

# Chapter 1

# Probability

## { Sample Space and Events .

- Sample Space ( $\Omega$ ) - Set of all possible outcomes of an experiment.  
 $w \in \Omega$  are called sample outcomes / realization or elements.
- Event ( $A$ ) - A subset of a sample space.

## { Probability

### Definition (Probability)

$P: \mathcal{P}(\Omega) \rightarrow \mathbb{R}$  is a probability distribution or a probability measure if

- i)  $P(A) \geq 0 \quad \forall A \in \mathcal{P}(\Omega)$
- ii)  $P(\Omega) = 1$
- iii) If  $A_1, A_2, \dots$  are disjoint,  $P(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$

↳ Countable additivity from measure theory ... ?!

Not all  $A \in \mathcal{P}(\Omega)$  are measurable  $\rightarrow \sigma\text{-algebra}$

## { View on Probability

◦ Frequentist.

◦ Bayesian

## § Independent Events

Definition (Independence of Event)

Two events  $A, B$  are independent if  $P(AB) = P(A)P(B)$

and we write  $A \perp\!\!\!\perp B$

A set of events  $\{A_i : i \in I\}$  is independent if  $P(\bigcap_{i \in J} A_i) = \prod_{i \in J} P(A_i)$ .

for all  $J \subseteq I$ , and denote

## § Conditional Probability

Definition (Conditional Probability)

If  $P(B) > 0$ , the conditional probability of  $A$  given  $B$  is

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

•  $P(\cdot|B)$  satisfies axiom of probability

$A, B$  are independent  $\Leftrightarrow P(A|B) = P(A)$

$$P(A_1 \cdots A_n) = P(A_1)P(A_2|A_1) \cdots P(A_n|A_1 \cdots A_{n-1})$$

## § Bayes' Theorem

Theorem (The Law of the Probability)

Let  $A_1, \dots, A_k$  be a partition of  $\Omega$ . Then for any event  $B$ ,

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

(pf) Let  $C_j = BA_j$ . Then,  $C_1, \dots, C_k$  are disjoint,  $B = \bigcup_{j=1}^k C_j$

$$\text{Hence, } P(B) = P\left(\bigcup_{j=1}^k C_j\right) = \sum_{j=1}^k P(C_j) = \sum_j P(BA_j) = \sum_j P(B|A_j)P(A_j).$$

Theorem (Bayes' Theorem).

Let  $A_1, \dots, A_k$  be a partition of  $\Omega$  such that  $P(A_i) > 0, P(B) > 0$ .

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

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

Remark  $P(A_i)$  - Prior probability of  $A_i$ .

$P(A_i|B)$  - Posterior probability of  $A_i$  given  $B$ .

$$P(A|B) = \frac{0.7 \times 0.9}{0.7 \times 0.9 + 0.2 \times 0.01 + 0.1 \times 0.01}$$
$$= \frac{0.63}{0.633}$$

1  
211

Random Variable

# Random Variables.

## Definition (Random Variable)

A random variable is a mapping  $X : \Omega \rightarrow \mathbb{R}$  that assigns a real nbr  $X(w)$  for each outcome  $w$

Given random variable  $X$  and  $A \subset \mathbb{R}$ ,  
let  $X^{-1}(A) := \{w \in \Omega \mid X(w) \in A\}$ .

Define  $P(X \in A) = |P(X^{-1}(A))|$ .

$P(X = x) = P(X^{-1}(x))$

Example Flip a coin twice  $X$ : nbr of heads.

$$P(X=0) = P(\{\text{TT}\}) = 1/4.$$

$$P(X=1) = P(\{\text{HT}, \text{TH}\}) = 1/2$$

$$P(X=2) = P(\{\text{HH}\}) = 1/4.$$

## Definition (Cumulative Distribution Function (CDF))

The CDF of a random variable  $X$  is  $F_X : \mathbb{R} \rightarrow [0, 1]$  by

$$F_X(x) := P(X \leq x) = P(X \in (-\infty, x])$$

Example Flip a coin twice,  $X$ : nbr of heads.

$$P(X=0) = \frac{1}{4}, P(X=1) = \frac{1}{2}, P(X=2) = \frac{1}{4}$$

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

Dynkin's theorem  
 $F: \mathbb{R} \rightarrow [0, 1]$

Theorem Let  $X$  have a CDF  $F$  and  $Y$  have a CDF  $G$ .

If  $F(x) = G(x) \ \forall x$ , then  $P(X \in A) = P(Y \in A) \ \forall A$ .

(pf) For  $a < b$ ,  $P(X \in (a, b)) = P(X^{-1}(a, \infty)) - P(X^{-1}(b, \infty))$

Let  $A_n = (at_n^{-1}, \infty) \Rightarrow \lim X^{-1}(A_n) = X^{-1}(-\infty, b)$ .

$\lim Y^{-1}(A_n) = Y^{-1}(-\infty, b)$ .

$$\begin{aligned} \text{Thus, } P(X^{-1}(a, \infty)) &= \lim P(X^{-1}(A_n)) = \lim P(X > a) \\ &= \lim 1 - P(X \leq a) = \lim 1 - P(Y \leq a) \\ &= \lim P(Y > a) = \lim P(Y^{-1}(A_n)) \\ &= P(Y^{-1}(a, \infty)). \end{aligned}$$

Similarly,  $P(X^{-1}(b, \infty)) = P(Y^{-1}(b, \infty))$ . Hence,  $P(X \in (a, b)) = P(Y \in (a, b))$

For  $A_n \in \mathcal{X}$ , suppose  $P(X \in A_n) = P(Y \in A_n) \ \forall n$ . Let  $B = \bigcup A_n$ ,  $C = \bigcap A_n$

$$P(X \in A^c) = 1 - P(X \in A) = 1 - P(Y \in A) = P(Y \in A^c).$$

$$P(X \in B) = \sum P(X \in B_n) =$$

$$A_n = \bigcup_{i \in n} A_i$$

for  $A \in \mathcal{X}$  if  $P(X \in A) = P(Y \in A)$

$$P(X \in A \cap A_e) =$$

For any open  $U \subseteq \mathbb{R}$ ,  $P(X \in U) = P(Y \in U)$ .

Let  $\mathcal{B}$  be collection of measurable sets  $A \in \mathcal{X}$  s.t.

$$P(X \in A) = P(Y \in A). \text{ (it suff. to show)}$$

$\mathcal{B}$  is a  $\sigma$ -algebra

$\mathcal{B}$  is closed under complement

Either closeness under countable intersection or union suff.

If closed under intersection b/w two sets, done.

$$\mathcal{I} := \{(a, b), [a, b), (a, b], [a, b]; \quad a < b\} \text{ . (closed under finite intersection)}$$

$$\mathcal{I} \subseteq \mathcal{D}, \quad \mathcal{B} = \sigma(\mathcal{I}) \subseteq \mathcal{D}$$

Theorem A function  $F: \mathbb{R} \rightarrow [0, 1]$  is a CDF for some probability  $P$  if and only if

$F$  satisfies (i)  $F$  is non-decreasing

$$(ii) \lim_{x \rightarrow -\infty} F(x) = 0, \lim_{x \rightarrow \infty} F(x) = 1$$

(iii)  $F$  is right continuous.

$$(f \Rightarrow) (i) P(X \leq x_1) = P(X^{-1}(-\infty, x_1]) \leq P(X^{-1}(-\infty, x_2)) = P(X \leq x_2)$$

$$f_x^{II}(x_1)$$

$$f_x^{II}(x_2)$$

$$\forall x_1 \leq x_2$$

(ii) Consider  $A_n = (-\infty, -n]$ .  $A_1 \supseteq A_2 \supseteq \dots$

$\cap A_n = \emptyset$ . Thus,  $\lim P(X \in A_n) = P(X \notin \emptyset) = 0$ .

i.e. Given  $\varepsilon > 0$ ,  $\exists N$  s.t.  $|P(X \in A_n)| < \varepsilon$

$$\Rightarrow 0 \leq F(n) = P(X \in A_n) < \varepsilon \Rightarrow |F(x)| < \varepsilon \quad \forall x \leq N$$

$$\rightarrow \lim_{x \rightarrow -\infty} F(x) = 0$$

$$\text{Similarly, } \lim_{x \rightarrow \infty} F(x) = 1$$

(iii) Take  $x_0 \in \mathbb{R}$ . If  $F(x_0) < 1$ , let  $2\varepsilon < 1 - F(x_0)$  be given.

Let  $x_1 = \inf \{x \mid F(x) > F(x_0) + \varepsilon\}$ .  $x_1 \geq x_0$ .

For any  $(x_n) \rightarrow x_0$  with  $x_n \geq x_0$ ,  $A_n := (-\infty, x_n]$ .

Then,  $P(X \in A_n) \rightarrow P(X \in A)$ . where  $A = (-\infty, x_0)$

$$\text{i.e. } F(x_n) \rightarrow F(x_0)$$

Definition (Discrete Random Var / Probability Mass Function)

$X$  is discrete if it takes countably many values.

PMF for  $X$  is  $f_X(x) = P(X=x)$

- $f_X(x) \geq 0 \quad \forall x \in \mathbb{R}, \sum_i f_X(x_i) = 1$

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

Definition (Continuous Random Variable / Probability distribution function)

A random variable  $X$  is continuous if  $\exists$  a function  $f_X$  such that  $f_X(x) \geq 0 \quad \forall x, \int_{-\infty}^{\infty} f_X(x) dx = 1$  and for  $a \leq b$ ,  $P(a < X < b) = \int_a^b f_X(x) dx$ .

- We have  $F_X(x) = \int_{-\infty}^x f_X(t) dt$ .

and  $f_X(x) = F'_X(x) \quad \forall x$  at which  $F_X$  is diff.

Lemma Let  $F$  be the CDF for a random variable  $X$ .

1.  $P(X=x) = F(x) - F(x^-)$

$$\text{IP}(a \leq x \leq b)$$

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

$$\text{IP}(\frac{1}{1} \leq x \leq b)$$

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

$$\text{IP}(a \leq x \leq b)$$

4. If  $X$  is continuous,  $F(b) - F(a) = P(a < X \leq b)$

## Definition (Inverse CDF / Quantile Function)

$X$ : random var. with CDF  $F$ .

$$F^{-1}(g) = \inf \{x; F(x) > g\}$$

$$F^{-1}(g) \text{ named } \begin{cases} \text{first quantile} & g = 1/4 \\ \text{second} & g = 1/2 \\ \text{third} & g = 3/4 \end{cases}$$

## Example (Discrete Random Variables)

### 1. Point mass distribution

$$X \sim f_a \text{ if } F(x) = \begin{cases} 0 & x < a \\ 1 & x \geq a \end{cases}$$

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

### 2. Discrete Uniform distribution, $k \geq 1$ .

$$f(x) = \begin{cases} 1/k & \text{for } x=1, \dots, k \\ 0 & \text{otherwise} \end{cases}$$

### 3. Bernoulli Distribution.

$P(X=1) = p$ ,  $P(X=0) = 1-p$ . for some  $p \in [0, 1]$

$X \sim \text{Bernoulli}(p)$ .

$$f(x) = p^x (1-p)^{1-x} \text{ for } x \in \{0, 1\}.$$

### 4. Binomial Distribution

$$f(x) = \begin{cases} \binom{n}{x} p^x (1-p)^{n-x} & \text{for } x=0, \dots, n \\ 0 & \text{otherwise} \end{cases}$$

$f(x) = P(X=x)$  be PMF. Then,  $X \sim \text{Binomial}(n, p)$ .

• If  $X_1 \sim \text{Binomial}(n_1, p)$ ,  $X_2 \sim \text{Binomial}(n_2, p)$ ,

then  $X_1 + X_2 \sim \text{Binomial}(n_1 + n_2, p)$

$$\therefore P(X_1 + X_2 = x) = \sum_{y=0}^x P(X_1 + X_2 = x | X_1 = y) P(X_1 = y)$$

### 5. Geometric Distribution

$X$  has a geometry distribution with  $p \in (0, 1)$ ,

denoted as  $X \sim \text{Geom}(p)$  if

$$P(X=k) = p(1-p)^{k-1} \quad \forall k \geq 1.$$

$$\therefore \sum P(X=k) = \sum p(1-p)^{k-1} = \frac{p}{1-(1-p)} = 1.$$

Means nbr of flips to the first head

## 6. Poisson Distribution

$X$  has a Poisson distribution with parameter  $\lambda$ , denoted as

$$X \sim \text{Poisson}(\lambda) \text{ if } P(X=x) = f(x) = e^{-\lambda} \frac{\lambda^x}{x!} \quad \theta_{x \geq 0}, x \in \mathbb{Z}$$

- $\sum_{x=0}^{\infty} f(x) = e^{-\lambda} \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} = e^{-\lambda} e^\lambda = 1$ .

- If  $X_1 \sim \text{Poisson}(\lambda_1)$ ,  $X_2 \sim \text{Poisson}(\lambda_2)$ , then  $X_1 + X_2 \sim \text{Poisson}(\lambda_1 + \lambda_2)$ .

- $E[X] = \sum x f(x) = e^{-\lambda} \sum_{x=1}^{\infty} \frac{\lambda^x}{(x-1)!} = \lambda e^{-\lambda} \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} = \lambda e^{-\lambda} e^\lambda = \lambda$ .

## § 2.4 Some Important Continuous Random Variable .

### 1. Uniform Distribution

$$X \sim \text{Uniform}(a, b) \text{ if } X \text{ has a PDF } f(x) = \begin{cases} \frac{1}{b-a} & x \in [a, b] \\ 0 & \text{otherwise.} \end{cases}$$

$$\Rightarrow \text{CDF } F(x) = \begin{cases} 0 & x < a \\ \frac{x-a}{b-a} & x \in [a, b] \\ 1 & x > b \end{cases}$$

PDF/CDF of  $Z$  is denoted

as  $f(z)$ ,  $F(z)$  respectively.

### 2. Normal(Gaussian) Distribution

$$X \sim N(\mu, \sigma^2) \text{ if } X \text{ has a PDF } f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2\sigma^2} (x-\mu)^2}. \quad x \in \mathbb{R}$$

- $X$  has a standard normal distribution if  $\mu=0, \sigma=1$ .

Standard Normal random variable is denoted by  $Z$ .

$\circ \mu$ : mean,  $\sigma$ : standard deviation.

Proposition / If  $X \sim N(\mu, \sigma^2)$ , then  $Z = (X - \mu) / \sigma \sim N(0, 1)$

(ii) If  $Z \sim N(0, 1)$ , then  $\mu + \sigma Z \sim N(\mu, \sigma^2)$

(iii) If  $X_i \sim N(\mu_i, \sigma_i^2)$ ,  $i=1, \dots, n$  are independent, then

$$\sum_{i=1}^n X_i \sim N\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2\right)$$

### 3. Exponential Distribution

$$X \sim \text{Exp}(\beta) \text{ if } f(x) = \frac{1}{\beta} e^{-\frac{x}{\beta}}, \forall x \geq 0$$

$$\int_0^\infty f(x) = \left[ e^{-\frac{x}{\beta}} \right]_0^\infty = 1.$$

$$\begin{aligned} \int_0^\infty x f(x) &= \int_0^\infty \frac{1}{\beta} x e^{-\frac{x}{\beta}} dx = \left[ x e^{-\frac{x}{\beta}} - \beta e^{-\frac{x}{\beta}} \right]_0^\infty \\ &= \beta \end{aligned}$$

### 4. Gamma Distribution

$$\text{For } \alpha > 0, P(x) := \int_0^\infty y^{\alpha-1} e^{-y} dy$$

$$X \sim \text{Gamma}(\alpha, \beta) \text{ if } f(x) = \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-\frac{x}{\beta}}, x > 0$$

## 5. Beta Distribution

$$X \sim \text{Beta}(\alpha, \beta) \text{ if } f(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}, \quad 0 < x < 1$$

## 6. t and Cauchy Distribution

$$X \sim t_v \text{ if } f(x) = \frac{\Gamma(\frac{v+1}{2})}{\Gamma(\frac{v}{2})} \frac{1}{(1 + \frac{x^2}{v})^{(v+1)/2}}$$

◦ Normal is t with  $v=\infty$ .

$$\text{Cauchy is t with } v=1. \quad f(x) = \frac{1}{\pi(1+x^2)}$$

## 7. $\chi^2$ Distribution

$$X \sim \chi_p^2 \text{ if } f(x) = \frac{1}{\Gamma(p/2) 2^{p/2}} x^{(p/2)-1} e^{-x/2}, \quad x > 0.$$

If  $Z_1, \dots, Z_p$  are independent standard normal random variables,  
then  $\sum_{i=1}^p Z_i^2 \sim \chi_p^2$

## § 2.5 Bivariate Distributions.

Definition (Joint CDF/Distribution).

- $F_{X,Y}(x,y) = P(X \leq x, Y \leq y)$ .

- Probability function Discrete: Joint mass f.m.  $f(x,y) = P(X=x, Y=y)$ .

Continuous: Joint PDF  $P((X,Y) \in A) = \iint_A f(x,y) dx dy$ .

## § 2.6 Marginal Distribution.

Definition (Marginal Distribution)

- If  $(X, Y)$  have joint distribution with mass function  $f_{X,Y}$ ,  
marginal mass function for  $X$  is  $f_X(x) = P(X=x) = \sum_y P(X=x, Y=y) = \sum_y f_{X,Y}(x,y)$   
 $Y$  is  $f_Y(y) = P(Y=y) = \sum_x P(X=x, Y=y) = \sum_x f_{X,Y}(x,y)$ .

- For continuous r.v., the marginal densities are

$$f_X(x) = \int f_{X,Y}(x,y) dy, \quad f_Y(y) = \int f_{X,Y}(x,y) dx.$$

Corresponding marginal distribution functions are denoted by  $F_X, F_Y$ .

## § 2.7 Independent Random Variables.

Definition (Independent RV).

Two RV  $X, Y$  are independent if  $\forall A, B$ ,

$$P(X \in A, Y \in B) = P(X \in A) P(Y \in B).$$

and write  $X \perp\!\!\!\perp Y$ . Otherwise  $X, Y$  are dependent:  $X \text{ and } Y$ .

Theorem Let  $X, Y$  have joint PDF  $f_{X,Y}$ . Then,  $X \perp\!\!\! \perp Y$  if and only if

$$f_{X,Y}(x,y) = f_X(x)f_Y(y) \quad \text{for all } x, y.$$

( $\Rightarrow$ )  $A_n = [x - \frac{1}{n}, x + \frac{1}{n}]$ ,  $B_n = [y - \frac{1}{n}, y + \frac{1}{n}]$ .

$$P(X \in A_n, Y \in B_n) = \iint_{A_n \times B_n} f_{X,Y}(x,y) dx dy.$$

$$P(X \in A_n) = \int_{A_n} f_X(x) dx \quad P(Y \in B_n) = \int_{B_n} f_Y(y) dy.$$

$$f_{X,Y}(x,y) = f_X(x)f_Y(y) \quad \text{a.e.}$$

$$\begin{aligned} (\Leftarrow) \quad P(X \in A, Y \in B) &= \iint_{A \times B} f_{X,Y}(x,y) dx dy \\ &= \iint_{A \times B} f_X(x)f_Y(y) dx dy = \int_A f_X(x) dx \int_B f_Y(y) dy \\ &= P(X \in A)P(Y \in B). \end{aligned}$$

Theorem Suppose range of  $X, Y$  is a (possibly infinite) rectangle.

If  $f_{X,Y}(x,y) = g(x)h(y)$  for some  $g, h$ , then  $X, Y$  are independent.

$$(P) \quad \iint_{\mathbb{R}^2} f dA = 1 \Rightarrow \int_{\mathbb{R}} g(x) \int_{\mathbb{R}} h(y) dy = 1.$$

Let  $P = \int_{\mathbb{R}} g(x) dx$ .  $Q = \int_{\mathbb{R}} h(y) dy$ .

$$f_x(x) = \int_y f(x, y) dy = \int_y g(x) h(y) dy = Qg(x)$$

$f_y(y) = P h(y)$  similarly.

$$\begin{aligned} f(x, y) &= g(x) h(y) = PQ g(x) h(y) \\ &= f_x(x) f_y(y). \end{aligned}$$

By thm 2.30, done.

## §2.8 Conditional Distribution

Definition (Conditional Probability Mass Function / PDF)

$$f_{x|y}(x|y) := P(X=x, Y=y) = \frac{P(X=x, Y=y)}{P(Y=y)} = \frac{f_{x,y}(x,y)}{f_y(y)} \quad \text{if } f_y(y) > 0.$$

$$f_{x|y}(x|y) := \frac{f_{x,y}(x,y)}{f_y(y)} \quad \text{assuming } f_y(y) > 0$$

$$\text{Then, } P(X \in A | Y=y) = \int_A f_{x|y}(x|y) dx.$$

## §2.9 Multivariate Distributions and IID Samples

Definition (IID)

If  $X_1, \dots, X_n$  are independent and has same marginal distribution with CDF  $F$ , we say  $X_1, \dots, X_n$  are IID and write  $X_1, \dots, X_n \sim F$ .

If  $F$  has a density  $f$ , we also write  $X_1, \dots, X_n \sim f$ .

## § 2.10 Two Important Multivariate Distribution

### Definition (Multinomial Distribution)

For  $p = (p_1, \dots, p_k)$ ,  $p_j \geq 0$ ,  $\sum_{j=1}^k p_j = 1$ ,  $X = (X_1, \dots, X_k)$

where  $X_j = \#$  of color  $j$  appear from  $n$  draw of ball

where probability of color  $j$  appearing is  $p_j$  on each draw.

We denote  $X \sim \text{Multinomial}(n, p)$ .

$$f(x) = \binom{n}{x_1, \dots, x_k} p_1^{x_1} \cdots p_k^{x_k} \text{ where } \binom{n}{x_1, \dots, x_k} = \frac{n!}{x_1! \cdots x_k!}$$

Lemma Suppose  $X \sim \text{Multinomial}(n, p)$ ,  $X = (X_1, \dots, X_k)$ ,

$P = (p_1, \dots, p_k)$ . The marginal distribution of  $X_j$  is Binomial( $n, p_j$ ).

(pf) WLOG,  $j = k$ .

$f_k(x_k) = \sum_{\substack{x_1, \dots, x_{k-1} \\ x_1 + \dots + x_k = n}} f(x)$ , which is probability # of color  $k = x_k$  after  $n$  draw with each draw

$$\begin{aligned} 1 &= 1^n = (\sum p_i)^n = (p_k + \sum_{i \neq k} p_i)^n \text{ probability of color } k = p_k, \\ &= (p_k + (1-p_k))^n. \end{aligned}$$

Coeff. of  $p_k^{x_k}$     coeff. of  $p^{x_k}$

## Definition (Multivariate Normal)

For vector  $\mu$ , matrix  $\Sigma \in S_+$ , let  $X = \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix}$

$$X \sim N(\mu, \Sigma) \text{ if } f(x; \mu, \Sigma) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

where  $|\Sigma|$  is the determinant of  $\Sigma$

$\mu=0, \Sigma=I \rightarrow \text{Standard Normal}$ .

Theorem If  $Z \sim N(0, I)$  and  $X = \mu + \Sigma^{1/2}Z$ , then  $X \sim N(\mu, \Sigma)$

Conversely, if  $X \sim N(\mu, \Sigma)$ , then  $\Sigma^{-1/2}(X-\mu) \sim N(0, I)$

$$S = r'(y) = \Sigma^{-1/2}(x-\mu) \quad f_Y(y) = f_X(S(y)) \left| \frac{dS(y)}{dy} \right|$$

(PFS)  $X = \mu + \Sigma^{1/2}Z = r(Z)$

$$P(X \in A) = P(\mu + \Sigma^{1/2}Z \in A) = P(Z \in S(A))$$

$$= \iint_{S(A)} f(z; 0, I) dz = \iint_{S(A)} \frac{1}{(2\pi)^{k/2}} e^{-\frac{1}{2}z^T z} dz$$

$$= \iint_{S(H)} \frac{1}{(2\pi)^{k/2}} e^{-\frac{1}{2}S(x)^T S(x)} \left| \frac{dS}{dx} \right| dx$$

$$= \iint_A \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)} dx$$

$$\Rightarrow X \sim N(\mu; \Sigma)$$

Similar for converse case

Theorem Let  $X = (X_a, X_b)$ ,  $\mu = (\mu_a, \mu_b)$ ,  $\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$

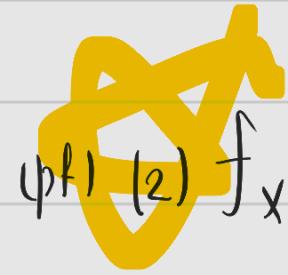
Then, (1) Marginal distribution of  $X_a$  is  $X_a \sim N(\mu_a, \Sigma_{aa})$

(2) Conditional distribution of  $X_b$  given  $X_a = x_a$  is

$$X_b | X_a = x_a \sim N(\mu_b + \Sigma_{ba}\Sigma_{aa}^{-1}(x_a - \mu_a), \Sigma_{bb} - \Sigma_{ba}\Sigma_{aa}^{-1}\Sigma_{ab})$$

(3) If  $a$  is a vector, then  $a^T X \sim N(a^T \mu, a^T \Sigma a)$

$$(4) V = (X - \mu)^T \Sigma^{-1} (X - \mu) \sim \chi_k^2$$



$$\begin{aligned}
 (1) \quad (2) \quad f_{X_b | X_a}(x_b | x_a) &= \frac{f_X(x)}{f_{X_a}(x_a)} \\
 &\stackrel{(1)}{=} \frac{(2\pi)^{k/2} |\Sigma|^{1/2}}{(2\pi)^{k_a/2} |\Sigma_a|^{1/2}} e^{-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)} \\
 &\qquad \qquad \qquad \frac{(2\pi)^{k_b/2} |\Sigma_b|^{1/2}}{(2\pi)^{k_b/2} |\Sigma_{ba}|^{1/2}} e^{-\frac{1}{2} (x_a - \mu_a)^T \Sigma_a^{-1} (x_a - \mu_a)}
 \end{aligned}$$

## § 2.11 Transformations of Random Variable.

Observation Let  $Y = r(X)$ , let  $A_Y := \{x \mid r(x) \leq y\}$ .

$$\begin{aligned}F_Y(y) &= P(Y \leq y) = P(r(X) \leq y) \\&= P(\{x \mid r(x) \leq y\}) = \int_{A_Y} f_X(x) dx.\end{aligned}$$

If  $r$  is strictly monotone,  $f_Y(y) = f_X(s(y)) \left| \frac{ds(y)}{dy} \right|$

## § 2.12 Transformations of Several Random Variables

Observation  $Z = r(X, Y)$ . Let  $A_Z = \{(x, y) \mid r(x, y) \leq z\}$ .

$$\begin{aligned}F_Z(z) &= P(Z \leq z) = P(r(X, Y) \leq z) \\&= P(\{(x, y) \mid r(x, y) \leq z\}) = \iint_{A_Z} f_{X,Y}(x, y) dx dy\end{aligned}$$

$$f_Z(z) = F_Z'(z)$$

Expectation

## § 3.1 Expectation of a Random Variable.

### Definition (Expectation Value)

The expectation of  $X$  is

$$E(X) = \int x dF(x) = \begin{cases} \sum_x xf(x) & \text{if } X \text{ is discrete} \\ \int xf(x) dx & \text{if } X \text{ is continuous.} \end{cases}$$

Assuming the sum (or integral) is well defined.

We denote  $E(X) = EX = \int x dF(x) = \mu = \mu_X$ .

$k^{\text{th}}$  central moment:  $E((X-\mu)^k)$

Example 1.  $X \sim \text{Bernoulli}(p)$ .  $E(X) = \sum_{x=0}^1 xf(x) = 0 \times (1-p) + 1 \times p = p$

2. Flip a fair coin twice,  $X = \# \text{ of heads}$ .

$$\begin{aligned} E(X) &= \sum_x xf_x(x) = 0 \times f_0 + 1 \times f_1 + 2 \times f_2 \\ &= 0 \times \frac{1}{4} + 1 \times \frac{1}{2} + 2 \times \frac{1}{4} = 1. \end{aligned}$$

$$\begin{aligned} 3. X \sim \text{Uniform}(-1, 3). \quad E(X) &= \int x dF_X(x) = \int x f_X(x) dx \\ &= \frac{1}{4} \int_{-1}^3 x dx = 1 \end{aligned}$$

## Theorem (Linearity of the Expectation)

If  $X_1, \dots, X_n$  are RV,  $a_1, \dots, a_n$  are constants,

$$\mathbb{E}(\sum a_i X_i) = \sum a_i \mathbb{E}(X_i)$$

(pf) By linearity of integration.

## Theorem (Multiplicativity of the Expectation)

Let  $X_1, \dots, X_n$  be independent. Then,

$$\mathbb{E}(\prod_{i=1}^n X_i) = \prod_{i=1}^n \mathbb{E}(X_i)$$

$$(pf) \quad \mathbb{E}(\prod X_i) = \int \prod x_i f_{X_i}(x_1, \dots, x_n) dx \quad \begin{matrix} x = (x_1, \dots, x_n) \\ X = (X_1, \dots, X_n) \end{matrix}$$

$$= \int \prod x_i \prod f_{X_i}(x_i) dx \\ = \int \prod x_i f_{X_i}(x_i) dx = \prod \int x_i f_{X_i}(x_i) dx_i.$$

## Theorem (Expectation of Transformed RV)

Let  $Y = r(x)$ . Then,  $\mathbb{E}(Y) = \mathbb{E}(r(x)) = \int r(x) dF_X(x)$

$$(pf) \quad \mathbb{E}(Y) = \mathbb{E}(r(x)) = \int y f_Y(y) dy = \int r(x) \left| \frac{ds(y)}{dy} \right| f_X(s(y)) dy \\ = \int r(x) f_X(x) dx.$$

### § 3.3 Variance and Covariance.

#### Definition (Variance)

Let  $X$  be a RV with mean  $\mu$ . The Variance of  $X$  denoted by  $\sigma^2$  or  $\text{Var}(X)$  or  $\text{IV}(X)$  is

$$\sigma^2 := \mathbb{E}((X-\mu)^2) = \int (x-\mu)^2 dF(x), \text{ assuming the expectation exists.}$$

The standard deviation is  $\text{sd}(X) := \sqrt{\text{IV}(X)}$  denoted by  $\sigma$  and  $\sigma_X$

Theorem Assuming variance is well defined,

$$1. \text{IV}(X) = \mathbb{E}(X^2) - \mu^2$$

$$2. \text{If } a, b \text{ are constants, then } \text{IV}(ax+b) = a^2 \text{IV}(X)$$

$$3. \text{If } X_1, \dots, X_n \text{ are independent, } a_1, \dots, a_n \text{ are constants,}\\ \text{then } \text{IV}\left(\sum_{i=1}^n a_i X_i\right) = \sum_{i=1}^n a_i^2 \text{IV}(X_i)$$

$$(P1) 1. \text{IV}(X) = \mathbb{E}(X-\mu)^2 = \mathbb{E}(X^2 - 2\mu X + \mu^2) = \mathbb{E}(X^2) - 2\mu \mathbb{E}(X) + \mu^2 \\ = \mathbb{E}(X^2) - 2\mu^2 + \mu^2 = \mathbb{E}(X^2) - \mu^2$$

$$2. \mathbb{E}(ax+b) = a\mathbb{E}(X) + b = a\mu + b$$

$$\text{IV}(ax+b) = \mathbb{E}((ax+b) - (a\mu+b))^2 = \mathbb{E}(a^2(X-\mu)^2) \\ = a^2 \mathbb{E}(X-\mu)^2 = a^2 \text{IV}(X)$$

$$3. \text{IV}\left(\sum X_i\right) = \mathbb{E}\left(\sum X_i - \sum \mu_i\right)^2 = \mathbb{E}\left(\left(\sum X_i - \sum \mu_i\right)^2\right) \\ = \sum \text{IV}(X_i) \text{ by induction on } n$$

(Use P2)

Example  $X \sim \text{Binomial}(n, p)$ ,  $X = \sum X_i$ ,  $X_i = 1$  if toss  $i$  = head  
 $0$  otherwise.

$X = \sum X_i$ , and  $X_i$  are independent.  $P(X_i = 1) = p$ ,  $P(X_i = 0) = 1 - p$ .

$$E(X_i) = p$$

$$E(X_i^2) = p \times 1^2 + (1-p) \times 0^2 = p$$

$$\text{Thus, } V(X_i) = E(X_i^2) - p^2 = p - p^2 = p(1-p).$$

$$V(X) = V(\sum X_i) = \sum p(1-p) = np(1-p).$$

## Definition (Sample Mean / Sample Variance)

If  $X_1, \dots, X_n$  are random variables, then

$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  is the sample mean

$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  is the sample variance

↳ To ensure unbiased estimator (following theorem)

Theorem Let  $X_1, \dots, X_n$  be IID.  $\mu = E(X_i)$ ,  $\sigma^2 = V(X_i)$ . Then,

$$E(\bar{X}_n) = \mu, V(\bar{X}_n) = \frac{\sigma^2}{n}, E(S_n^2) = \sigma^2$$

$$\text{pf)} E(\bar{X}_n) = E\left(\frac{1}{n} \sum X_i\right) = \frac{1}{n} \sum E(X_i) = \frac{1}{n} n \mu = \mu$$

$$V(\bar{X}_n) = V\left(\frac{1}{n} \sum X_i\right) = \sum \frac{1}{n^2} V(X_i) = n \times \frac{\sigma^2}{n^2} = \frac{\sigma^2}{n}$$

$$E(S_n^2) = E\left(\frac{1}{n-1} \sum (X_i - \bar{X}_n)^2\right) = \frac{1}{n-1} \sum E(X_i - \bar{X}_n)^2$$

$$X_n - \bar{X}_n = \frac{n-1}{n} X_n - \frac{1}{n} \sum_{i=1}^{n-1} X_i, E(X_n - \bar{X}_n) = 0.$$

$$\begin{aligned}
 \Rightarrow \mathbb{E}(X_n - \bar{X}_n)^2 &= \mathbb{V}(X_n - \bar{X}_n)^2 = \mathbb{V}\left(\frac{n-1}{n}X_n - \frac{1}{n}\sum_{i=1}^{n-1}X_i\right) \\
 &= \mathbb{V}\left(\frac{n-1}{n}X_n - \frac{n-1}{n}\frac{1}{n-1}\sum_{i=1}^{n-1}X_i\right) \\
 &= \left(\frac{n-1}{n}\right)^2 \left(\mathbb{V}(X_n) + \mathbb{V}\left(\frac{1}{n-1}\sum_{i=1}^{n-1}X_i\right)\right) \\
 &= \left(\frac{n-1}{n}\right)^2 \left(6^2 + \frac{6^2}{n-1}\right) = \frac{n-1}{n}6^2.
 \end{aligned}$$

$$\begin{aligned}
 \mathbb{E}(S_n^2) &= \frac{1}{n-1} \sum_{i=1}^n \mathbb{E}(X_i - \bar{X}_n)^2 \\
 &= \frac{1}{n-1} \sum_{i=1}^n \frac{n-1}{n}6^2 = 6^2
 \end{aligned}$$

## Definition (Covariance)

For random variables  $X, Y$  with mean  $\mu_X, \mu_Y$ , and standard deviation  $\sigma_X, \sigma_Y$  respectively, the covariance btwn  $X$  and  $Y$  is

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

and the correlation is

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

Theorem  $\text{Cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y)$   
 $-1 \leq \rho(X, Y) \leq 1$ .

(pf).  $\text{Cov}(X, Y) = \mathbb{E}((X - \mu_X)(Y - \mu_Y))$

$$\begin{aligned}
 &= E(XY - \mu_x Y - \mu_y X + \mu_x \mu_y) \\
 &= E(XY) - \mu_x \mu_y - \mu_y \mu_x + \mu_x \mu_y \\
 &= E(XY) - \mu_x \mu_y = E(XY) - E(X)E(Y).
 \end{aligned}$$

$$\begin{aligned}
 |E|(X-\mu_x)(Y-\mu_y)| &= \int |(X-\mu_x)(Y-\mu_y)| f_{XY}(x,y) dx dy \\
 &\leq \left( \int (X-\mu_x)^2 dx \right)^{1/2} \left( \int (Y-\mu_y)^2 dy \right)^{1/2} \text{ by C-S ineq.} \\
 &= \sigma_x \sigma_y.
 \end{aligned}$$

Thus,  $|Cov(X,Y)| \leq \sigma_x \sigma_y$ . i.e.,  $\rho_{X,Y} = \frac{Cov(X,Y)}{\sigma_x \sigma_y} \leq 1$

• If  $Y = aX + b$ ,  $\rho(X,Y) = 1$  if  $a > 0$   
 $\rho(Y,Y) = -1$  if  $a < 0$ .

•  $X, Y$  are independent  $\Rightarrow Cov(X,Y) = \rho = 0$ .

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

More generally  $V(\sum a_i X_i) = \sum a_i V(X_i) + 2 \sum_{i < j} a_i a_j \text{Cov}(X_i, X_j)$

$$\begin{aligned} (\text{pf}) \quad V(ax+by) &= E((ax+by) - (a\mu_x+b\mu_y))^2 \\ &= E((a(x-\mu_x)+b(y-\mu_y))^2) \\ &= a^2 E((x-\mu_x)^2) + b^2 E((y-\mu_y)^2) + 2ab E((x-\mu_x)(y-\mu_y)) \\ &= a^2 V(X) + b^2 V(Y) + 2ab \text{Cov}(X, Y). \end{aligned}$$

Use induction with fact that

$$\begin{aligned} \text{Cov}(X, Y+Z) &= E((X-\mu_X)(Y+Z - (\mu_Y+\mu_Z))) \\ &= E((X-\mu_X)(Y-\mu_Y)) + E((X-\mu_X)(Z-\mu_Z)) \\ &= \text{Cov}(X, Y) + \text{Cov}(X, Z) \end{aligned}$$

## §3.5 Conditional Expectation

### Definition (Conditional Expectation)

The conditional expectation of  $X$  given  $Y=y$  is

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

If  $r(x,y)$  is a fn of  $x,y$ ,

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

•  $y \mapsto \mathbb{E}(X|Y=y)$  is a fn of  $y$ .

$\mathbb{E}(X|Y)$  is the random variable whose value is  $\mathbb{E}(X|Y=y)$  when  $Y=y$ .

R.V.  $X: \Omega \rightarrow E$ ,  $\mathbb{P}: \mathcal{F}(\Omega) \rightarrow \mathbb{R}$

$$\mathbb{P}(X \in A) = \mathbb{P}(X^{-1}(A))$$

### Theorem (The Rule of Iterated Expectations).

For R.V.  $X, Y$ , assuming the expectations exists,

$$\mathbb{E}(\mathbb{E}(Y|X)) = \mathbb{E}(Y), \quad \mathbb{E}(\mathbb{E}(X|Y)) = \mathbb{E}(X).$$

More generally

$$\mathbb{E}(\mathbb{E}(r(X,Y)|X)) = \mathbb{E}(r(X,Y))$$

$$\begin{aligned}
 (\text{pf}) \quad \mathbb{E}(\mathbb{E}(Y|X)) &= \int \mathbb{E}(Y|X=x) f_X(x) dx \\
 &= \int \int y f_{Y|X}(y|x) dy f_X(x) dx \\
 &= \iint y f_{Y|X}(y|x) f_X(x) dy dx \\
 &= \iint y f(x,y) dy dx = \mathbb{E}(Y).
 \end{aligned}$$

### Definition (Conditional Variance)

The conditional variance is

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

where  $\mu(x) = \mathbb{E}(Y|X=x)$

Theorem  $V(Y) = \mathbb{E} V(Y|X) + V \mathbb{E}(Y|X)$

★

$$\begin{aligned}
 (\text{pf}) \quad \mathbb{E} V(Y|X) &= \mathbb{E} \int (y - \mu(X))^2 f(y|X) dy \\
 &= \int \int (y - \mu(x))^2 f(y|x) dy f_X(x) dx \\
 &= \iint (y - \mu(x))^2 f(x,y) dx dy \quad (y - \mu_y)^2
 \end{aligned}$$

$$V \mathbb{E}(Y|X) = V \int y f(y|X) dy \quad \mathbb{E} \mathbb{E}(Y|X) = \mathbb{E}(Y)$$

$$\begin{aligned} &= \int \left( \int y f(y|x) dy - \mu_y \right)^2 f_x(x) dx \\ &= \int (\mu(x) - \mu_y)^2 f_x(x) dx \\ &= \int \left( \int (\mu(x) - \mu_y) f(y|x) dy \right) f_x(x) dx \end{aligned}$$

## §3.6 Moment Generating Functions

### Definition (Moment Generating Function)

The moment generating function (MGF) or Laplace transform, of  $X$  is  $\psi_X(t) = \mathbb{E}(e^{tx}) = \int e^{tx} dF(x)$ .

where  $t \in \mathbb{R}$

**Observation**  $\psi'(0) = \left[ \frac{d}{dt} \mathbb{E}[e^{tx}] \right]_{t=0} = \mathbb{E}[xe^{tx}]_{t=0} = \mathbb{E}(X)$ .  
 $\psi^{(k)}(0) = \mathbb{E}(X^k)$  ( $k^{\text{th}}$  moment).

### Lemma (Properties of the MGF)

(1) If  $Y = ax + b$ , then  $\psi_Y(t) = e^{bt} \psi_X(at)$

(2) If  $X_1, \dots, X_n$  are independent and  $Y = \sum X_i$ , then  $\psi_Y(t) = \prod \psi_{X_i}(t)$

where  $\psi_{X_i}$  is the MGF of  $X_i$ .

$$(1) \quad \psi_Y(t) = \mathbb{E} e^{tY} = \mathbb{E} e^{t(ax+b)} = e^{bt} \mathbb{E} e^{atX} = e^{bt} \psi_X(at)$$

$$(2) \quad \psi_Y(t) = \mathbb{E}(e^{t\sum X_i}) = \mathbb{E}(\prod e^{tX_i}) = \prod \mathbb{E}(e^{tX_i}) = \prod \psi_{X_i}(t)$$

Example  $X \sim \text{Binomial}(n, p)$ ,  $X = \sum_{i=1}^n X_i$ , where  $P(X_i=1)=p$ ,  $P(X_i=0)=1-p$ .

$$\Psi_x(t) = \mathbb{E} e^{X_i t} = p e^t + (1-p) \cdot 1 = p e^t + q \quad \text{where } q = 1-p.$$

$$\Psi_X(t) = \prod \Psi_x(t) = (p e^t + q)^n.$$



Theorem Let  $X, Y$  be R.V. If  $\Psi_X(t) = \Psi_Y(t)$  for  $\forall t$  in some

open interval around 0, then  $X \stackrel{d}{=} Y$  (i.e. they have the same distrib.)

Example.  $X_1 \sim \text{Binomial}(n_1, p)$ ,  $X_2 \sim \text{Binomial}(n_2, p)$ ,  $Y = X_1 + X_2$ .

Then  $\Psi_Y(t) = \Psi_{X_1}(t)\Psi_{X_2}(t) = (pe^t + q)^{n_1} (pe^t + q)^{n_2} = (pe^t + q)^{n_1 + n_2}$ .

$Y$  has MGF of Binomial( $n_1 + n_2, p$ ).

Hence by thm,  $Y \sim \text{Binomial}(n_1 + n_2, p)$

[Note] (Moment Generation Functions for Some Common Distributions)

Distribution

MGF  $\Psi(t)$

Bernoulli ( $p$ )

$$pe^t + (1-p)$$

Binomial( $n, p$ )

$$(pe^t + (1-p))^n$$

Poisson( $\lambda$ )

$$e^{\lambda(e^t - 1)}$$

Normal ( $\mu, \sigma$ )

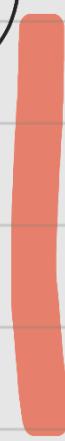
$$e^{\mu t + \frac{\sigma^2 t^2}{2}}$$

Gamma ( $\alpha, \beta$ )

$$\left(\frac{1}{1-\beta t}\right)^\alpha \quad \text{for } t < \frac{1}{\beta}.$$

$\therefore$  ) Poisson( $\lambda$ ):

L Inequalities



## § 4.1 Probability Inequality

[Theorem] (Markov's Inequality)

Let  $X$  be a nonnegative R.V. and  $\exists E(X)$ . For any  $t > 0$ ,

$$P(X > t) \leq \frac{E(X)}{t}.$$

$$\begin{aligned} (\text{pf}) \quad E(X) &= \int_0^\infty xf(x)dx = \int_0^t xf(x)dx + \int_t^\infty xf(x)dx \\ &\geq \int_t^\infty xf(x)dx \geq t \int_t^\infty f(x)dx = t P(X > t). \end{aligned}$$

[Theorem] (Chernoff's Inequality)

Let  $X$  be a R.V. then,  $P(X > \varepsilon) \leq \inf_{t > 0} e^{-t\varepsilon} E(e^{tx})$ .

(pf)  $t > 0 \quad P(X > \varepsilon) = P(e^{tx} > e^{t\varepsilon}) \leq e^{-t\varepsilon} E(e^{tx}) \quad (\text{by Markov's ineq.})$

Thus,  $P(X > \varepsilon) \leq \inf_{t > 0} e^{-t\varepsilon} E(e^{tx})$

[Theorem] (Chebyshev's Inequality)

Let  $\mu = E(X)$ ,  $\sigma^2 = V(X)$ . Then

$P(|X-\mu| \geq t) \leq \frac{\sigma^2}{t^2}$  and  $P(|Z| \geq k) \leq \frac{1}{k^2}$ , where  $Z = \frac{X-\mu}{\sigma}$ .

(pf) Second ineq. follows from the first by setting  $t=k\sigma$ .

$$P(|X-\mu| \geq t) = P((X-\mu)^2 \geq t^2) \leq \frac{\mathbb{E}(X-\mu)^2}{t^2} = \frac{6^2}{t^2}.$$

$\hookrightarrow$  Markov's ineq.

[Example]  $V(\bar{X}_n) = \frac{V(X_i)}{n} = \frac{P(1-p)}{n}$ .

By Chebyshev,  $P(|\bar{X}_n - p| > \varepsilon) \leq \frac{V(\bar{X}_n)}{\varepsilon^2} = \frac{P(1-p)}{n\varepsilon^2} \leq \frac{1}{4n\varepsilon^2}$ .  
 $\hookrightarrow p(1-p) \leq \frac{1}{4}$  if  $p=0.5$

$\rightarrow$  Used for bounded random variable.

[Theorem] (Hoeffding's Inequality)

Let  $Y_1, \dots, Y_n$  be independent such that  $E(Y_i) = 0$ ,  $a_i \leq Y_i \leq b_i$ .

Let  $\varepsilon > 0$ . Then for any  $t > 0$ ,

$$P\left(\sum_{i=1}^n Y_i \geq \varepsilon\right) \leq e^{-t\varepsilon} \prod_{i=1}^n e^{t(b_i-a_i)/8}$$

$$\text{In particular, } P\left(\sum_{i=1}^n Y_i \geq \varepsilon\right) \leq \inf_{t>0} e^{-t\varepsilon} \prod_{i=1}^n e^{t^2(b_i-a_i)^2/8}$$

(pf)  $P(\sum Y_i \geq \varepsilon) \leq e^{-t\varepsilon} E(e^{t\sum Y_i})$  (Chernoff ineq.)  
 $= e^{-t\varepsilon} \prod_{i=1}^n E e^{tY_i}$ .

$$\text{Let } Y_i = \alpha b_i + (1-\alpha) a_i, \alpha = \frac{Y_i - a_i}{b_i - a_i}.$$

$$e^{tY_i} = e^{t(\alpha b_i + (1-\alpha)a_i)} \leq \alpha e^{tb_i} + (1-\alpha) e^{ta_i} = \frac{Y_i - a_i}{b_i - a_i} e^{tb_i} + \frac{b_i - Y_i}{b_i - a_i} e^{ta_i}$$

$$E(e^{tY_i}) \leq \frac{-a_i}{b_i - a_i} e^{tb_i} + \frac{b_i}{b_i - a_i} e^{ta_i} = e^{g(u)}, \text{ where}$$

$$u = t(b_i - a_i), g(u) = -\gamma u + \log(1 - \gamma + \gamma e^u), \gamma = -\frac{a_i}{b_i - a_i}.$$

$$g(0) = g'(0) = 0, g''(u) \leq \frac{1}{4} \Theta u > 0. \text{ By Taylor, } \exists \xi \in (0, u) \text{ s.t.}$$

$$g(u) = \frac{u^2}{2}, g''(\xi) \leq \frac{u^2}{8} = \frac{t^2(b_i - a_i)^2}{8}.$$

$$\text{Thus, } E(e^{tY_i}) \leq e^{g(u)} = e^{t^2(b_i - a_i)^2/8}$$

**Theorem** Let  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ . Then for  $\varepsilon > 0$ ,

$$P(|\bar{X}_n - p| > \varepsilon) \leq 2e^{-2n\varepsilon^2},$$

$$\text{where } \bar{X}_n = n^{-1} \sum_{i=1}^n X_i.$$

(pf)  $P(\bar{X}_n - p > \varepsilon) = P(\sum Y_i > n\varepsilon)$ , where  $Y_i = X_i - p$  has  $EY_i = 0$   
 $-p \leq Y_i \leq 1-p$ .

$$\begin{aligned} \text{By Hoeffding, } P(\bar{X}_n - p > \varepsilon) &\leq e^{-t n \varepsilon} \prod e^{t^2/8} \\ &\leq e^{-4n\varepsilon^2} \cdot e^{2n\varepsilon^2} = e^{-2n\varepsilon^2} (t=4\varepsilon) \end{aligned}$$

Similarly,  $P(\bar{X}_n - p < \varepsilon) = e^{-2n\varepsilon^2}$  so that  $P(|\bar{X}_n - p| > \varepsilon) \leq 2e^{-2n\varepsilon^2}$ .

**Theorem** (Mill's Inequality)

$$\text{Let } Z \sim N(0, 1). \text{ Then } P(|Z| > t) \leq \sqrt{\frac{2}{\pi}} \frac{e^{-t^2/2}}{t}$$

$$(pf) p(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}, \quad p'(z) = -z p(z).$$

$$P(Z > t) = \int_t^\infty p(z) dz = \int_t^\infty -\frac{p'(z)}{z} dz \leq -\frac{1}{t} \int_t^\infty p'(z) dz$$

$$= -\frac{1}{t} p(z) \Big|_t^\infty = \frac{1}{t} p(t) = \frac{1}{\sqrt{2\pi}} \frac{e^{-t^2/2}}{t} = P(Z < -t).$$

$$\text{Thus, } P(|Z| > t) = P(Z > t) + P(Z < -t) = \sqrt{\frac{2}{\pi}} \frac{e^{-t^2/2}}{t}.$$

## § 4.2 Inequalities For Expectations

**Theorem** (Cauchy-Schwarz Inequality)

If  $X, Y$  have finite variance, then

$$E|XY| \leq \sqrt{E(X^2) E(Y^2)}$$

$$(pf) E((X+Y))^2 = t^2(E(X^2) + 2tE(XY) + E(Y^2)) \geq 0 \quad \forall t.$$

$$\text{Thus, } (E(XY))^2 \leq E(X^2) E(Y^2)$$

$$\Rightarrow E|XY| \leq \sqrt{E(X^2) E(Y^2)}$$

**Theorem** (Jensen's Inequality)

If  $g$  is convex, then  $Eg(x) \geq g(E(x))$

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

(pf) Let  $g$  be conv.  $L(x) = a + bx$  be a line tangent to  $g(x)$  at  $E(x)$ .

Then,  $Eg(x) \geq EL(x) = E(a + bx) = a + bE(x) = L(E(x)) = g(E(x))$ .

Corollary  $\mathbb{E}(X^2) \geq (\mathbb{E} X)^2$ ,  $|\mathbb{E} X| \leq \mathbb{E}|X|$ .

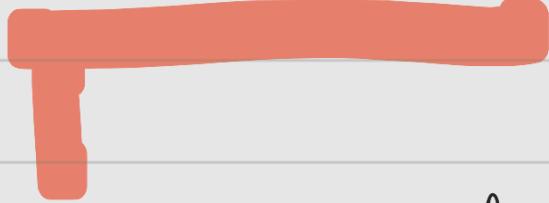
If  $X > 0$ ,  $\mathbb{E}(1/X) \geq 1/\mathbb{E}(X)$

$\mathbb{E}(\log X) \leq \log \mathbb{E}(X)$ .

Definition (Kullback - Leibler (KL) divergence)

KL divergence b/wn densities  $p, q$  is

$$D(p||q) = \int p(z) \log \left( \frac{p(z)}{q(z)} \right) dz.$$



# Convergence of Random Variables



## § 5.2 Types of Convergence

### Definition (Types of Convergence)

Let  $X_1, X_2, \dots$  be a sequence of random variables and let  $X$  denote another R.V. Let  $F_n$  denote the CDF of  $X_n$ ,  $F$  be CDF of  $X$ .

1.  $X_n$  converges to  $X$  in probability, written  $X_n \xrightarrow{P} X$ , if

$$\forall \epsilon > 0, \quad P(|X_n - X| > \epsilon) \rightarrow 0 \text{ as } n \rightarrow \infty$$

2.  $X_n$  converges to  $X$  in distribution, written  $X_n \rightsquigarrow X$ , if

$$\lim_{n \rightarrow \infty} F_n(t) = F(t)$$

at all  $t$  for which  $F$  is continuous

### Definition (Convergence in Quadratic Mean / L<sub>2</sub> / as)

- $X_n$  converges to  $X$  in quadratic mean (convergence in  $L_2$ ), written  $X_n \xrightarrow{q.m.} X$  if  $E(X_n - X)^2 \rightarrow 0$
- $X_n$  converges to  $X$  in  $L_1$ , written  $X_n \xrightarrow{L_1} X$  if  $E|X_n - X| \rightarrow 0$  as  $n \rightarrow \infty$
- $X_n$  converges almost surely to  $X$ , written  $X_n \xrightarrow{as} X$  if  $P(\{s : X_n(s) \rightarrow X(s)\}) = 1$ .

Theorem The following s hold:

(1)  $X_n \xrightarrow{d} X$  implies  $X_n \xrightarrow{P} X$

(2)  $X_n \xrightarrow{P} X$  implies  $X_n \rightsquigarrow X$

(3) If  $X_n \rightsquigarrow X$  and  $P(X=c)=1$  for some  $c$ , then  $X_n \xrightarrow{P} X$ .

$$(pf) (1) P(|X_n - X| > \varepsilon) = P(|X_n - X|^2 > \varepsilon^2) \leq \frac{1}{\varepsilon^2} E(X_n - X)^2 \rightarrow 0$$

(2) Let  $F$  be conti. at  $x$ . Fix  $x$ .

$$\begin{aligned} F_n(x) &= P(X_n \leq x) = P(X_n \leq x, X \leq x+\varepsilon) + P(X_n \leq x, X > x+\varepsilon) \\ &\leq P(X \leq x+\varepsilon) + P(|X_n - X| > \varepsilon) \\ &= F(x+\varepsilon) + P(|X_n - X| > \varepsilon) \end{aligned}$$

$$\begin{aligned} F(x-\varepsilon) &= P(X \leq x-\varepsilon) = P(X \leq x-\varepsilon, X_n \leq x) + P(X \leq x-\varepsilon, X_n > x) \\ &\leq P(X_n \leq x) + P(|X_n - X| > \varepsilon) \\ &= F_n(x) + P(|X_n - X| > \varepsilon) \end{aligned}$$

$$\text{Thus, } F(x-\varepsilon) - P(|X_n - X| > \varepsilon) \leq F_n(x) \leq F(x+\varepsilon) + P(|X_n - X| > \varepsilon).$$

$n \rightarrow \infty$  gives  $F(x-\varepsilon) \leq \liminf_n F_n(x) \leq \limsup_n F_n(x) \leq F(x+\varepsilon)$ .

By  $\varepsilon \rightarrow 0$ ,  $\lim_{n \rightarrow \infty} F_n(x) = F(x)$ .

$$(3) F(x) = \begin{cases} 0 & \text{if } x < c \\ 1 & \text{if } x \geq c \end{cases} \quad \text{Thus, } \lim_{n \rightarrow \infty} F_n(x) = \begin{cases} 0 & \text{if } x < c \\ 1 & \text{if } x > c \end{cases}$$

$$\text{Let } \varepsilon > 0. \quad P(|X_n - c| > \varepsilon) = P(X_n > c+\varepsilon) + P(X_n < c-\varepsilon)$$

$$\leq 1 - F_n(c+\varepsilon) + F_n(c-\varepsilon)$$

$$\rightarrow 1 - F(c+\varepsilon) + F(c-\varepsilon) = 1 - 1 + 0 = 0.$$

Theorem (1)  $X_n \xrightarrow{as} X$  implies  $X_n \xrightarrow{P} X$ .

(2)  $X_n \xrightarrow{fm} X$  implies  $X_n \xrightarrow{L} X$

(3)  $X_n \xrightarrow{L} X$  implies  $X_n \xrightarrow{P} X$ .

$$(\text{pf } | \text{(1)} \quad P(|X_n - X| > \varepsilon) = 1 - P(|X_n - X| < \varepsilon))$$

Let  $A_{N,\varepsilon} = \{s \mid |X_n(s) - X(s)| < \varepsilon \quad \forall n \geq N\}$ . Then,  $A_{1,\varepsilon} \subseteq A_{2,\varepsilon} \subseteq \dots$

$A_\varepsilon = \bigcup A_{N,\varepsilon} = \{s \mid \exists N \text{ s.t. } |X_n(s) - X(s)| < \varepsilon \quad \forall n \geq N\}$ .

$A = \bigcap A_\varepsilon$ .

By  $X_n \xrightarrow{as} X$ ,  $P(A) = 1$ .

$$\begin{aligned} P(|X_n - X| > \varepsilon) &= 1 - P(|X_n - X| < \varepsilon) \\ &\leq 1 - P(|X_n - X| < \varepsilon \quad \forall n \geq N) \\ &= 1 - P(A_{\varepsilon,N}) \end{aligned}$$

$$\lim_N P(A_{\varepsilon,N}) = P(A_\varepsilon) \geq P(A) = 1.$$

$\xrightarrow{\text{Theorem}}$

Thus,  $\lim_N P(|X_n - X| > \varepsilon) = 0$ .

$$(2) \quad \mathbb{E} (X_n - X)^2 \rightarrow 0.$$

$$(\mathbb{E} |X_n - X|)^2 \leq \mathbb{E} (X_n - X)^2 \rightarrow 0. \text{ Thus, } \mathbb{E} |X_n - X| \rightarrow 0.$$

$$(3) \quad P(|X_n - X| < \varepsilon) \leq \frac{1}{\varepsilon} \mathbb{E} |X_n - X| \rightarrow 0. \text{ Thus, } X_n \xrightarrow{P} X.$$

Note

almost surely  
 $L_2 \rightarrow L_1 \rightarrow \text{Probability} \xrightarrow{\quad} \text{distribution}$   
 $\nwarrow$  point mass distribution

Counter example for direction with no arrows:

## Theorem (Convergence Rules)

Let  $X_n, X, Y_n, Y$  be distribution,  $g$  be continuous.

- (1) If  $X_n \xrightarrow{P} X$  and  $Y_n \xrightarrow{P} Y$ , then  $X_n + Y_n \xrightarrow{P} X + Y$
- (2) If  $X_n \xrightarrow{\text{gm}} X$  and  $Y_n \xrightarrow{\text{gm}} Y$ , then  $X_n + Y_n \xrightarrow{\text{gm}} X + Y$
- (3) If  $X_n \rightsquigarrow X$  and  $Y_n \rightsquigarrow Y$ , then  $X_n + Y_n \rightsquigarrow X + Y$
- (4) If  $X_n \xrightarrow{P} X$  and  $Y_n \xrightarrow{P} Y$ , then  $X_n Y_n \xrightarrow{P} XY$
- (5) If  $X_n \rightsquigarrow X$  and  $Y_n \rightsquigarrow Y$ , then  $X_n Y_n \rightsquigarrow XY$
- (6) If  $X_n \xrightarrow{P} X$ , then  $g(X_n) \xrightarrow{P} g(X)$
- (7) If  $X_n \rightsquigarrow X$ , then  $g(X_n) \rightsquigarrow g(X)$ .

(pf) (1), (2), (4), follows if convergence in p, gm are

convergence in topological space induced by

$$d(X, Y) = \sup_{\varepsilon > 0} \mathbb{P}(|X - Y| > \varepsilon)$$

$$d'(X, Y) = [\mathbb{E}(X - Y)^2]^{\frac{1}{2}}$$
 respectively.

↪ Have to check it is norm.

## §5.3 The Law of Large Numbers

[Theorem] (The Weak Law of Large Numbers)

If  $X_1, \dots, X_n$  are IID, then  $\bar{X}_n \xrightarrow{P} \mu$ . ( $\mu = E(X_i)$ ).

(pf) Assume  $0 < \alpha$ . (Just to simplify proof)

By Chebyshev,  $P(|\bar{X}_n - \mu| > \varepsilon) \leq \frac{V(\bar{X}_n)}{\varepsilon^2} = \frac{\sigma^2}{n\varepsilon^2} \rightarrow 0$  as  $n \rightarrow \infty$ .

[Theorem] (The Strong Law of Large Numbers)

Let  $X_1, \dots, X_n$  be IID. If  $\mu = E|X_i| < \infty$ , then  $\bar{X}_n \xrightarrow{as} \mu$

(pf)