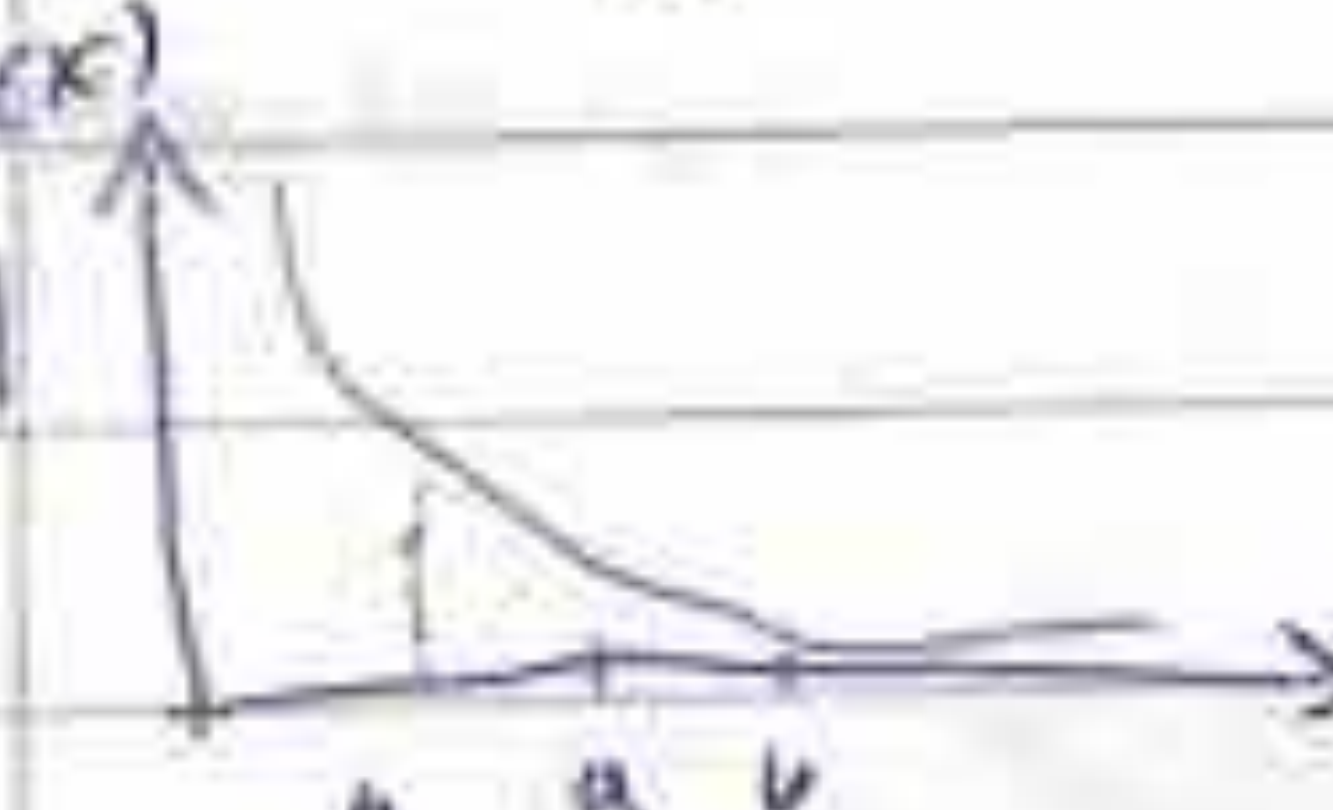
# Probability Distributions

Binomial Distribution  
 $P(X) = \binom{n}{x} \cdot p^x \cdot q^{n-x}$  getting  $x$  successes in  $n$  tries  $p$  - success,  $q$  - fail =  $1-p$   
 $P(X) = \frac{n!}{(n-x)!x!} \cdot p^x \cdot q^{n-x}$   $\mu = n \cdot p$   $\sigma = \sqrt{npq}$

2. Geometric Distribution  
 The probability that the first success will happen on the  $n^{\text{th}}$  time - the  $n^{\text{th}}$  event will succeed  
 $P(X=n) = q^{n-1} \cdot p$   $P(X > n) = q^n$   $P(X \geq n) = q^{n-1}$   
 $P(X \leq n) = 1 - q^n$   $P(X < n) = 1 - q^{n-1}$   
 $\sigma^2 = \frac{1}{p} \left( \frac{1}{p} - 1 \right) = \frac{1-p}{p^2}$   
 $m = \frac{1}{p}$

3. Poisson Distribution  
 Number of times an event occurs in a given interval.  
 $P(X=n) = \frac{m^n}{n!} e^{-m}$  or  $\frac{\lambda^n \cdot e^{-\lambda}}{n!}$   $m = \lambda = n \cdot p$   $\sigma^2 = np = \lambda$   
 $P(X > n) = 1 - e^{-m} \left[ \sum_{k=0}^n \frac{m^k}{k!} \right]$   
 $P(X \leq n) = e^{-m} \left[ \sum_{k=0}^n \frac{m^k}{k!} \right]$

4. Uniform Distribution  
  
 Area = 1 = Base  $\times$  Height =  $(b-a) \cdot f(x) \Rightarrow f(x) = \frac{1}{b-a}$   
 $P(a \leq x \leq c) = B \cdot H = \frac{c-a}{b-a}$   
 $m = \frac{a+b}{2}$   $\sigma = \frac{b-a}{\sqrt{12}}$

5. Exponential Distribution  
  
 $\lambda$  - rate parameter  
 $m = \frac{1}{\lambda}$   $\sigma^2 = \frac{1}{\lambda^2}$   
 $A_L = P(X \leq n) = 1 - e^{-\lambda n}$   
 $f(x) = \lambda \cdot e^{-\lambda x}$   
 $A_R = P(X \geq n) = \int_n^{\infty} \lambda e^{-\lambda x} dx = e^{-\lambda n}$   
 $\int_0^1 \lambda e^{-\lambda x} dx = 1 - e^{-\lambda}$

6. Bernoulli Distribution  
 One trial with two possible discrete outcomes  
 $p(x) = \begin{cases} p, & \text{for } x=1 \\ 1-p, & \text{for } x=0 \end{cases}$   $P(x) = p^x (1-p)^{1-x}$   $m=p$   $\sigma^2 = p(1-p)$

Mutually exclusive events: they cannot happen the same time eg. get head & tails by flipping one coin  $P(A \cap B) = 0$

Independent: the occurrence of one event doesn't affect the occurrence of another

$P(A \cap B) = P(A) \cdot P(B)$   $P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A) \cdot P(B)}{P(B)} = P(A)$

eg. Two independent fair coin tosses.  $H_1$  - 1st toss is H,  $H_2$  - 2nd toss is H

$P(H_1) = P(H_2) = 1/2$  Let  $C$  - 2 tosses with same result.  $\{HH\}, \{TT\}$   $P(C) = 1/2$

Check for Pairwise independence  $P(H_1, C) = P(H_1 \cap H_2) = 1/4$   $P(H_1) \cdot P(C) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$  } pairwise independent

⊗ if mutually independent  $\Rightarrow$  always pairwise independent (not the other way around)



Student Planning Exam

10 yes/no questions

$N \sim (10, 1/10)$  <sup>count answers</sup>

guess the answer with equal probability  $Z \sim (10 - N, 1/2)$  # of correct answers

$Y = N + Z$   $Y$ : number of total correct answers

deterministic threshold  $T$  for passing such that  $Y \geq T$ ,  $T \in \{0, 1, \dots, 10\}$

for each  $T$  compute the probability that the student knows less than 5 correct answers given that he passed ( $N \leq 5$ )

Check for independence  $P(H_1 \cap H_2 \cap C) = P(H_1 H_2) = 1/4$

$$P(H_1) \cdot P(H_2) \cdot P(C) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{8}$$

they are not independent

Ex. 2 Let 4 balls and events.  $A = \{1, 2\}$ ,  $B = \{1, 3\}$  and  $C = \{1, 4\}$

to be pairwise independent should  $P(A \cap B) = P(A)P(B)$  because

$$P(\text{any ball is } 1) = \frac{1}{4}$$

$$A = \{1, 2\} \text{ and } B = \{1, 3\} \quad P(\{1\}) = \frac{1}{4} \quad P(A) = \frac{2}{4} = \frac{1}{2} \quad P(B) = \frac{2}{4} = \frac{1}{2}$$

$$P(A \cap B) = P(A)P(B) \Rightarrow \text{Holds}$$

same holds for  $P(A \cap C) = P(A)P(C)$  and  $P(C \cap B) = P(C)P(B)$

Joint Independence is a stronger condition:  $P(A \cap B \cap C) = P(A)P(B)P(C)$

$$P(A \cap B \cap C) = P(\{1\}) = \frac{1}{4}$$

$$P(A)P(B)P(C) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{8} \quad \} \neq \Rightarrow \text{not jointly independent}$$

Multivariate Distributions: (Ex. 1) Find  $P(X_1 < X_2 < X_3)$  for  $X_1, X_2, X_3$  with joint pdf

$$f(x_1, x_2, x_3) = 1 \quad 0 < x_1 < 1, 0 < x_2 < 1, 0 < x_3 < 1$$

$$P(X_1 < X_2 < X_3) = \int_0^1 \int_0^{x_1} \int_0^{x_2} 1 \, dx_3 dx_2 dx_1 = \frac{1}{6}$$

Ex. 2) Find  $F(x_1, x_2, x_3)$  for  $X_1, X_2, X_3$  with joint pdf  $f(x_1, x_2, x_3) = 6e^{-x_1 - 2x_2 - 3x_3}$   
 $x_1 > 0, x_2 > 0, x_3 > 0$

$$\begin{aligned} \text{CDF: } F(x_1, x_2, x_3) &= P(X_1 \leq x_1, X_2 \leq x_2, X_3 \leq x_3) = \int_0^{x_1} \int_0^{x_2} \int_0^{x_3} 6e^{-w_1 - 2w_2 - 3w_3} \, dw_3 dw_2 dw_1 \\ &= \int_0^{x_1} 6e^{-w_1} \, dw_1 + \int_0^{x_2} 6e^{-2w_2} \, dw_2 + \int_0^{x_3} 6e^{-3w_3} \, dw_3 \\ &= (1 - e^{-x_1})(1 - e^{-2x_2})(1 - e^{-3x_3}) \end{aligned}$$

$$F(x_1, x_2, x_3) = \begin{cases} 0 & , x_1 \leq 0, x_2 \leq 0, x_3 \leq 0 \\ (1 - e^{-x_1})(1 - e^{-2x_2})(1 - e^{-3x_3}) & x_1 > 0, x_2 > 0, x_3 > 0 \end{cases}$$

Ex. Show that  $X_1, X_2, X_3$  are pairwise independent but not mutually independent

joint pdf:  $f(x_1, x_2, x_3) = \frac{1}{4}$   $(x_1, x_2, x_3) \in \{(1, 0, 0), (0, 1, 0), (0, 0, 1), (1, 1, 1)\}$

Let consider  $X_1, X_2$ : steps find  $F(x_1, x_2)$ ,  $F(x_1, x_2) = \frac{1}{4}$   $x_1 \in \{0, 1\}, x_2 \in \{0, 1\}$

marginal dist.  $f_{X_1}(x_1) = \frac{1}{2}$   $x_1 \in \{0, 1\}$  since  $F(x_1, x_2) = f_{X_1}(x_1) f_{X_2}(x_2) = 1$   $X_1, X_2$ : pairwise independent

Mutually independent? No, since  $f_{X_1}(x_1) f_{X_2}(x_2) f_{X_3}(x_3) = \frac{1}{8} \neq \frac{1}{4}$



## after 6: Models, Statistical Inference and Learning:

Parametric & Non-parametric Models. A statistical model  $f$  is a set of distributions (density or regression functions). A parametric model is a set  $f$  that can be parameterised by a finite number of parameters.

If data come from Normal distribution then the model is:

$$f = \left\{ f(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\} : \mu \in \mathbb{R}, \sigma > 0 \right\} \quad (6.1)$$

↳ this is a two-parameter model → we have written the density as  $f(x; \mu, \sigma)$  →  $x$  is a value of the RV and  $\mu, \sigma$  parameters [2-parameter model]

Parametric model takes the form:  $f = \{f(x; \theta) : \theta \in \Theta\}$  where  $\theta$  is an unknown parameter that takes values from the parameter space  $\Theta$ . If  $\theta$  is a vector <sup>but</sup> we are interested in one component of  $\theta$  we call the remaining parameters nuisance parameters.

A nonparametric model is a set  $f$  that cannot be parameterised by a finite number of parameters. eg.  $f_{all} = \{all\}$  (CDF's) is nonparametric.

Ex 1 (One-dimensional Parametric Estimation). Let  $X_1, \dots, X_n$  be independent Bernoulli ( $p$ ) observations. The problem is to estimate the parameter  $p$ .

Ex 2 (Two-dimensional Parametric Estimation). Suppose that  $X_1, \dots, X_n \sim F$  and we assume that the PDF  $f \in \mathcal{F}$  where  $\mathcal{F}$  is given in (6.1) → there are two parameters  $\mu$  and  $\sigma$ . The goal is to estimate the parameters from the data. If we are only interested in estimating  $\mu$  then  $\mu$  is the parameter of interest and  $\sigma$  is a nuisance parameter.

Ex 3 Nonparametric estimation of functionals. Let  $X_1, \dots, X_n \sim F$ . Suppose we want to estimate  $\mu = E(X_1) = \int x dF(x)$  assuming only that  $\mu$  exists. The mean  $\mu$  may be thought of as a functional of  $F$  called statistical functional:  $\mu = T(F) = \int x dF(x)$ .

$$\text{variance } T(F) = \int x^2 dF(x) - (\int x dF(x))^2$$

Ex 4 Suppose we observe pairs of data  $(X_1, Y_1), \dots, (X_n, Y_n)$ . Perhaps  $X_i$  is the blood pressure of subject  $i$  and  $Y_i$  how long they live.  $X$  is called a predictor / regressor / feature / independent var.  $Y$  is the outcome, response, dependent var.  $R(x) = E(Y|X=x)$  regression function.

The goal of predicting  $Y$  for a new patient based on their  $X$  value is predicting. If  $Y$  is discrete ⇒ classification. We need to estimate the function  $r$  → this is called regression or curve estimation.  $Y = r(X) + \epsilon$  where  $E(\epsilon) = 0$ .

let  $\epsilon = Y - r(X)$  and hence  $Y = r(X) + \epsilon$  →  $r(X) = r(X) + \epsilon$  → moreover  $E(\epsilon) = E(\epsilon|X) = 0$   
 $E(E(Y - r(X))|X) = E(E(Y|X) - r(X)) = E(r(X) - r(X)) = 0$

Point Estimation refers to providing a single "best guess" of some quantity.

$\hat{\theta}_1, \dots, \hat{\theta}_n$  point estimates of  $\theta$  →  $\theta$  is an unknown quantity and  $\hat{\theta}$  depends on the data so  $\hat{\theta}$  is random variable.

Let  $X_1, \dots, X_n$  be  $n$  iid data points from some distribution  $F$ .  
→ A point estimator  $\hat{\theta}_n$  of a parameter  $\theta$  is some function of  $X_1, \dots, X_n$ :

$$\hat{\theta}_n = g(X_1, \dots, X_n) \quad \text{The bias of an estimator is defined by}$$

$$\text{Bias}(\hat{\theta}_n) = E(\hat{\theta}_n) - \theta. \quad \hat{\theta}_n \text{ unbiased if } E(\hat{\theta}_n) = \theta$$

A requirement of an estimator: it should converge to the true parameter as  $n \rightarrow \infty$  as we collect more & more data.

A point estimator  $\hat{\theta}_n$  of a parameter  $\theta$  is consistent if  $\hat{\theta}_n \xrightarrow{P} \theta$

→ the distribution of  $\hat{\theta}_n$  : sampling distribution

→ standard error:  $se = \text{se}(\hat{\theta}_n) = \sqrt{V(\hat{\theta}_n)}$  → estimation SE



Ex 5 Let  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$  and let  $\hat{p}_n = n^{-1} \sum_{i=1}^n X_i$ . Then  $E(\hat{p}_n) = p$   
 so  $\hat{p}_n$  is unbiased. The  $se = \sqrt{V(\hat{p}_n)} = \sqrt{p(1-p)/n} \rightarrow \hat{se} = \sqrt{\hat{p}(1-\hat{p})/n}$

The quality of a point estimate is assessed by the  $MSE = E(\hat{\theta}_n - \theta)^2$

$$MSE = \text{bias}^2(\hat{\theta}_n) + V(\hat{\theta}_n)$$

Theorem: If  $\text{bias} \rightarrow 0$  and  $se \rightarrow 0$  as  $n \rightarrow \infty$  then  $\hat{\theta}_n$  is consistent that  
 $\hat{\theta}_n \xrightarrow{P} \theta$

An estimator is asymptotically Normal if  $\frac{\hat{\theta}_n - \theta}{se} \xrightarrow{d} N(0,1)$

Confidence sets: A  $1-\alpha$  confidence interval for a parameter  $\theta$  is an interval  $C_n = (a, b)$  where  $a = a(X_1, \dots, X_n)$  are functions of the data s.t.

$P_\theta(\theta \in C_n) \geq 1-\alpha$  for all  $\theta \in \Theta$ . In other words,  $(a, b)$  traps  $\theta$  with probability  $1-\alpha$  coverage of confidence interval.

\* If  $\theta$  is vector  $\Rightarrow$  use of confidence set (such as a sphere/ellipse) instead of an interval.

Hypothesis Testing: Null hypothesis: check if the data provide sufficient evidence to reject the theory

Ex 6 Let  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$  be  $n$  independent coin flips. Test if a coin is fair. Let  $H_0$  hypothesis that the coin is fair and let  $H_1$  hypothesis that the coin is not fair.  $H_1$  alternative hypothesis.  $H_0: p = 1/2$  vs  $H_1: p \neq 1/2$ . It is reasonable to reject  $H_0$  if  $T = |\hat{p}_n - 1/2|$  is large.

$$MSE = \text{bias}^2(\hat{\theta}_n) + V(\hat{\theta}_n)$$

$$\begin{aligned} \text{Proof } E(\hat{\theta}_n - \theta)^2 &= E(\hat{\theta}_n - \bar{\theta}_n + \bar{\theta}_n - \theta)^2 \\ &= E(\hat{\theta}_n - \bar{\theta}_n)^2 + 2(\bar{\theta}_n - \theta)E(\hat{\theta}_n - \bar{\theta}_n) + E(\bar{\theta}_n - \theta)^2 \\ &= (\bar{\theta}_n - \theta)^2 + E(\hat{\theta}_n - \bar{\theta}_n)^2 \\ &= \text{bias}^2(\hat{\theta}_n) + V(\hat{\theta}_n) \end{aligned}$$

$$\text{where } E(\hat{\theta}_n - \bar{\theta}_n) = \bar{\theta}_n - \bar{\theta}_n = 0$$

$1-\alpha$ : coverage of the confidence interval

CI: if you repeat the experiment over & over the interval will contain the parameter  $95\%$  of the time.



# For 9 Parametric Inference: parametric models $\mathcal{F} = \{f(x; \theta) : \theta \in \Theta\}$

interested in a function  $T(\theta)$  eg  $X \sim N(\mu, \sigma^2)$  then the parameter is  $\theta = (\mu, \sigma)$   
 goal is to estimate  $\mu = T(\theta)$  parameter of interest,  $\sigma$  nuisance parameter  
 we have a complicated function of  $\theta$

Let  $X_1, \dots, X_n \sim \text{Normal}(\mu, \sigma^2)$  the parameter is  $\theta = (\mu, \sigma)$  and the parameter space is  
 $\Theta = \{(\mu, \sigma) : \mu \in \mathbb{R}, \sigma > 0\}$  Suppose that  $X_i$  is the outcome of a blood test and suppose we  
 are interested in  $\tau$  (the fraction of the population whose test score is larger than 1. Let  
 $Z$  denote a standard Normal RV. Then:

$$\tau = P(X > 1) = 1 - P(X < 1) = 1 - P\left(\frac{X - \mu}{\sigma} < \frac{1 - \mu}{\sigma}\right) = 1 - P(Z < \frac{1 - \mu}{\sigma}) = 1 - \Phi\left(\frac{1 - \mu}{\sigma}\right)$$

The parameter of interest is  $\tau = T(\mu, \sigma) = 1 - \Phi((1 - \mu)/\sigma)$

Method of Moments for generating parametric estimators (usually are not optimal but they are easy to compute)

Suppose, parameter  $\theta = (\theta_1, \dots, \theta_K)$  has  $K$  components. For  $j \in \{1, \dots, K\}$ , the  $j$ th moment

$$a_j = a_j(\theta) = E_\theta(X^j) = \int X^j dF_\theta(x) \quad \text{and the } j\text{-th sample moment: } \hat{a}_j = \frac{1}{n} \sum_{i=1}^n X_i^j$$

The method of moments estimator  $\hat{\theta}_n$  is defined to be the value of  $\theta$  such that:

$$a_1(\hat{\theta}_n) = \hat{a}_1, \quad a_2(\hat{\theta}_n) = \hat{a}_2, \quad \dots, \quad a_K(\hat{\theta}_n) = \hat{a}_K$$

Ex 1. Let  $X_1, X_2, \dots, X_n \sim \text{Bernoulli}(p)$ . Then  $a_1 = E_p(X) = p$  and  $\hat{a}_1 = \frac{1}{n} \sum X_i$

Ex 2. Let  $X_1, \dots, X_n \sim \text{Normal}(\mu, \sigma^2)$ . Then  $a_1 = E_\theta(X_1) = \mu$  and  $a_2 = E_\theta(X_1^2) = \text{Var}(X_1) + (E_\theta(X_1))^2 = \sigma^2 + \mu^2$

we need to solve 2 equations:  $\hat{\mu} = \frac{1}{n} \sum X_i, \quad \hat{\sigma}^2 + \hat{\mu}^2 = \frac{1}{n} \sum X_i^2$

$\Rightarrow$  solution  $\hat{\mu} = \bar{X}_n$   
 $\hat{\sigma}^2 = \frac{1}{n} \sum (X_i - \bar{X}_n)^2$

Theorem: Let  $\hat{\theta}_n$  denote the method of moments estimator. Under appropriate conditions on the model the following statements hold:

1. The estimate  $\hat{\theta}_n$  exists with prob tending to 1
2. The estimate is consistent  $\hat{\theta}_n \xrightarrow{P} \theta$
3. The estimate is asymptotically Normal  $\sqrt{n}(\hat{\theta}_n - \theta) \rightarrow N(0, \Sigma)$

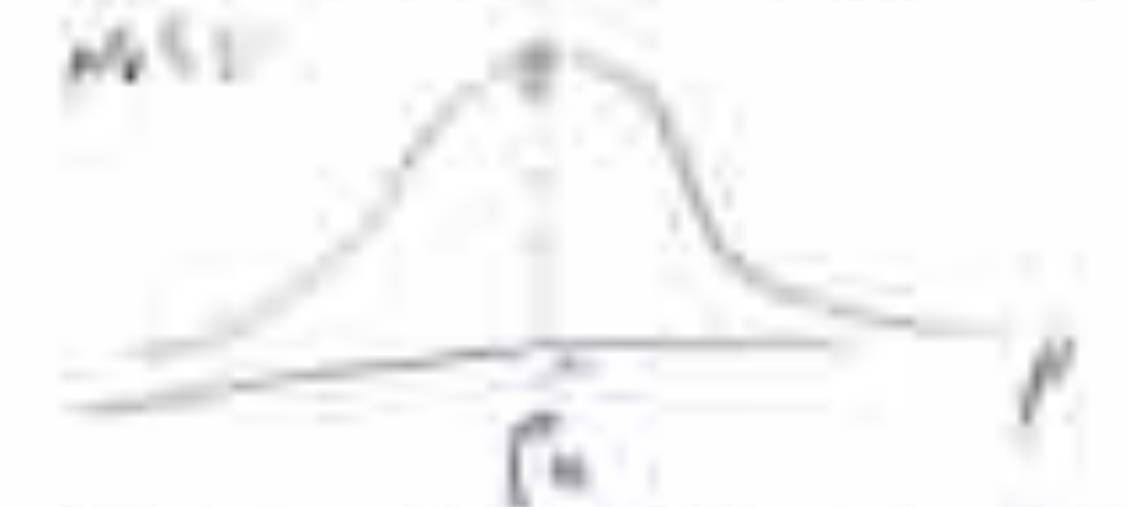
where  $\Sigma = g E_\theta(Y Y^T) g^T$ ,  $Y = (X, X^2, \dots, X^K)^T$ ,  $g = (g_1, \dots, g_K)$  and  $g_j = \partial a_j / \partial \theta$

Maximum Likelihood method for estimating parameters in a parametric model. Let  $X_1, \dots, X_n$  be iid with PDF  $f(x; \theta)$ . The Likelihood function is defined by:

$$L_n(\theta) = \prod_{i=1}^n f(X_i; \theta). \quad \text{The log-likelihood function is defined by } \ell_n(\theta) = \log L_n(\theta)$$

The likelihood function is the joint density of the data except that we treat it as a function of the parameter  $\theta$ . Thus  $L_n: \Theta \rightarrow [0, \infty]$ . This function is not a density function  $\int L_n(\theta) d\theta$  doesn't integrate to 1 (not  $\theta$ ).

The Maximum Likelihood Estimator [MLE]  $\hat{\theta}_n$  is the value that max.  $L_n(\theta)$  (or  $\ell_n(\theta)$ ). The same for  $\ell_n(\theta)$ .  
 \* If we multiply  $L_n(\theta)$  by any positive constant  $c$  (depending on  $n$ ) then this will not change the MLE. So we drop often constants in the function.



Ex: Suppose  $X_1, X_2, \dots, X_n \sim \text{Bernoulli}(p)$ . The probability function is  $f(x; p) = p^x (1-p)^{1-x}$  for  $x \in \{0, 1\}$ .  
 $L_n(p) = \prod_{i=1}^n f(X_i; p) = \prod_{i=1}^n p^{X_i} (1-p)^{1-X_i} = p^S (1-p)^{n-S}$ ,  $S = \sum X_i$

Hence  $\ell_n(p) = S \log p + (n-S) \log(1-p)$ . Take the derivative and set to zero MLE is  $\hat{p}_n = S/n$



2.2) Let  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ . The parameter is  $\theta = (\mu, \sigma)$  and  $n$  (ignoring some constants) is:

$$L_n(\mu, \sigma) = \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2}(X_i - \mu)^2\right\} = \frac{1}{\sigma^n} \exp\left\{-\frac{1}{2\sigma^2} \sum (X_i - \mu)^2\right\}$$

$$= \sigma^{-n} \exp\left\{-\frac{n\sigma^2}{2\sigma^2}\right\} \cdot \exp\left\{-\frac{n(\bar{X} - \mu)^2}{2\sigma^2}\right\} \quad \text{where } \bar{X} = n^{-1} \sum X_i$$

$$S^2 = n^{-1} \sum (X_i - \bar{X})^2$$

$$\sum (X_i - \mu)^2 = n S^2 + n(\bar{X} - \mu)^2 = \sum (X_i - \bar{X} + \bar{X} - \mu)^2$$

The log-likelihood is

$$\ell(\mu, \sigma) = -n \log \sigma - \frac{n S^2}{2\sigma^2} - n \frac{(\bar{X} - \mu)^2}{2\sigma^2}$$

$$\frac{\partial \ell(\mu, \sigma)}{\partial \mu} = 0 \quad \text{and} \quad \frac{\partial \ell(\mu, \sigma)}{\partial \sigma} = 0$$

$\hat{\mu} = \bar{X}$  and  $\hat{\sigma}^2 = S^2$  : global maximum of likelihood

Properties of MLE:

- 1) MLE is consistent:  $\hat{\theta}_n \xrightarrow{P} \theta_0$  where  $\theta_0$  = true value of parameter  $\theta$
  - 2) MLE is equivariant:  $\hat{\theta}_n$  is the MLE of  $\theta$  then  $g(\hat{\theta}_n)$  is the MLE of  $g(\theta)$
  - 3) MLE asymptotically Normal:  $(\hat{\theta}_n - \theta_0)/\hat{se} \xrightarrow{d} N(0, 1)$  also the estimated  $\hat{se}$  can often be computed analytically
  - 4) MLE asymptotically optimal/efficient: this means that among all well-behaved estimators the MLE has the smallest variance (at least for large samples)
  - 5) MLE is approximately the Bayes estimator
- \* Only hold if the model satisfies certain regularity conditions (smoothness conditions on  $f(x, \theta)$ )

- Consistency of MLE: If  $f$  and  $g$  are PDF's define the Kullback-Leibler distance between them:

$$D(f, g) = \int f(x) \log\left(\frac{f(x)}{g(x)}\right) dx$$

It can be shown that  $D(f, g) \geq 0$  and  $D(f, f) = 0$

- Equivariance of the MLE: Let  $\tau = g(\theta)$  be a function of  $\theta$ . Let  $\hat{\theta}_n$  be the MLE of  $\theta$ . The  $\hat{\tau}_n = g(\hat{\theta}_n)$  is the MLE of  $\tau$ .

ex. Let  $X_1, \dots, X_n \sim N(0, 1)$ . The MLE for  $\theta$  is  $\hat{\theta}_n = \bar{X}_n$ . Let  $\tau = e^{\theta}$ . Then the MLE for  $\tau$  is  $\hat{\tau}_n = e^{\hat{\theta}_n} = e^{\bar{X}_n}$ .

- Asymptotic Normality: The score function is defined to be  $s(X; \theta) = \frac{\partial \log f(X; \theta)}{\partial \theta}$ . The Fisher information is  $I_\theta(\theta) = V_\theta\left(\sum s(X_i; \theta)\right) = \sum V_\theta(s(X_i; \theta))$ .

ex. Let  $X_1, \dots, X_n \sim \text{Poisson}(\theta)$ . Then  $\hat{\theta}_n = \bar{X}_n$  and some calculations show that  $I_\theta(\theta) = \frac{n}{\theta}$  so  $\hat{se} = \frac{1}{\sqrt{n I(\hat{\theta}_n)}} = \sqrt{\frac{\hat{\theta}_n}{n}}$

Approximate 1- $\alpha$  CI:  $\hat{\theta}_n \pm z_{\alpha/2} \sqrt{\hat{\theta}_n/n}$

- Optimality: Suppose that  $X_1, \dots, X_n \sim N(\theta, \sigma^2)$ . The MLE is  $\hat{\theta}_n = \bar{X}_n$ . Another reasonable estimator of  $\theta$  is the sample median  $\tilde{\theta}_n$ . The MLE satisfies  $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, \sigma^2)$ .

The Delta Method: If  $\tau = g(\theta)$  where  $g$  differentiable and  $g'(\theta) \neq 0$  then  $\frac{(\hat{\tau}_n - \tau)}{\hat{se}(\hat{\tau}_n)} \xrightarrow{d} N(0, 1)$  where  $\hat{\tau}_n = g(\hat{\theta}_n)$  and  $\hat{se}(\hat{\tau}_n) = |g'(\hat{\theta}_n)| \cdot \hat{se}(\hat{\theta}_n)$ . Hence if

$C_n = (\hat{\tau}_n - z_{\alpha/2} \hat{se}(\hat{\tau}_n), \hat{\tau}_n + z_{\alpha/2} \hat{se}(\hat{\tau}_n))$  then  $P_\theta(\tau \in C_n) \rightarrow 1 - \alpha$  as  $n \rightarrow \infty$

ex. Let  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$  and let  $\psi = g(p) = \log(p/(1-p))$ . The Fisher information function is  $I(p) = 1/(p(1-p))$  so the estimated s.e. of MLE  $\hat{p}_n$  is  $\hat{se} = \sqrt{\frac{\hat{p}_n(1-\hat{p}_n)}{n}}$ . Since  $g'(p) = 1/(p(1-p))$  according to the delta method:  $\hat{se}(\hat{\psi}_n) = |g'(\hat{p}_n)| \hat{se}(\hat{p}_n) = \frac{1}{(\hat{p}_n(1-\hat{p}_n))} \sqrt{\frac{\hat{p}_n(1-\hat{p}_n)}{n}}$ . An approximate 95% CI:  $\hat{\psi}_n \pm \frac{2}{\sqrt{n \hat{p}_n(1-\hat{p}_n)}}$



## Parameter Models

$(\theta_1, \dots, \theta_k)$  and  $\theta = (\theta_1, \dots, \theta_k)$  be the MLE. Let  $l_n = \sum_{i=1}^n \log L(x_i; \theta)$ ,  
 $\frac{\partial^2 l_n}{\partial \theta_j^2}$  and  $H_{jk} = \frac{\partial^2 l_n}{\partial \theta_j \partial \theta_k}$

Define the Fisher Information Matrix by:

$$I_n(\theta) = - \begin{bmatrix} E_0(H_{11}) & E_0(H_{12}) & \dots & E_0(H_{1k}) \\ E_0(H_{21}) & E_0(H_{22}) & \dots & E_0(H_{2k}) \\ \vdots & \vdots & \ddots & \vdots \\ E_0(H_{k1}) & E_0(H_{k2}) & \dots & E_0(H_{kk}) \end{bmatrix} \quad \text{Let } J_n(\theta) = I_n^{-1}(\theta) \text{ be the inverse of } I_n.$$

Under appropriate regularity conditions  $(\hat{\theta} - \theta) \approx N(0, J_n)$

Also if  $\hat{\theta}_j$  is the  $j$ -th component of  $\hat{\theta}$  then  $\frac{(\hat{\theta}_j - \theta_j)}{s_{\hat{\theta}_j}} \rightarrow N(0, 1)$

where  $s_{\hat{\theta}_j}^2 = J_n(j, j)$  is the  $j$ th diagonal element of  $J_n$ . The approximate covariance of  $\hat{\theta}_j$  and  $\hat{\theta}_k$  is  $\text{cov}(\hat{\theta}_j, \hat{\theta}_k) \approx J_n(j, k)$

→ There is also a multiparameter delta method. Let  $\tau = g(\theta_1, \dots, \theta_k)$  be a function and let

$$\nabla g = \begin{pmatrix} \frac{\partial g}{\partial \theta_1} \\ \vdots \\ \frac{\partial g}{\partial \theta_k} \end{pmatrix} \text{ be the gradient of } g.$$

Suppose that  $\nabla g$  evaluated at  $\theta$  is not 0. Let  $\hat{\tau} = g(\hat{\theta})$ . Then  $\frac{(\hat{\tau} - \tau)}{s_{\hat{\tau}}} \rightarrow N(0, 1)$

where  $s_{\hat{\tau}}(\hat{\tau}) = \sqrt{(\hat{\nabla} g)^T J_n(\hat{\theta}) (\hat{\nabla} g)}$ ,  $J_n = J_n(\hat{\theta})$  and  $\hat{\nabla} g$  is  $\nabla g$  evaluated at  $\theta = \hat{\theta}$ .

Ex. Let  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ . Let  $\tau = g(\mu, \sigma) = \sigma/\mu$ . In Exercise 8 you will show that  $I_n(\mu, \sigma) = \begin{bmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{2n}{\sigma^2} \end{bmatrix}$ . Hence,  $J_n = I_n^{-1}(\mu, \sigma) = \frac{1}{n} \begin{bmatrix} \sigma^2 & 0 \\ 0 & \frac{\sigma^2}{2} \end{bmatrix}$ . The gradient of  $g$

$$\nabla g = \begin{pmatrix} -\frac{\sigma}{\mu^2} \\ \frac{1}{\mu} \end{pmatrix} \text{ Thus, } s_{\hat{\tau}}(\hat{\tau}) = \sqrt{(\hat{\nabla} g)^T J_n(\hat{\theta}) (\hat{\nabla} g)} = \frac{1}{\sqrt{n}} \sqrt{\frac{1}{\hat{\mu}^4} + \frac{\hat{\sigma}^2}{2\hat{\mu}^2}}$$

Parametric Bootstrap: For parametric models we can use this technique to estimate SE and CI. In the nonparametric bootstrap we sampled  $X_1^*, \dots, X_n^*$  from the empirical distribution  $\hat{F}_n$ . In the parametric bootstrap we sample instead from  $f(x, \hat{\theta}_n)$ . Hence  $\hat{\theta}_n$  could be MLE or method of moments estimator.

Ex. Consider previous example. To get the bootstrap standard error, simulate  $X_1, \dots, X_n^* \sim N(\hat{\mu}, \hat{\sigma}^2)$  compute  $\hat{\mu}^* = n^{-1} \sum X_i^*$  and  $\hat{\sigma}^{2*} = n^{-1} \sum (X_i^* - \hat{\mu}^*)^2$ . Then compute  $\hat{\tau}^* = g(\hat{\mu}^*, \hat{\sigma}^{2*}) = \hat{\sigma}^*/\hat{\mu}^*$ . Repeat this  $B$  times yields bootstrap replications  $\hat{\tau}_1^*, \dots, \hat{\tau}_B^*$  and:

$$s_{\hat{\tau}} = \sqrt{\frac{\sum (\hat{\tau}_b^* - \hat{\tau})^2}{B}} \quad \text{Bootstrap easier than delta method. Delta's method advantage is the closed form expression for the standard error.}$$

Checking Assumptions: Good idea to check that your data come from a parametric model, by inspecting the data. Eg. if a histogram of the data have very bimodal then the assumption of normality might be questionable. A formal way to test a parametric model is to use a goodness-of-fit test.



## ADS (Chapter 22: Classification)

Consider  $n$  data  $(x_1, y_1), \dots, (x_n, y_n)$  where  $x_i = (x_{i1}, \dots, x_{id}) \in \mathcal{X} \subset \mathbb{R}^d$  and  $y_i$  takes values in some finite set  $\mathcal{Y}$ . A classification rule is a function  $h: \mathcal{X} \rightarrow \mathcal{Y}$ . When we observe a new  $x$  we predict  $y$  to be  $h(x)$ .

Our goal is to find a classification rule  $h$  that makes accurate predictions.

The true error rate of a classifier  $h$  is  $L(h) = P(\{h(X) \neq Y\})$  and the empirical error rate is  $\hat{L}(h) = \frac{1}{n} \sum I(h(x_i) \neq y_i)$ .

First we consider the special case when  $\mathcal{Y} = \{0, 1\}$ . Let  $r(x) = E(Y|X=x) = P(Y=1|X=x)$  (regression function).

From Bayes' theorem we have that  $r(x) = \frac{f_1(x) P(Y=1)}{f_1(x) P(Y=1) + f_0(x) P(Y=0)}$  where  $f_0(x) = f(x|Y=0)$  and  $f_1(x) = f(x|Y=1)$ .

$\pi = P(Y=1)$

The Bayes' classification rule  $h^*$  is  $h^*(x) = \begin{cases} 1 & \text{if } r(x) > \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$ . The set  $D(h) = \{x | P(Y=1|X=x) > P(Y=0|X=x)\}$  is called the decision boundary.

The Bayes rule is optimal, that is if  $h$  is any other classification rule then  $L(h) \geq L(h^*)$ . 3 main approaches for approximation to the Bayes rule:

1. Empirical Risk Minimization: choose a set of classifiers  $\mathcal{H}$  and find  $h \in \mathcal{H}$  that min some estimate of  $L(h)$ .

2. Regression: find an estimate  $\hat{r}$  of the regression function  $r$  and define  $\hat{h}(x) = \begin{cases} 1 & \text{if } \hat{r}(x) > \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$ .

3. Density Estimation: estimate  $f_i$  from the  $x_i$ 's for which  $y_i = 0$ , estimate  $f_1$  from the  $x_i$ 's for which  $y_i = 1$  and let  $\hat{\pi} = n^{-1} \sum y_i$ . Define:

$$\hat{r}(x) = \hat{P}(Y=1|X=x) = \frac{\hat{\pi} \hat{f}_1(x)}{\hat{\pi} \hat{f}_1(x) + (1-\hat{\pi}) \hat{f}_0(x)} \quad \text{and} \quad \hat{h}(x) = \begin{cases} 1 & \text{if } \hat{r}(x) > \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

Suppose that  $\mathcal{Y} = \{1, \dots, K\}$ . The optimal rule is:

$$h(x) = \arg \max_k P(Y=k|X=x) = \arg \max_k \pi_k f_k(x)$$

$$\text{where } P(Y=k|X=x) = \frac{f_k(x) \pi_k}{\sum_{r=1}^K f_r(x) \pi_r}, \quad \pi_r = P(Y=r), \quad f_r(x) = f(x|Y=r) \quad \text{and } \arg \max_k$$

means: the value of  $k$  that max the expression.

Gaussian & Linear Classifiers: Suppose that  $\mathcal{Y} = \{0, 1\}$  and that  $f_0(x) = f(x|Y=0)$ ,  $f_1(x) = f(x|Y=1)$  are both multivariate Gaussians:  $f_k(x) = \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k)\right\}$ ,  $k=0,1$ .

Thus  $X|Y=0 \sim N(\mu_0, \Sigma_0)$  and  $X|Y=1 \sim N(\mu_1, \Sigma_1)$ . Then the Bayes rule is:  $h^*(x) = \begin{cases} 1 & \text{if } r_1^2 < r_0^2 + \log\left(\frac{\pi_1}{\pi_0}\right) + \log\left(\frac{|\Sigma_0|}{|\Sigma_1|}\right) \\ 0 & \text{otherwise} \end{cases}$

where  $r_i^2 = (x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i)$ ,  $i=0,1$  is the Mahalanobis distance.

equivalent Bayes rule:  $h^*(x) = \arg \max_k \delta_k(x)$

where  $\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k) + \log \pi_k$

and  $|A|$  denotes the determinant of a matrix  $A$ .



Pr: decision boundary if the classifier is quadratic → Quadratic Discriminant Analysis

sample estimates of  $\pi, \mu, \Sigma$  in place of the true

$$\pi_0 = \frac{1}{n} \sum_{i=1}^n (1 - Y_i), \quad \pi_1 = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$\mu_0 = \frac{1}{n_0} \sum_{i=1}^n X_i, \quad \mu_1 = \frac{1}{n_1} \sum_{i=1}^n X_i$$

$$\Sigma_0 = \frac{1}{n_0} \sum_{i=1}^n (X_i - \mu_0)(X_i - \mu_0)^T, \quad \Sigma_1 = \frac{1}{n_1} \sum_{i=1}^n (X_i - \mu_1)(X_i - \mu_1)^T$$

where  $n_0 = \sum_{i=1}^n (1 - Y_i)$  and  $n_1 = \sum_{i=1}^n Y_i$

A simplification occurs if we assume that  $\Sigma_0 = \Sigma_1 = \Sigma$

$$h^*(x) = \arg \max_k \delta_k(x), \quad \delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$$

$$\text{The MLE is } \hat{\delta} = \frac{n_0 \delta_0 + n_1 \delta_1}{n_0 + n_1}, \quad h^*(x) = \begin{cases} 1 & \text{if } \delta_1(x) > \delta_0(x) \\ 0 & \text{otherwise} \end{cases}$$

$$\text{where } \delta_j(x) = x^T \Sigma^{-1} \mu_j - \frac{1}{2} \mu_j^T \Sigma^{-1} \mu_j + \log \pi_j$$

There is the discriminant function. The decision boundary  $\{x : \delta_0(x) = \delta_1(x)\}$  is linear and called Linear Discrimination Analysis.

Linear & Logistic Regression: assume we have an estimator  $\hat{f}(x)$

$$\text{classification rule: } \hat{h}(x) = \begin{cases} 1 & \text{if } \hat{f}(x) > \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

Linear regression model:

$$Y = r(X) + \epsilon = b_0 + \sum_{j=1}^d b_j X_j + \epsilon \quad \text{where } E[\epsilon] = 0$$

Least squares estimate of  $\beta = (b_0, b_1, \dots, b_d)^T$  minimizes the residual sum of squares:

$$RSS(\beta) = \sum_{i=1}^n (Y_i - b_0 - \sum_{j=1}^d X_{ij} b_j)^2$$

$X$  is a  $N \times (d+1)$  matrix

$$X = \begin{bmatrix} 1 & X_{11} & \dots & X_{1d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & \dots & X_{nd} \end{bmatrix} \quad Y = (Y_1, \dots, Y_n)^T$$

$$RSS(\beta) = (Y - X\beta)^T (Y - X\beta)$$

$$\text{and model is } Y = X\beta + \epsilon, \quad \epsilon = (\epsilon_1, \dots, \epsilon_n)^T$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y \Rightarrow \hat{Y} = X \hat{\beta}$$

$$\text{for logistic regression: } r(x) = p(Y=1|X=x) = \frac{e^{b_0 + \sum_{j=1}^d b_j x_j}}{1 + e^{b_0 + \sum_{j=1}^d b_j x_j}}$$

Relationship Between Logistic Regression & LDA

The difference is in how we estimate the parameters: The joint density of a single observation is  $P(X, Y) = P(X|Y)P(Y) = P(Y|X)P(X)$ . In LDA we estimated the whole joint distribution by maximizing the likelihood:

$$\prod P(X_i, Y_i) = \underbrace{\prod P(X_i|Y_i)}_{\text{Gaussian}} \underbrace{\prod P(Y_i)}_{\text{Bernoulli}}$$

In logistic regression we maximized the conditional likelihood  $\prod P(Y_i|X_i)$  but we ignored the second term.

Since classification only requires knowing  $P(Y|X)$  we don't need to estimate the whole joint distribution. By the regression learner the marginal distribution  $P(X)$  is unspecified so it is more nonparametric than LDA (advantage over LDA).

The naive Bayes classifier is popular when  $X$  is high-dimensional & discrete.

The naive Bayes classifier is popular when  $X$  is high-dimensional & discrete.



Recall:

Let  $\hat{p}_j(j) = \frac{\sum I(Y_i = j, X_i \in A_j)}{\sum I(X_i \in A_j)}$

for  $j=1,2$  and  $j=0,1$

The impurity of the split  $t$  is:  $I(t) = \sum p_j \log p_j$

where  $g = 1 - \sum_{j=0}^2 \hat{p}_j(j)^2$ . This is known as the Gini Index (split to min impurity)

\* Usually the training error rate  $\hat{L}_n(h)$  is an estimate of the true error rate because biased downward.

2 ways to estimate the error rate: 1) cross validation and 2) probability inequalities

CV: splitting the data into: Training set  $T$  and Revalidation set  $V$

we estimate  $\hat{L}(h) = \frac{1}{m} \sum_{i \in V} I(h(X_i) \neq Y_i)$

or you can use k-fold CV: 1) Randomly divide data into k chunks of approx equal size. A common choice is  $k=10$ . 2) for  $k=1$  to  $k$ : a) delete chunks  $k$  from the data b) compute the classifier  $\hat{h}_k$  from the rest of the data c) use  $\hat{h}_k$  to predict the data in chunk  $k$ . Let  $\hat{L}_k$  denote the observed error rate:  $\hat{L}(h) = \frac{1}{k} \sum \hat{L}_k$

Probability Inequalities: This method is useful in the context of empirical risk minimization. Let  $\mathcal{H}$  set of linear classifiers.

$\hat{h} = \arg \min_{h \in \mathcal{H}} \hat{L}_n(h) = \arg \min_{h \in \mathcal{H}} (\frac{1}{n} \sum I(h(X_i) \neq Y_i))$

Use Hoeffding's Inequality: if  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$  for any  $\epsilon > 0$   $P(|\hat{p} - p| > \epsilon) \leq 2e^{-2n\epsilon^2}$

Theorem: Uniform convergence: Assume  $\mathcal{H}$  is finite and has  $m$  elements. Then  $P(\max_{h \in \mathcal{H}} |\hat{L}_n(h) - L(h)| > \epsilon) \leq 2m e^{-2n\epsilon^2}$

Theorem: Let

$\epsilon = \sqrt{\frac{2}{n} \log(\frac{2m}{\alpha})}$  Then  $\hat{L}_n(\hat{h}) \pm \epsilon$  is a  $1-\alpha$  CI for  $L(\hat{h})$

\* The larger is  $\mathcal{H}$ , the larger the CI for  $L(\hat{h})$   
 ↳ more likely to overfit but compensate with larger CI.

Support Vector Machines (SVM) - class of linear classifiers.

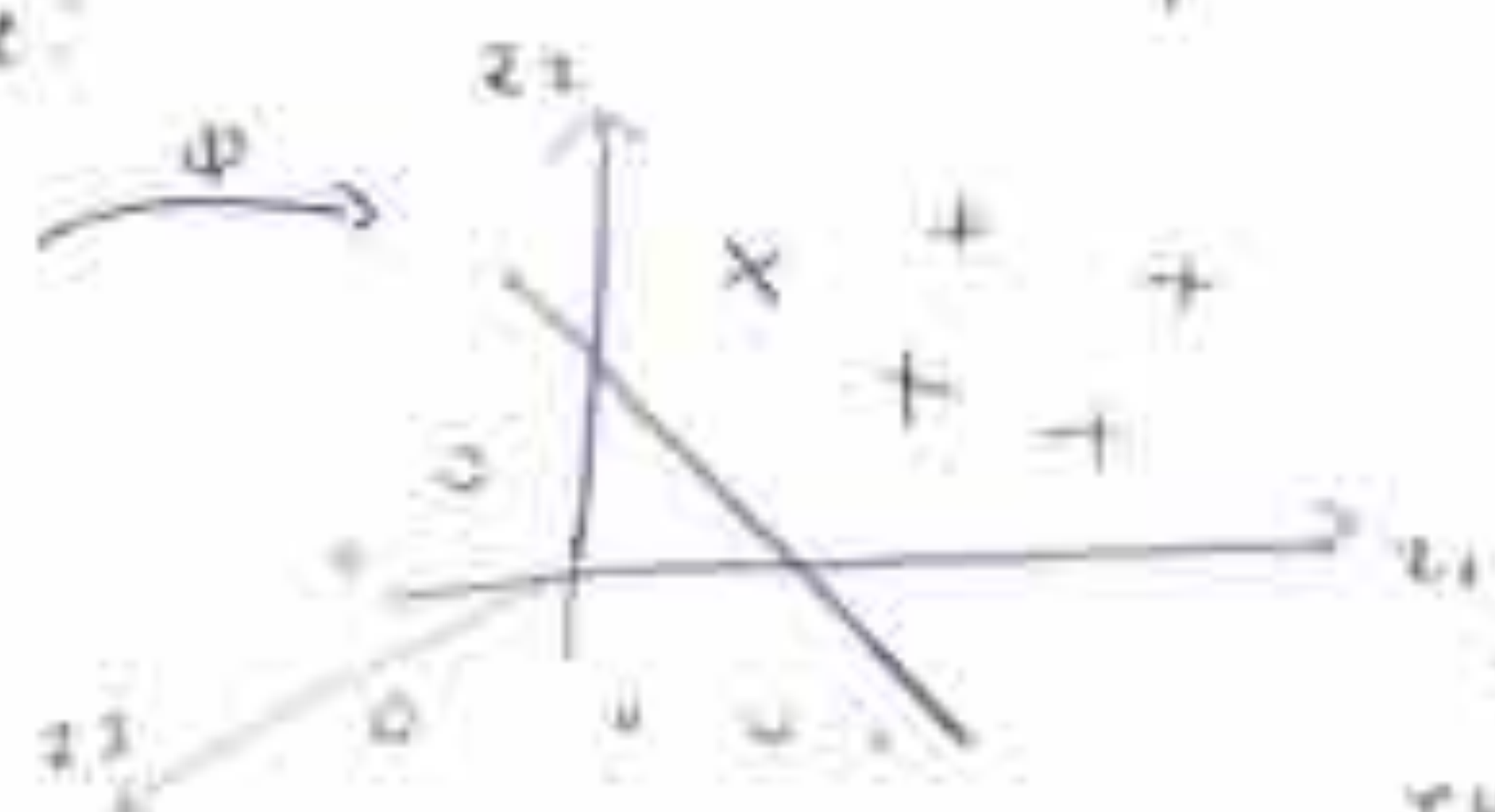
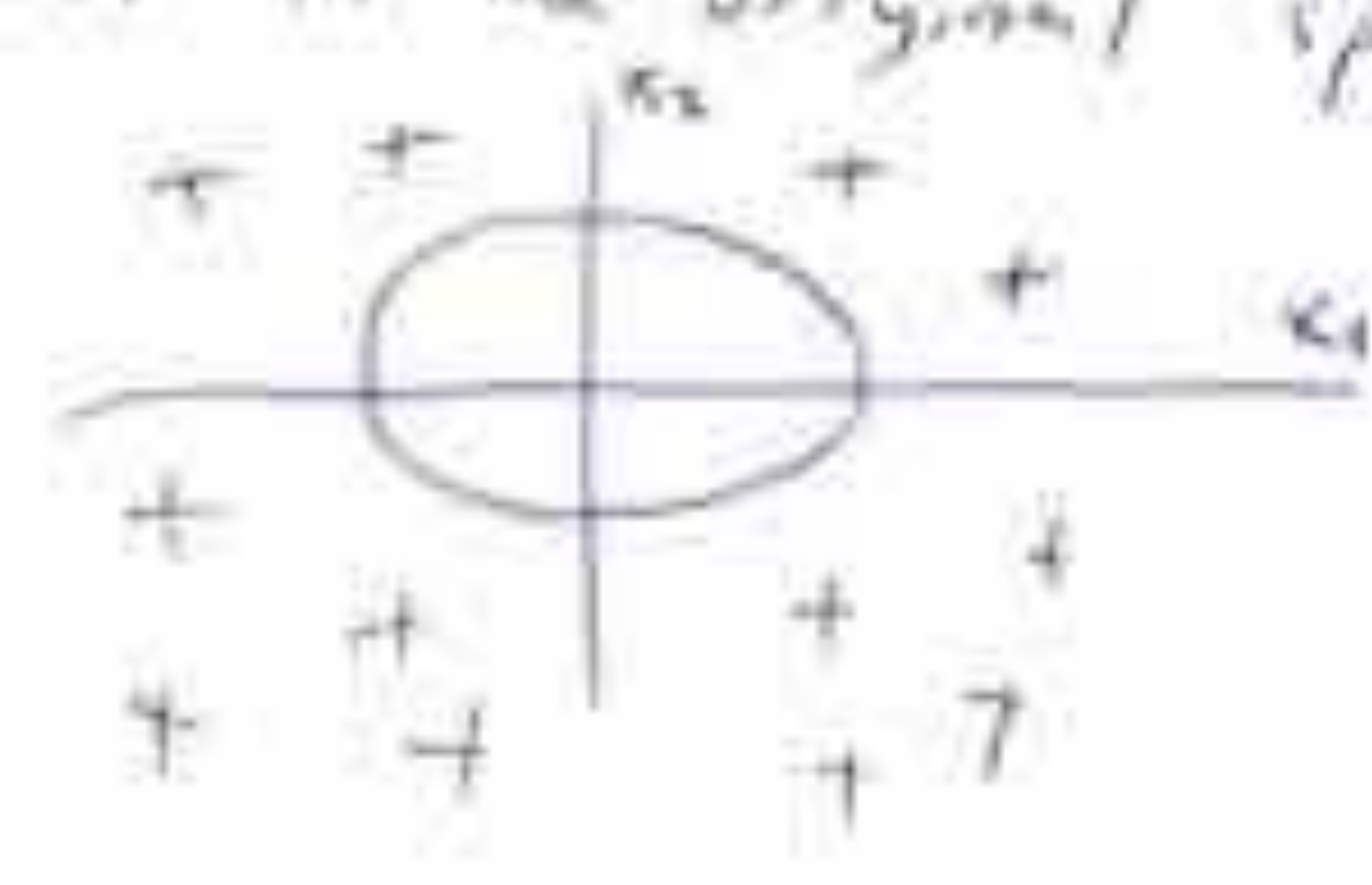
↳ Assume  $Y$  binary  $\in \{-1, 1\}$

Linear classifier:  $h(x) = \text{sign}(H(x))$  where  $x = (x_1, \dots, x_d)$

$H(x) = a_0 + \sum a_i x_i$  and  $\text{sign}(z) = \begin{cases} -1 & \text{if } z < 0 \\ 0 & \text{if } z = 0 \\ 1 & \text{if } z > 0 \end{cases}$

\* If data are linearly separable, there exists a  $2d$ -plane that perfectly separates the 2 classes.  $H(x) = a_0 + \sum a_i x_i$  s.t.  $Y_i H(x_i) \geq 1, i=1, \dots, n$   
 The margin is given by minimizing  $\frac{1}{2} \sum a_i^2$

Kernelisation: use to improve computationally a simple classifier  $h$ .  
 a linear classifier in a higher-dimensional space corresponds to a non-linear classifier in the original space.



$\langle \tilde{z}, \tilde{z} \rangle = \langle \phi(x), \phi(\tilde{x}) \rangle$   
 $= x_1^2 \tilde{x}_1^2 + 2x_1 \tilde{x}_1 x_2 \tilde{x}_2 + x_2^2 \tilde{x}_2^2$   
 $= (\langle x, \tilde{x} \rangle)^2 = k(x, \tilde{x})$

\* we can compute  $\langle \tilde{z}, \tilde{z} \rangle$  without ever computing  $\tilde{z}_i = \phi(x_i)$

Find a mapping  $\phi: X \rightarrow Z$



used kernels:

linear  $K(x, \tilde{x}) = (\langle x, \tilde{x} \rangle + a)^c$

sigmoid  $K(x, \tilde{x}) = \tanh(a \langle x, \tilde{x} \rangle + b)$

Gaussian  $K(x, \tilde{x}) = \exp(-\|x - \tilde{x}\|^2 / (2\sigma^2))$

for classifiers

- NN very simple

Bagging: method for reducing the variability of a classifier. Useful for highly nonlinear classifier such as trees

Boosting: start with a simple classifier and gradually improve it by refitting the data giving higher weight to misclassified samples. eg. AdaBoost

## Chapter 23: Probability Redux: Stochastic Processes

Imagine a sequence of DEPENDENT Random Variables. (eg temperature between days)  
A stochastic process  $\{X_t : t \in T\}$  is a collection of RV. The variables  $X_t$  take values in some set  $X$  called state space. The set  $T$  is called the index set (can be thought of time) and it can be discrete or continuous.  $\{0, 1, 2, 3, \dots\}$  or  $[0, +\infty)$

Markov Chains: A stochastic process for which the distribution of  $X_t$  depends only on  $X_{t-1}$ . Assuming that the state space is discrete,  $X = \{1, \dots, N\}$  or  $X = \{1, 2, \dots\}$  and the index set is  $T = \{1, 2, \dots\}$

→ The process  $\{X_n : n \in T\}$  is a Markov chain if  $P(X_n = x | X_0, \dots, X_{n-1}) = P(X_n = x | X_{n-1})$  for all  $n$  and for all  $x \in X$ .

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

$$X_0 \rightarrow X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_n \rightarrow \dots$$

Each variable has a single parent (the previous observation)

Transition Probabilities: jumping from one state to another. A MC is homogeneous if  $P(X_{n+1} = j | X_n = i)$  does not change with time. Thus, a homogeneous MC:

$P_{ij} = P(X_{n+1} = j | X_n = i)$ : transition probabilities. The matrix  $P$  whose  $(i, j)$  element is  $P_{ij}$  is called Transition matrix.

Properties of  $P$ : a)  $P_{ij} \geq 0$  and  $\sum_j P_{ij} = 1$  (each row is a PDF)

$$P_{m+n} = P_m P_n, \quad P_1 = P, \quad P_2 = P_{1+1} = P_1 P_1 = P \cdot P = P^2 \quad \text{so}$$
$$P_n = P^n = \underbrace{P \cdot P \cdot P \dots P}_{n \text{ times}}$$

\* We say that  $i$  reaches  $j$  if  $P_{ij}^{(n)} > 0$  for some  $n$  and we write  $i \rightarrow j$ . If  $i \rightarrow j$  and  $j \rightarrow i$  then we write  $i \leftrightarrow j$  and say  $i$  and  $j$  communicate.

\* If all states communicate with each other then the chain is called irreducible. A set of states is closed if once you enter that set of states you never leave. A closed set consisting of a single state is called an absorbing state.

eg.  $P = \begin{pmatrix} \frac{1}{2} & \frac{2}{3} & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & 0 & 0 \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \\ 1 & 0 & 0 & 0 \end{pmatrix}$  state 4 is an absorbing state.



State  $i$  is recurrent/persistent if  $P(X_n = i \text{ for some } n \geq 1 | X_0 = i) = 1$   
 otherwise state  $i$  is transient.  $\sum_n P_{ii}(n) = \infty$

\* If  $i$  is recurrent  $a_i = 1$ . The chain will eventually return to  $i$  and once it does, argue again that state  $a_i = 1$  the chain will return to state  $i$  again.  $[E(Y | X_0 = i) = \infty]$   
 \* If  $i$  is transient  $a_i < 1$ . When  $i$  is in state  $i$  there is a probability  $1 - a_i > 0$  that we never return to state  $i$ . Thus the prob that the chain is in state  $i$  exactly  $n$  times is  $a_i^{n-1}(1-a_i)$  geometric distribution with finite mean.

Facts about recurrence

1. If state  $i$  is recurrent and  $i \leftrightarrow j$  then  $j$  is recurrent.
2. If state  $i$  is transient and  $i \leftrightarrow j$  then  $j$  is transient.
3. A finite MC must have at least one recurrent state.
4. The states of a finite, irreducible MC are all recurrent.

Convergence of MC: Suppose that  $X_0 = i$ . Define the recurrence time:  $T_{ij} = \min\{n > 0 : X_n = j\}$  assuming  $X_n$  ever returns to state  $i$ , otherwise define  $T_{ij} = \infty$ . The mean recurrence time of a recurrent state  $i$  is:  $m_i = E(T_{ii}) = \sum_n n f_{ii}(n)$   
 where  $f_{ij}(n) = P(X_1 \neq j, X_2 \neq j, \dots, X_{n-1} \neq j, X_n = j | X_0 = i)$

A recurrent state is null if  $m_i = \infty$  otherwise it is called non-null/positive.

Formally the period of state  $i$  is  $d$  if  $P_{ii}(n) = 0$  whenever  $n$  is not divisible by  $d$  and  $d$  is the largest integer with this property. Thus  $d = \gcd\{n : P_{ii}(n) > 0\}$  where  $\gcd$ : "greater common divisor". State  $i$  is periodic if  $d(i) > 1$  and aperiodic if  $d(i) = 1$ .

Let  $\pi = (\pi_i; i \in X)$  be a vector of non-negative numbers that sum up to one. (can be thought as PDF). We say  $\pi$  is a stationary distribution if  $\pi = \pi P$ .

\* Just because a chain has a stationary distribution it doesn't mean it converges.

Poisson Processes:

$$P(X=x) = p(x, \lambda) = \frac{e^{-\lambda} \cdot \lambda^x}{x!}, \quad x=0, 1, 2, \dots, \quad E(X) = \lambda = V(X)$$

A Poisson process is a stochastic process  $\{X_t; t \in [0, \infty)\}$  with state space  $X = \{0, 1, 2, \dots\}$   
 s.t. 1)  $X(0) = 0$ , 2) For any  $0 = t_0 < t_1 < t_2 < \dots < t_n$  the increments

$X(t_1) - X(t_0), X(t_2) - X(t_1), \dots, X(t_n) - X(t_{n-1})$  are independent.

There is a function  $\lambda(t)$  s.t.

$$P(X(t+h) - X(t) = 1) = \lambda(t) \cdot h + o(h) \quad , \lambda(t) : \text{intensity function}$$

$$P(X(t+h) - X(t) \geq 2) = o(h)$$

→ A Poisson process with intensity function  $\lambda(t) = \lambda$  for some  $\lambda > 0$  is called a homogeneous Poisson process with rate  $\lambda$ . In this case,  $X(t) \sim \text{Poisson}(\lambda t)$ .

ITDS: Chapter 6 - Random Variable Generation

Pseudorandom: Consider a finite set  $M = \{0, 1, \dots, M-1\}$  and consider the sequence  $u_0, u_1, \dots$ . For every  $a \in M$  define  $N_a(n)$  as the number of  $u_i = a$  for  $i = 0, 1, 2, \dots, n-1$ . We call the sequence  $u_0, u_1, \dots$  pseudorandom on  $M$  if and only if for every  $a \in M$

$$\frac{N_a(n)}{n} \rightarrow \frac{1}{M}$$



Linear Generators start with a number  $U_0$  (that we call a random seed) and apply a map  $D$  to get the next number in the sequence.

is fixed and let  $D$  be a map define the dynamical system  $U_i = D(U_{i-1})$ ,  $i=1, \dots$  to the period of  $D$  started at  $U_0$  the smallest positive integers  $n$  s.t.  $U_n = U_0$ . smallest period  $T$  for all admissible starting points  $U_0$  is called the period of  $D$ .

modulus  $m=16$

multiplier  $a=3$

increment  $c=0$

$$U_{n+1} = (3U_n + 0) \bmod 16$$

1.  $U_0 = 2$

2.  $U_1 = 3 \cdot 2 \bmod 16 = 6 \bmod 16 = 6$

3.  $U_2 = 3 \cdot 6 \bmod 16 = 18 \bmod 16 = 2$

4.  $U_3 = 3 \cdot 2$

sequence  
2, 6, 2, 6, ...

## Sampling

### Algorithm 1: Accept-Reject Sampler

• Target Density  $f(x)$ : PDF from which you want to generate samples

• Sampling Density  $g(x)$ : A simpler distribution that is easy to sample from & satisfies:

$$f(x) \leq M g(x) \quad \forall x \quad \text{where } M > 0$$

Steps: 1. Initialization - initial state  $X_0$  from the sampling density  $g(x)$

2. Main Iteration Loop (Repeat until desired # of samples)

2.1 Draw a candidate  $X$  from sampling distr.  $g(x)$

2.2 Compute the acceptance rate  $r(x) = \frac{f(x)}{Mg(x)}$

2.3 Generate a Uniform RV  $U \sim \text{Uniform}(0, 1)$

2.4 Accept/Reject the proposal  $X$ : if  $U < r(x)$  accept and set  $X_{\text{acc}} = X$  otherwise reject and return to step 2.1, resampling  $Z$  from  $g(x)$

3. Repeat until desired # of samples obtained

### ITDS: Chapter 7 Finite MC

Let  $\{X_n, n \in \mathbb{N}\}$  be a homogeneous MC. Let the state space  $X = \{s_1, \dots, s_N\}$  be countable and let  $p_0$  be the PMF of  $X_0$ . Then the PMF  $p_n$  for  $X_n$  is:  $p_n = p_0 \cdot P^n$

### ITDS: Chapter 8 Pattern Recognition

In many ML text books that are practically oriented you will see the recommendation that the training/testing split should be 70/30. In pattern recognition problems this doesn't make much sense, it is better to use that to choose the number of determinations with large probability you want the bound to hold and use that to choose the number of samples to reach a 5% conf. interval.

precision:  $P(Y=1 | g(X)=1)$

Recall:  $P(g(X)=1 | Y=1)$

often used in medical testing and are then called sensitivity

### ITDS: Chapter 11 Dimensionality Reduction

#### SVD

Consider a line given by the unit vector  $v$  and consider a point  $x$  then the projection of  $x$  onto  $v$  is given by  $(v \cdot x) v$ . Let  $v$  be a unit vector. Consider the projection of each  $X_i$  onto  $v$ , but only consider the proportions  $Y_i = (X_i \cdot v)$

$$\bar{Y}_n = \frac{1}{n} \sum Y_i = \frac{1}{n} \sum X_i \cdot v = 0 \quad (\text{assumed zero empirical mean})$$



$$v_1 = \arg \max_{\|v\|=1} \frac{1}{n} \sum (Y_i - \bar{Y}_n)^2 = \arg \max_{\|v\|=1} \sum |X_i \cdot v|^2$$

let  $A$   $n \times n$  matrix with the rows  $x_i$  :  $\sum_{i=1}^n |X_i \cdot v|^2 = \|Av\|^2$   
 $\Rightarrow \arg \max_{\|v\|=1} \|Av\|$

\* The singular vectors are not necessarily unique, in fact if  $v$  is a singular vector, then so is  $-v$ . We can also have ties, we arbitrarily pick one

$$A = UDV^T$$

left singular vectors  $U$   $\rightarrow$  diagonal matrix with  $\sigma_i$

The power method:

$$A^T A = (UDV^T)^T (UDV^T) = (VDU^T UDV^T) = VD^2 V^T \quad V^T V = I \text{ (columns orthonormal)}$$

Principal Component Analysis:

It is a coordinate transformation from the original coordinates to the coordinate system given by the singular vectors.

$$PCA(A) = AV = UDV^T V = UD$$

Explained Variance: is how much % of the total variance is captured by our singular vectors.