# Chapter 5. Joint Probability Distributions and Random Samples

Math 3670 Spring 2025

Georgia Institute of Technology

# Section 1.
# Jointly Distributed Random Variables

# Two Discrete Random Variables

## Definition

Let $X$ and $Y$ be two discrete RVs defined on the sample space $\mathcal{S}$ of an experiment.
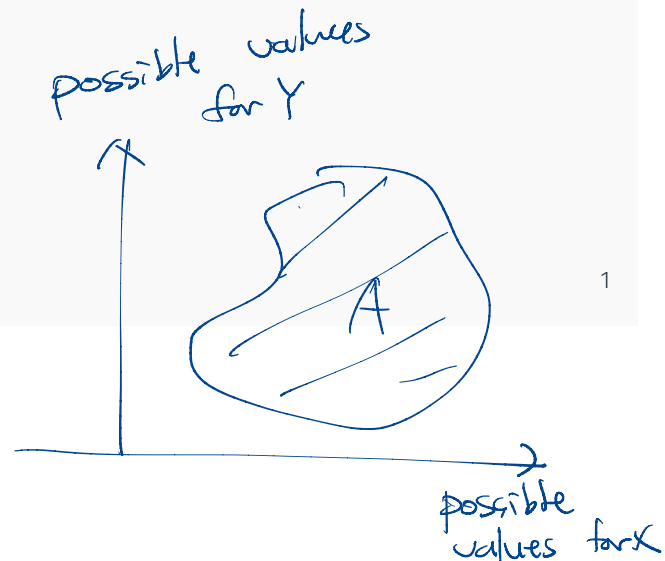
The joint probability mass function $p(x, y)$ is defined by

$$p(x, y) = \mathbb{P}(X = x, Y = y)$$

$$\hookleftarrow \text{AND}$$

The joint PMF satisfies

1. $p(x, y) \geq 0$
2. $\sum_{x,y} p(x, y) = 1$
3. $\mathbb{P}((X, Y) \in A) = \sum_{(x,y) \in A} p(x, y)$

possible values for $Y$

possible values for $X$

$A$

# Two Discrete Random Variables

## Example

A large insurance agency services a number of customers who have purchased both a homeowner's policy and an automobile policy from the agency. For an automobile policy, the choices are $100 and $250, whereas for a homeowner's policy, the choices are 0, $100, and $200.

Suppose an individual with both types of policy is selected at random from the agency's files. Let $X$ be deductible amount on the auto policy and $Y$ deductible amount on the homeowner's policy.

| X \ Y | 0 | 100 | 200 |
|-------|------|------|-----|
| 100   | 0.2  | 0.1  | 0.2 |
| 250   | 0.05 | 0.15 | 0.3 |

$$p(100, 100) = P(X = 100, Y = 100) = 0.1$$

$$P(250, 0) = P(X = 250, Y = 0) = 0.05$$

$$P_X(100) = P(X = 100) = P(X = 100, Y = 0)$$
$$+ P(X = 100, Y = 100)$$
$$+ P(X = 100, Y = 200)$$
$$= 0.2 + 0.1 + 0.2 = 0.5$$

## Definition

For a given joint PMF $p(x, y)$ of random variables $X$ and $Y$, the ==marginal probability mass function of== $X$ is given by

Marginal PMFs

$$p_X(x) := \mathbb{P}(X = x) = \sum_{\text{all } y's} p(x, y)$$

$$p_Y(y) = \mathbb{P}(Y = y) = \sum_{\text{all } x's} p(x, y)$$

Knowing Joint PMF $\implies$ Can get $P_X$, $P_Y$

$\Longleftarrow$ in general

# Two Discrete Random Variables

## Example

A large insurance agency services a number of customers who have purchased both a homeowner's policy and an automobile policy from the agency. For an automobile policy, the choices are $100 and $250, whereas for a homeowner's policy, the choices are 0, $100, and $200.

Suppose an individual with both types of policy is selected at random from the agency's files. Let $X$ be deductible amount on the auto policy and $Y$ deductible amount on the homeowner's policy.

|     | 0    | 100  | 200 |
| --- | ---- | ---- | --- |
| 100 | 0.2  | 0.1  | 0.2 |
| 250 | 0.05 | 0.15 | 0.3 |

## Definition

Let *X* and *Y* be two continuous RVs.

The joint probability density function $f(x, y)$ is defined by

$$f(x, y) \neq \qquad \mathbb{P}(X = x, Y = y)$$

The joint PDF satisfies

1. $f(x, y) \geq 0$
2. $\int \int f(x, y)$ $\qquad \int_{\mathbb{R}} \int_{\mathbb{R}} f(x, y)\, dx\, dy = 1$
3. $\mathbb{P}((X, Y) \in A) =$

$$\mathbb{P}((X, Y) \in A) = \int \int_A f(x, y)\, dx\, dy$$

5

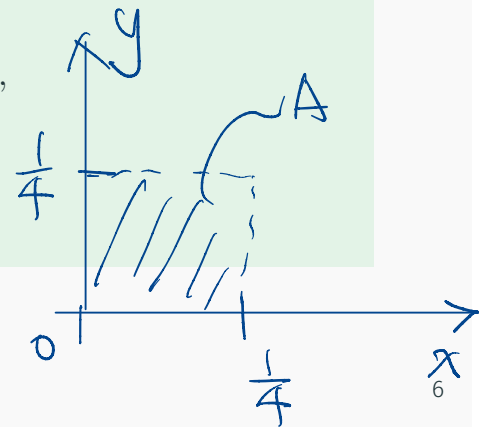## Two Continuous Random Variables

### Example

A bank operates both a drive-up facility and a walk-up window.

On a randomly selected day, let $X$ be the proportion of time that the drive-up facility is in use (at least one customer is being served or waiting to be served) and $Y$ the proportion of time that the walk-up window is in use.

The joint PDF is given by

$$f(x,y) = \begin{cases} \frac{6}{5}(x+y^2), & 0 \le x \le 1, 0 \le y \le 1, \\ 0 & \text{otherwise.} \end{cases}$$

Find $\mathbb{P}(0 \le X \le \frac{1}{4}, 0 \le Y \le \frac{1}{4})$.

$$= \mathbb{P}((X,Y) \in A)$$

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

$$= \int_0^{\frac{1}{4}}\int_0^{\frac{1}{4}} \frac{6}{5}(x+y^2)\, dx\, dy = \int_0^{\frac{1}{4}} \frac{6}{5}\left[\frac{1}{2}x^2 + y^2 \cdot x\right]_0^{\frac{1}{4}} dy$$

$$= \int_0^{\frac{1}{4}} \frac{6}{5}\left(\frac{1}{32} + \frac{y^2}{4}\right) dy = \frac{6}{5}\left[\frac{1}{32}\cdot y + \frac{y^3}{12}\right]_0^{\frac{1}{4}}$$
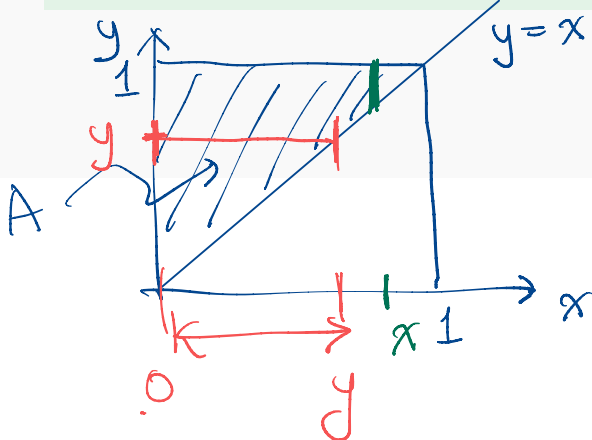
$$= \quad -- \cdots$$

A bank operates both a drive-up facility and a walk-up window.

On a randomly selected day, let $X$ be the proportion of time that the drive-up facility is in use (at least one customer is being served or waiting to be served) and $Y$ the proportion of time that the walk-up window is in use.

The joint PDF is given by

$$f(x, y) = \begin{cases} \frac{6}{5}(x + y^2), & 0 \leq x \leq 1, 0 \leq y \leq 1, \\ 0 & \text{otherwise.} \end{cases}$$

Find ~~$P(0 \leq X \leq \frac{1}{4}, 0 \leq Y \leq \frac{1}{4})$~~

$$\mathbb{P}(\ X \leq Y)$$

$$= \mathbb{P}(\ (X, Y) \in A)$$

$$= \int_0^1 \int_0^y \frac{6}{5}(x + y^2)\, dx\, dy$$

$$= \int_0^1 \int_x^1 \qquad dy\, dx$$



6

Joint PMF : $p(x,y) = P(X=x, Y=y)$

Joint PDF : $f(x,y)$ such that

(i) $f(x,y) \geq 0$

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

(iii) $P((X,Y) \in A) = \iint_A f(x,y) \, dx \, dy$.

## Two Continuous Random Variables

**Definition**

For a given joint PDF $f(x,y)$ of random variables $X$ and $Y$, the marginal probability density function of $X$ is given by

$$f_X(x) := \int_{-\infty}^{\infty} f(x,y) \, dy$$

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

### Example

A bank operates both a drive-up facility and a walk-up window.

On a randomly selected day, let $X$ be the proportion of time that the drive-up facility is in use (at least one customer is being served or waiting to be served) and $Y$ the proportion of time that the walk-up window is in use.

The joint PDF is given by

$$f(x,y) = \begin{cases} \frac{6}{5}(x+y^2), & 0 \le x, y \le 1 \\ 0 & \text{otherwise.} \end{cases}$$

$$0 \le x \le 1$$
$$0 \le y \le 1$$

Find the marginal PDFs.

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

$$= \int_{0}^{1} \frac{6}{5}(x+y^2)\, dy$$

$$= \frac{6}{5}\left[ xy + \frac{1}{3}y^3 \right]_0^1 = \frac{6}{5}\left(x+\frac{1}{3}\right)$$

$$f_X(x) = \begin{cases} \frac{6}{5}\left(x+\frac{1}{3}\right), & 0 \le x \le 1 \\ 0 & \text{otherwise} \end{cases}$$

### Example

The joint PDF is given by

$$f(x,y) = \begin{cases} 24xy, & 0 \le x \le 1, 0 \le y \le 1, x+y \le 1, \\ 0 & \text{otherwise.} \end{cases}$$

Find the marginal PDFs and $\mathbb{P}(X+Y \le 1/2)$.

$f_X \qquad f_Y$

Inequalities

region

look at __equalities__

gives bdry

$y=1$

$y=0$

$x$

$x=0 \qquad x=1$

9

fixed $0 \le x \le 1$

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

$$= \int_{0}^{1-x} 24xy \; dy$$

$$= \left[ 12 xy^2 \right]_{0}^{1-x}$$

$$= 12 x (1-x)^2 .$$

$$f_Y(y) = 12 y (1-y)^2$$

# Two Continuous Random Variables

### Example

The joint PDF is given by

$$f(x,y) = \begin{cases} 24xy, & 0 \leq x \leq 1, 0 \leq y \leq 1, x+y \leq 1, \\ 0 & \text{otherwise.} \end{cases}$$
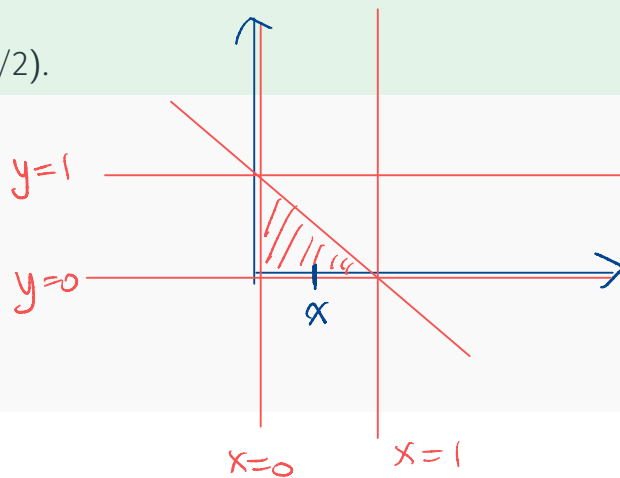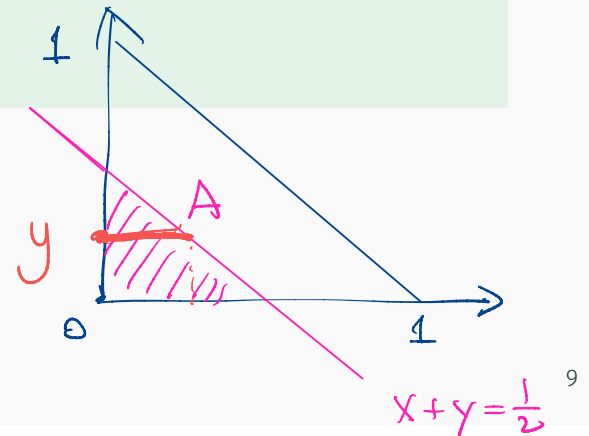
Find the marginal PDFs and $\mathbb{P}(X+Y \leq 1/2)$.

$$P(X+Y \leq \tfrac{1}{2})$$

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

$$= \iint_A f(x,y)\, dx\, dy$$

$$= \int_0^{\frac{1}{2}} \int_0^{\frac{1}{2}-y} 24xy\, dx\, dy$$



$1$

$y$

$A$

$0$

$1$

$x+y=\tfrac{1}{2}$

Recall    Two events   A, B  are  Indep.

iff          $P(A \wedge B) = P(A) \cdot P(B)$.

## Independent Random Variables

### Definition

Two random variables $X$ and $Y$ are said to be independent if

for  any  intervals   $I$, $J$  in  $\mathbb{R}$,

$\{X \in I\}$,   $\{Y \in J\}$   indep.

---

① If $X, Y$ discrete,

$X, Y$  indep  $\Longleftrightarrow$  $p(x, y) = P_X(x) \cdot P_Y(y)$

② If $X, Y$ continuous,

$X, Y$  indep  $\Longleftrightarrow$  $f(x, y) = f_X(x) \cdot f_Y(y)$

**Example**

The joint PDF is given by

$$f(x, y) = \begin{cases} 24xy, & 0 \leq x \leq 1, 0 \leq y \leq 1, x + y \leq 1 \\ 0 & \text{otherwise.} \end{cases}$$

Are *X* and *Y* independent?

$$f_X(x) = 12 \, x \, (1-x)^2$$
$$f_Y(y) = 12 \, y \, (1-y)^2$$

previous example.

11

$$24xy = f(x, y) \neq f_X(x) \cdot f_Y(y)$$

$$= 12 \, x \, (1-x)^2 \cdot 12 \cdot y \, (1-y)^2.$$

$$\Rightarrow X, Y \quad \text{Not} \quad \text{Indep.}$$

**Example**

Suppose that the lifetimes of two components are independent of one another and that the first lifetime, $X_1$, has an exponential distribution with parameter $\lambda_1$, whereas the second, $X_2$, has an exponential distribution with parameter $\lambda_2$.

Find the joint PDF.

$$f(x_1, x_2) = f_{X_1}(x_1) \cdot f_{X_2}(x_2)$$

$$= \begin{cases} \lambda_1 e^{-\lambda_1 x_1} \cdot \lambda_2 e^{-\lambda_2 x_2} & , \; x_1, x_2 \geq 0 \\ 0 & , \; o.w. \end{cases}$$

12

# Independent Random Variables

## Definition

The random variables $X_1, X_2, \cdots, X_n$ are said to be **independent** if

$$\{ X_1 \in I_1 \}, \{ X_2 \in I_2 \}, \cdots, \{ X_n \in I_n \} \quad \text{Indep.}$$

$$\text{for any choice of } I_1, I_2, \cdots, I_n.$$

## Conditional Distributions

### Definition

Let $X$ and $Y$ be two continuous RVs with joint PDF $f(x, y)$.

Then for any $x$ for which $f_X(x) > 0$, the conditional probability density function of $Y$ given that $X = x$ is

$$A \quad New \quad RV: \qquad Y \mid X = x \qquad with \quad \overset{Condi.}{PDF}$$

$$f_{Y|X}(y|x) = \frac{f(x,y)}{f_X(x)}$$

**Example**

A bank operates both a drive-up facility and a walk-up window.

On a randomly selected day, let $X$ be the proportion of time that the drive-up facility is in use (at least one customer is being served or waiting to be served) and $Y$ the proportion of time that the walk-up window is in use.

The joint PDF is given by

$$f(x, y) = \begin{cases} \frac{6}{5}(x + y^2), & 0 \leq x \leq 1, 0 \leq y \leq 1, \\ 0 & \text{otherwise.} \end{cases}$$

Find the conditional PDF of $Y$ given $X = 0.8$.

Compute $\mathbb{P}(Y \leq 0.5 | X = 0.8)$.

$$f_{Y|X}(y \mid 0.8) = \frac{f(0.8, y)}{f_X(0.8)} = \boxed{\frac{\frac{6}{5}(0.8 + y^2)}{\frac{6}{5}(0.8 + \frac{1}{3})}}$$

$$0 \leq y \leq 1.$$

$$f_X(x) = \frac{6}{5}\left(x + \frac{1}{3}\right)$$

$$f_X(0.8) = \frac{6}{5}\left(\frac{4}{5} + \frac{1}{3}\right) = \frac{104}{75}$$

$$\mathbb{P}(Y \leq 0.5 \mid X = 0.8) =$$

$$= \mathbb{P}(\,(Y | X = 0.8) \leq 0.5\,)$$

$$\frac{\mathbb{P}(Y \leq 0.5, X = 0.8)}{\mathbb{P}(X = 0.8)} = 0$$

15

$$= \int_0^{0.5} f_{Y|X}(y \mid 0.8) \, dy$$

$$= \int_0^{0.5} \frac{(0.8 + y^2)}{(0.8 + \frac{1}{3})} \, dy$$

$$E[Y \mid X = 0.8] = \int_0^1 y \cdot f_{Y|X}(y \mid 0.8) \, dy$$

## Exercise

(5.1-12) Two components of a minicomputer have the following joint PDF for their useful lifetimes $X$ and $Y$:

$$f(x, y) = \begin{cases} xe^{-x(y+1)}, & x \geq 0, y \geq 0, \\ 0 & \text{otherwise.} \end{cases}$$

1. What is the probability that the lifetime $X$ of the first component exceeds 3?
2. What are the marginal PDFs of $X$ and $Y$? Are the two lifetimes independent? Explain.
3. What is the probability that the lifetime of at least one component exceeds 3?

Section 2.
Expected Values, Covariance, and Correlation

**Proposition**

Let $X$ and $Y$ be jointly distributed RVs with PMF $p(x, y)$ or PDF $f(x, y)$ according to whether the variables are discrete or continuous.

Let $h(x, y)$ be a function of two variables, then we can define a new random variable $Z = h(X, Y)$.

The expectation of $Z$ is

$$\mathbb{E}[Z] = \mathbb{E}[h(X, Y)] = \begin{cases} \sum_x \sum_y h(x,y)\, p(x,y) \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} h(x,y)\, f(x,y)\, dx\, dy \end{cases}$$

$$\underline{\text{Example}} \quad h(x,y) = x$$
$$h(x,y) = y$$
$$h(x,y) = x \cdot y$$

17

$$\mathbb{E}[X] = \iint x \cdot f(x,y)\, dx\, dy = \int x \left( \int f(x,y)\, dy \right) dx$$

$$= \int x \cdot f_X(x)\, dx$$

$$\overset{\shortparallel}{f_X(x)}$$

$$\mathbb{E}[X \cdot Y] = \iint x \cdot y\, f(x,y)\, dx\, dy$$

# Expectation of a Function of Two Random Variables

**Example**

Five friends have purchased tickets to a certain concert.

If the tickets are for seats 1–5 in a particular row and the tickets are randomly distributed among the five, what is the expected number of seats separating any particular two of the five?

Let $X$ and $Y$ denote the seat numbers of the first and second individuals, respectively.
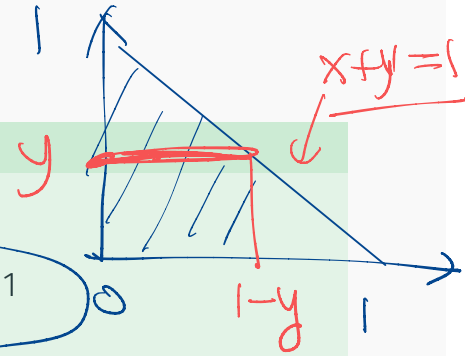
**Example**

The joint PDF is given by

$$f(x,y) = \begin{cases} 24xy, & 0 \le x \le 1, 0 \le y \le 1, x+y \le 1 \\ 0 & \text{otherwise.} \end{cases}$$

Find $\mathbb{E}[XY]$.

$$\mathbb{E}[XY] = \iint xy \cdot f(x,y) \, dx\, dy$$

$$= \int_0^1 \int_0^{1-y} 24 \, x^2 y^2 \, dx\, dy$$

$$= \int_0^1 \left[ 8 x^3 y^2 \right]_0^{1-y} dy$$

$$= \int_0^1 8(1-y)^3 \cdot y^2 \, dy$$

$$= \int_0^1 8(1-t)^2 t^3 \, dt \quad (1-y=t)$$

$x+y=1$

$1-y$

19

# Covariance

$$\mu_X = \mathbb{E}[X] \quad, \quad \mu_Y = \mathbb{E}[Y]$$

**Definition**

The covariance between two RVs $X$ and $Y$ is

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

$$= \iint (x - \mu_X)(y - \mu_Y)\, f(x, y)\, dx\, dy$$

$$= \sum_y \sum_x (x - \mu_X)(y - \mu_Y)\, p(x, y)$$

## Covariance

The joint PDF is given by

$$f(x,y) = \begin{cases} 24xy, & 0 \leq x \leq 1, 0 \leq y \leq 1, x+y \leq 1 \\ 0 & \text{otherwise.} \end{cases}$$

previous example.

Find the covariance of $X, Y$.

$$\text{Cov}(X,Y) = \underline{E[X \cdot Y]} - E[X] \cdot E[Y].$$

$$E[X] = \iint x \cdot f(x,y) \, dx \, dy$$

$$= \int_0^1 x \cdot \underbrace{12 x (1-x)^2}_{f_X(x)} \, dx$$

$$= 12 \int_0^1 x^2 (1 - 2x + x^2) \, dx$$

$$= 12 \cdot \left[ \frac{1}{3} - 2 \cdot \frac{1}{4} + \frac{1}{5} \right] = \frac{12}{30} = \frac{2}{5}$$

$$= E[Y] \quad \text{(by symmetry)}$$

Recall $\quad \text{Var}(X) = E\left[(X - \mu_X)^2\right] = \text{Cov}(X, X)$

$$= E[X^2] - (E[X])^2$$

### Proposition

$$\text{Cov}(X, Y) = E[XY] - E[X] \cdot E[Y]$$

4/8/2025

Recall    $X, Y$    joint PDF (or PMF)

- $\text{Cov}(X, Y) = E[(X - \mu_x) \cdot (Y - \mu_Y)] = \iint (x-\mu_x)(y-\mu_Y) f(x,y) \, dx \, dy$

where    $\mu_x = E[X]$,    $\mu_Y = E[Y]$

- $\text{Cov}(X, Y) = E[X \cdot Y] - E[X] \cdot E[Y]$.

- $\text{Cov}(X, Y) = \text{Cov}(Y, X)$,    $\text{Cov}(X, X) = \text{Var}(X)$

## Correlation Coefficient

**Definition**

The correlation coefficient of $X$ and $Y$ is defined by

$$\text{Corr}(X, Y) = \rho_{X,Y} = \rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)} \sqrt{\text{Var}(Y)}}$$

(rho)

- $\text{Corr}(X, Y) = \text{Corr}(Y, X)$

- $\text{Corr}(X, X) = 1$.

- $-1 \leq \text{Corr}(X, Y) \leq 1$.

23

In general,    $(E[X \cdot Y])^2 \leq E[X^2] \cdot E[Y^2]$

$\Rightarrow$    $\text{Cov}(X, Y)^2 \leq \text{Var}(X) \cdot \text{Var}(Y)$

$\Rightarrow$    $\text{Corr}(X, Y)^2 \leq 1$

$$\text{Var}(aX+b) = a^2 \text{Var}(X)$$
$$\text{Var}(cY+d) = c^2 \text{Var}(Y)$$

$$\text{Cov}(aX+b, cY+d) = a \cdot c \, \text{Cov}(X, Y)$$

## Correlation Coefficient

### Properties

1. For constants $a, b, c, d$,
$$\text{Corr}(aX+b, cY+d) = \left(\frac{ac}{\sqrt{a^2c^2}}\right) \text{Corr}(X, Y)$$

2. $-1 \leq Corr(X, Y) \leq 1$      $(\text{Cov}(X,Y)=0)$

3. If $X$ and $Y$ are independent, then $\text{Corr}(X, Y) = 0$. The converse does not hold in general.
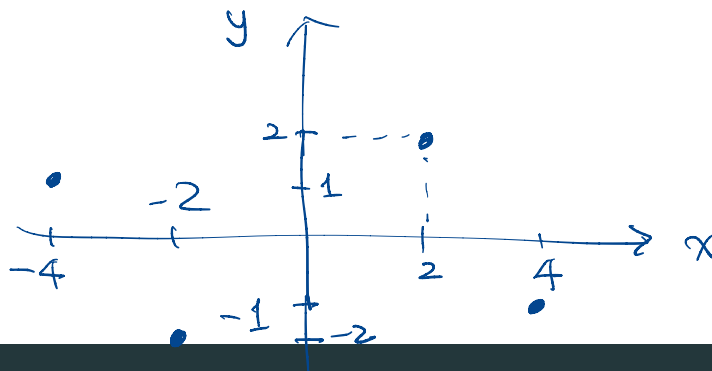
4. If $\text{Corr}(X, Y) = 1, -1$, then $Y = aX + b$ for some $a, b$.

- $X, Y$    indep     $\Rightarrow$     $\text{Cov}(X, Y) = 0$

                                  $\Rightarrow$     $\text{Corr}(X, Y) = 0.$

- $\text{Corr}(aX+b, cY+d)$

$$= \begin{cases} \text{Corr}(X, Y) & , \text{ when } \quad ac \geq 0 \\ -\text{Corr}(X, Y) & , \text{ when } \quad a \cdot c < 0 \end{cases}$$

## Correlation Coefficient

### Example

The joint PMF is given by

$$p(x, y) = \begin{cases} \frac{1}{4}, & (x, y) = (-4, 1), (4, -1), (2, 2), (-2, -2) \\ 0 & \text{otherwise.} \end{cases}$$

Find the covariance and the correlation coefficient.

$$X = \quad 4 \quad , 2 , \quad -2, \quad -4 \qquad \text{equally likely}$$

$$Y = \quad 2 \quad 1 \quad -1 \quad -2 \qquad\qquad \text{"}$$

$$E[X] = 0 \quad = \quad E[Y]$$

$$\text{Cov}(X, Y) \; = \; E[X \cdot Y] \; = \; \frac{1}{4} \left( (-4) \cdot 1 + 4 (-1) + 2 \cdot 2 + \right.$$
$$\left. (-2) \cdot (-2) \right)$$

$$= \quad 0 \quad = \quad \text{Corr}(X, Y)$$

## Exercise

(5.2-24) Six individuals, including A and B, take seats around a circular table in a completely random fashion.

Suppose the seats are numbered $1, 2, \cdots, 6$.

Let $X$ be A's seat number and $Y$ B's seat number.

If A sends a written message around the table to B in the direction in which they are closest,

how many individuals (including A and B) would you expect to handle the message?

$\underline{X, Y}$    with   joint   PDF $\sim$ PMF     $f_{Y|X}(y|x)$

a New RV     $Y | \underset{\smile}{X=x}$    with $\begin{cases} \text{PDF} & \dfrac{f(x,y)}{f_X(x)} \\[2em] \text{PMF} & \dfrac{p(x,y)}{p_X(x)} \end{cases}$

$= P_{Y|X}(y|x)$

$\mathbb{E}[h(Y) | X=x] = \int h(y) \cdot f_{Y|X}(y|x)\, dy$

---

## Conditional distribution

> **Definition**
>
> **The conditional expectation** of $Y$ given $X = x$ is defined by
> $$\mathbb{E}[Y|X = x] = \sum_y y f_{Y|X}(y|x).$$
>
> The conditional variance of $Y$ given $X = x$ is defined by
> $$\begin{aligned} \mathrm{Var}(Y|X = x) &= \mathbb{E}[(Y - \mathbb{E}[Y|X = x])^2 | X = x] \\ &= \mathbb{E}[Y^2|X = x] - (\mathbb{E}[Y|X = x])^2. \end{aligned}$$

$\mathbb{E}[Y | \underset{\smile}{X = x}] \Leftarrow$   a $\underline{number}$

- For each choice of $x$, $\mathbb{E}[Y|X=x]$ might give a diff #.

- $h(x) = \mathbb{E}[Y|X=x]$ a function of $x$

One can consider $\mathbb{E}[Y|X = x]$ as a function of $x$.

Say $h(x) = \mathbb{E}[Y|X = x]$

$x:$ real #

We define a random variable $\mathbb{E}[Y|X] = h(X)$.

$\underset{\text{RV}}{X}$

New RV $\quad h(X) = \mathbb{E}[Y \cdot | X]$

$X, Y$

$\Rightarrow$ New RV $\qquad Y \mid X = x \qquad$ for each $x$

$\Rightarrow h(x) = \mathbb{E}[Y \mid X = x]$

$\Rightarrow$ From $h$, define a new RV $h(X)$

$\qquad h(X) = \mathbb{E}[Y \mid X]$

## Contional expectation as a function and a random variable

$$\mathbb{E}[h(X)] = \mathbb{E}[Y] = \mathbb{E}(\mathbb{E}[Y \mid X])$$

**Theorem**

1. $\mathbb{E}[\mathbb{E}[Y|X]] = \mathbb{E}[Y]$
2. $\text{Var}(Y) = \mathbb{E}[\text{Var}(Y|X)] + \text{Var}(\mathbb{E}[Y|X])$

## Exercise

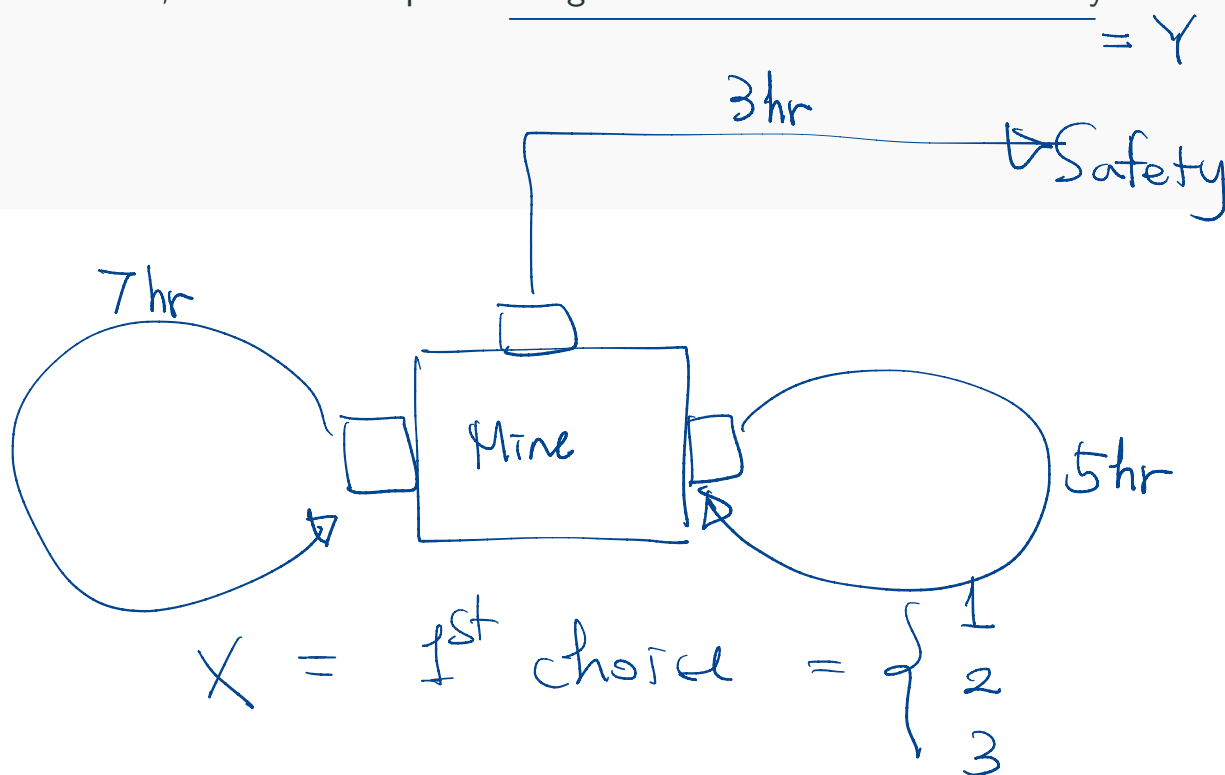$$\mathbb{E}[Y] = \mathbb{E}[\mathbb{E}[Y \mid X]]$$

A miner is trapped in a mine containing 3 doors.

The first door leads to a tunnel that will take him to safety after 3 hours of travel.

The second door leads to a tunnel that will return him to the mine after 5 hours of travel.

The third door leads to a tunnel that will return him to the mine after 7 hours.

If we assume that the miner is at all times equally likely to choose any one of the doors, what is the expected length of time until he reaches safety?

$= Y$

3 hr

↪ Safety

7 hr

Mine

5 hr

$$X = 1^{st} \text{ choice} = \begin{cases} 1 \\ 2 \\ 3 \end{cases}$$

If   X = 1   $E[Y|X=1] = \underline{3}$

   X = 2   $E[Y|X=2] = \underline{E[Y] \quad +5}$

   X = 3   $E[Y|X=3] = \overline{\underline{E[Y] + 7}}$

$$E(E[Y|X]) = \frac{1}{3} \cdot 3 + \frac{1}{3}(E[Y]+5)$$

$$+ \frac{1}{3}(E[Y]+7)$$

$$= E[Y]$$

$$\frac{1}{3} E[Y] = \frac{1}{3}(3+5+7)$$

$$E[Y] = 15.$$

# Section 4.
# The Distribution of the Sample Mean

joint PDF

$$f(x_1, x_2, \cdots, x_n)$$
$$= f_{x_1}(x_1) \cdot f_{x_2}(x_2) \cdots f_{x_n}(x_n)$$

equal functions

## Sample Mean

### Definition

The RVs $X_1, X_2, \cdots, X_n$ are said to form a (simple) random sample of size $n$ if

1. they are independent RVs, and
2. every $X_i$ has the same probability distribution.

] indep. identically distributed

$=$ i.i.d.

The sample mean is defined by

$$\bar{X} = \frac{X_1 + X_2 + \cdots + X_n}{n}$$

The sample total is defined by

$$T = X_1 + X_2 + \cdots + X_n$$

Simple case $\quad n = 2.$

$$T = X_1 + X_2$$
$$E[T] = E[X_1 + X_2] = E[X_1] + E[X_2]$$
$$Var(T) = E[(X_1+X_2)^2] - (E[X_1+X_2])^2$$
$$\underset{\uparrow}{=} Var(X_1) + Var(X_2)$$

because $X_1, X_2$ indep.

If not indep.    $\quad Var(T) = Var(X_1) + Var(X_2) + 2Cov(X_1, X_2)$

$$Var(X_1) = \cdots = Var(X_n) = \sigma^2$$

## Sample Mean

$$\mathbb{E}[X_1] = \cdots = \mathbb{E}[X_n] = \mu$$

### Proposition

i.i.d.

Let $X_1, \cdots, X_n$ be a random sample from a distribution with mean $\mu$ and variance $\sigma^2$. Then,

$$\mathbb{E}[\bar{X}] = \frac{1}{n} \cdot n \cdot \mu = \mu$$
$$Var(\bar{X}) = \frac{1}{n^2} \cdot n\sigma^2 = \sigma^2/n$$
$$\mathbb{E}[T] = \mathbb{E}[X_1] + \cdots + \mathbb{E}[X_n] = n \cdot \mu$$
$$Var(T) = Var(X_1) + \cdots + Var(X_n) = n \cdot \sigma^2$$

$$T = X_1 + X_2 + \cdots + X_n$$

$$\bar{X} = \frac{X_1 + \cdots + X_n}{n} = \frac{T}{n}$$

$$\mathbb{E}[\bar{X}] = \mathbb{E}\left[\frac{T}{n}\right] = \mathbb{E}\left[\left(\frac{1}{n}\right) \cdot T\right] = \frac{1}{n}\mathbb{E}[T]$$

$$Var(\bar{X}) = Var\left(\left(\frac{1}{n}\right)T\right) = \frac{1}{n^2} Var(T)$$

> **Example**
>
> In a notched tensile fatigue test on a titanium specimen, the expected number of cycles to first acoustic emission (used to indicate crack initiation) is $\mu = 28,000$, and the standard deviation of the number of cycles is $\sigma = 5000$.
>
> Let $X_1, X_2, \cdots, X_{25}$ be a random sample of size 25, where each $X_i$ is the number of cycles on a different randomly selected specimen.

$$\overline{X} = \frac{X_1 + \cdots + X_{25}}{25}$$

$$E[\overline{X}] = \mu = 28,000$$

$$\text{Var}(\overline{X}) = \frac{\sigma^2}{n} = \frac{(5000)^2}{25} = 1000,000$$

$$\sqrt{\text{Var}(\overline{X})} = \sqrt{1000000} = 1000$$

$$X_1, \cdots, X_n \qquad \text{Indep.} \qquad \underline{\text{Same distribution}}$$

# The Case of a Normal Population Distribution

**Proposition**

$$\text{indep.} \qquad X_1, \cdots, X_n \sim N(\mu, \sigma^2)$$

Let $X_1, \cdots, X_n$ be a random sample from **a normal distribution** with mean $\mu$ and variance $\sigma^2$. Then,

$$\boxed{\overline{X}} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$T \sim N(n\mu, n\sigma^2)$$

$$E[\overline{X}] = \mu \quad , \quad Var(\overline{X}) = \frac{\sigma^2}{n}$$

$$X_1, \cdots, X_n \sim N(\mu, \sigma^2) \implies \bar{X} \sim N(\mu, \frac{\sigma^2}{n})$$
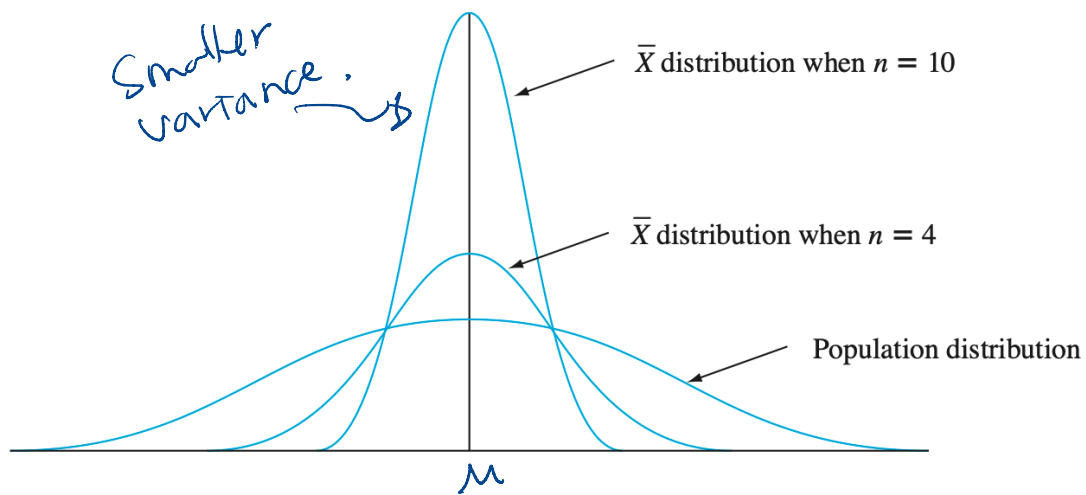
## The Case of a Normal Population Distribution

Smaller variance.

$\bar{X}$ distribution when $n = 10$

$\bar{X}$ distribution when $n = 4$

Population distribution

$\mu$

**Figure 5.14**   A normal population distribution and $\bar{X}$ sampling distributions

31

**Example**

_Indep._

The time that it takes a randomly selected rat of a certain subspecies to find its way through a maze is a normally distributed RV with $\mu = 1.5$ min and $\sigma = .35$ min.

Suppose five rats are selected.

Let $X_1, \ldots, X_5$ denote their times in the maze.

$X_1, X_2, \cdots, X_5 \sim N(\mu, \sigma^2)$

Assuming the $X_i$'s to be a random sample from this normal distribution,

what is the probability that the total time is between 6 and 8 min?

$$T = X_1 + \cdots + X_5 \sim N(5 \cdot (1.5), \, 5 \cdot (0.35)^2)$$

$$N(7.5, \, 0.6125)$$

$$P(6 < T < 8)$$

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

$$= P\left( \frac{6 - 7.5}{\sqrt{0.6125}} < \frac{T - 7.5}{\sqrt{0.6125}} < \frac{8 - 7.5}{\sqrt{0.6125}} \right)$$

$$= \Phi(\ ) - \Phi(\ ) = \text{use table !!}$$

# The Central Limit Theorem

(CLT)

## Theorem

i.i.d

finite    finite

Let $X_1, X_2, \cdots, X_n$ be a random sample from a distribution with mean $\mu$ and variance $\sigma^2$.

If $n$ is sufficiently large, $\bar{X}$ and $T$ have approximately normal distributions.

Rule of Thumb: If $n \geq 30$, the Central Limit Theorem can be used.

$$\bar{X} \approx N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$\frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}} \approx N(0, 1)$$

> **Example**
>
> A certain consumer organization customarily reports the number of major defects for each new automobile that it tests.
>
> Suppose the number of such defects for a certain model is a random variable with mean value 3.2 and standard deviation 2.4.
>
> Among 100 randomly selected cars of this model, how likely is it that the sample average number of major defects exceeds 4?

$$X_1, X_2, \cdots, X_{100} : \quad i.i.d.$$

$$\bar{X} \approx N\left(3.2, \frac{(2.4)^2}{100}\right)$$

$$\frac{\bar{X} - 3.2}{0.24} \approx N(0, 1)$$

$$P(\bar{X} > 4) = P\left(\frac{\bar{X} - 3.2}{0.24} > \frac{4 - 3.2}{0.24}\right)$$

$$\text{by CLT} \quad \approx \quad P(Z > 3.33)$$

$$= 1 - \Phi(3.33)$$

$X_1, X_2, \cdots, X_n$ : random sample of size $n$

(i.i.d.)

(i) have the same distribution

$$\left( \begin{array}{l} \mathbb{E}[X_1] = \mathbb{E}[X_2] = \cdots = \mathbb{E}[X_n] = \mu \\ \text{Var}(X_1) = \text{Var}(X_2) = \cdots = \text{Var}(X_n) = \sigma^2 \end{array} \right)$$

(ii) Independent.

$$\overline{X} = \frac{X_1 + \cdots + X_n}{n} = \frac{T}{n}$$

$$T = X_1 + \cdots + X_n$$

$$\left[ \begin{array}{ll} \mathbb{E}[\overline{X}] = \mu , & \text{Var}(\overline{X}) = \frac{\sigma^2}{n} \\ \mathbb{E}[T] = n\mu , & \text{Var}(T) = n \cdot \sigma^2 \end{array} \right] \quad \text{Why?}$$

In general, for $X_1, X_2$ general (not neccessarily i.i.d.)

$$\mathbb{E}[aX_1 + bX_2] = a\,\mathbb{E}[X_1] + b\,\mathbb{E}[X_2]$$

$$\text{Var}(X_1 + X_2) = \mathbb{E}[(X_1 + X_2)^2] - (\mathbb{E}[X_1 + X_2])^2$$

$$= \mathbb{E}[X_1^2 + 2X_1 \cdot X_2 + X_2^2]$$

$$\qquad - ((\mathbb{E}[X_1])^2 + 2\,\mathbb{E}[X_1]\mathbb{E}[X_2] + (\mathbb{E}[X_2])^2)$$

$$= (\mathbb{E}[X_1^2] - (\mathbb{E}[X_1])^2) + 2(\mathbb{E}[X_1 \cdot X_2] - \mathbb{E}[X_1]\mathbb{E}[X_2])$$

$$\qquad\qquad + (\mathbb{E}[X_2^2] - (\mathbb{E}[X_2])^2)$$

$$= \text{Var}(X_1) + 2\,\text{Cov}(X_1, X_2) + \text{Var}(X_2)$$

# ⟨Central Limit Theorem⟩

Suppose $\quad X_1, X_2, \cdots, X_n \quad$ i.i.d.

Assume $\quad \mu = \mathbb{E}[X_1] < \infty \quad , \quad \sigma^2 = \text{Var}(X_1) < \infty$

Then,

$\underline{\overline{X}} \quad$ is approximately normal

$$\boxed{\overline{X}} \approx N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$\frac{\overline{X} - \mu}{\sqrt{\sigma^2/n}} \approx N(0,1) \qquad n \to \infty$$

? 

"approximately" or "converges to"

$$\mathbb{P}\left( \frac{\overline{X} - \mu}{\sqrt{\sigma^2/n}} \leq x \right) \xrightarrow[\text{as } n \to \infty]{} \Phi(x)$$

## Example

$$X \sim \text{Bin}(n, p)$$

$$X = X_1 + X_2 + \cdots + X_n$$

$$X_1, \cdots, X_n \sim \text{Ber}(p) \quad \text{indep}$$

$$\frac{X}{n} = \overline{X} \implies \text{normal} \qquad \text{by CLT.}$$

$$\frac{\frac{X}{n} - p}{\sqrt{\frac{p(1-p)}{n}}} = \frac{\frac{X}{n} - \mu}{\sqrt{\sigma^2/n}} \implies N(0,1)$$

$$\mu = p$$
$$\sigma^2 = p(1-p)$$

$$= \frac{X - np}{\sqrt{np(1-p)}}$$

# The Central Limit Theorem

**Normal approximation to Binomial**

If $X \sim \mathsf{Bin}(n, p)$ and $n$ is large enough,

$X$ and $X/n$ have approximately normal distribution.

# Exercise

(5.4-56) A binary communication channel transmits a sequence of "bits" (0s and 1s). Suppose that for any particular bit transmitted, there is a 10% chance of a transmission error (a 0 becoming a 1 or a 1 becoming a 0).

Assume that bit errors occur independently of one another.

1. Consider transmitting 1000 bits. What is the approximate probability that at most 125 transmission errors occur?

2. Suppose the same 1000-bit message is sent two different times independently of one another.
   What is the approximate probability that the number of errors in the first transmission is within 50 of the number of errors in the second?

① $X = \#$ of errors in 1000 bits $\sim Bin(1000, 0.1)$

$$X \approx N(100, 90)$$
$$P(X \leq 125) \approx P\left(Z \leq \frac{125 - 100}{\sqrt{90}}\right)$$

with half-unit correction:
$$P(X \leq 125) = P(X \leq 125.5) \approx P\left(Z \leq \frac{125.5 - 100}{\sqrt{90}}\right)$$

② $X_1 \sim Bin(1000, 0.1)$
$X_2 \sim Bin(1000, 0.1)$ $\Big\}$ Indep.

$$\mathbb{P}\left(-50 \leq \underbrace{X_1 - X_2} \leq 50\right)$$

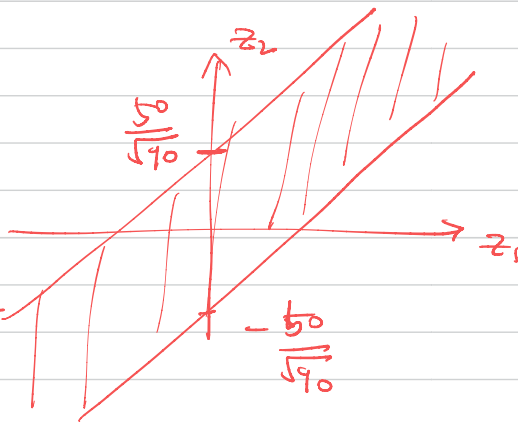$$\left(\frac{X_1-100}{\sqrt{90}} - \frac{X_2-100}{\sqrt{90}} = \frac{X_1-X_2}{\sqrt{90}}\right)$$

$$= \mathbb{P}\left(-\frac{50}{\sqrt{90}} \leq \underbrace{\left(\frac{X_1-100}{\sqrt{90}}\right)}_{\approx Z_1} - \underbrace{\left(\frac{X_2-100}{\sqrt{90}}\right)}_{\approx Z_2} \leq \frac{50}{\sqrt{90}}\right)$$

$$\approx \mathbb{P}\left(-\frac{50}{\sqrt{90}} \leq Z_1 - Z_2 \leq \frac{50}{\sqrt{90}}\right) \qquad Z_1, Z_2 \sim N(0,1) \text{ indep}$$

$$= \int_{-\infty}^{\infty} \int_{-\frac{50}{\sqrt{90}}+Z_2}^{\frac{50}{\sqrt{90}}+Z_2} \frac{1}{\sqrt{2\pi}} e^{-\frac{Z_1^2}{2}} \cdot \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{Z_2^2}{2}} \, dz_1 \, dz_2$$

$$\mathbb{P}\left((Z_1, Z_2) \in \underset{\sim}{A}\right)$$

$$= \iint_{\boxed{A}} f(z_1, z_2) \, dz_1 \, dz_2$$

# Section 5.
# The Distribution of a Linear Combination

**Definition**

Given a collection of $n$ random variables $X1, \cdots, X_n$ and constants $a_1, \cdots, a_n$,

$$Y = a_1X_1 + a_2X_2 + \cdots + a_nX_n$$

is called **a linear combination** of the $X_i$'s.

## Examples

① if $\quad a_1 = a_2 = \cdots = a_n = 1$

$$Y = X_1 + \cdots + X_n = T$$

② if $\quad a_1 = a_2 = \cdots = a_n = \frac{1}{n}$

$$Y = \frac{1}{n} X_1 + \frac{1}{n} X_2 + \cdots + \frac{1}{n} X_n = \overline{X}.$$

### Proposition

For a collection of $n$ random variables $X1, \cdots, X_n$ and constants $a_1, \cdots, a_n$, consider

$$Y = a_1 X_1 + a_2 X_2 + \cdots + a_n X_n$$

Then,

$$\mathbb{E}[Y] = a_1 \mathbb{E}[X_1] + a_2 \mathbb{E}[X_2] + \cdots + a_n \mathbb{E}[X_n]$$

$$\text{Var}(Y) =$$

In particular, if they are independent,

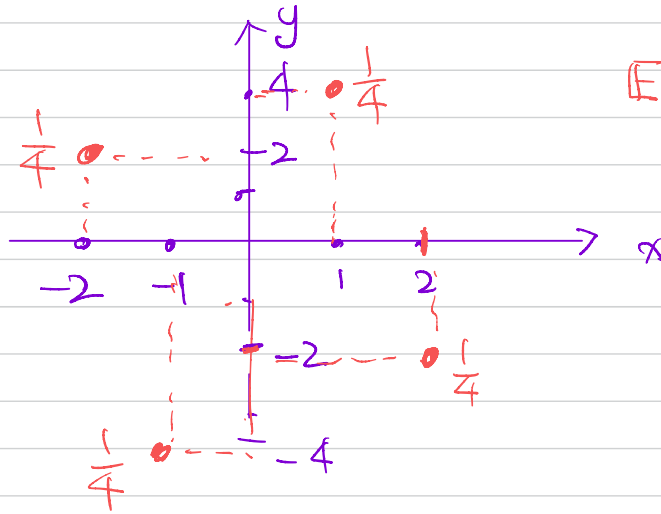$$\text{Var}(Y) = \sum_{i=1}^{n} a_i^2 \, \text{Var}(X_i)$$

38

$$\text{Var}(Y) = \sum_{i=1}^{n} \text{Var}(a_i X_i) + 2 \sum_{1 \leq i < j \leq n} \text{Cov}(X_i, X_j)$$

$a_i \quad a_j$ \qquad $\binom{n}{2}$ terms

$$= \sum_{i=1}^{n} a_i^2 \, \text{Var}(X_i) + 2 \sum_{1 \leq i < j \leq n} a_i a_j \, \text{Cov}(X_i, X_j)$$

⟨ Counter Example ⟩

$\underline{Cov(X, Y)} = 0$      but      Not Indep.
‖

$\mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$



$\mathbb{E}[X \cdot Y] = \frac{1}{4} \cdot (4 - 4 + 4 - 4)$

$= 0$

$\therefore Cov(X, Y) = 0$

$\mathbb{P}(X = 2, Y = -2) = \frac{1}{4}$

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

$\mathbb{P}(Y = -2) = \frac{1}{4}$

$X \sim Unif(-1, 1)$        $Y = |X| \sim Unif(0, 1)$

$\mathbb{P}(Y \leq t) = \mathbb{P}(|X| \leq t)$

$= \mathbb{P}(-t \leq X \leq t)$

$Cov(X, Y)$

$= \mathbb{E}[\widehat{XY}] - \underline{\mathbb{E}[X]}\,\mathbb{E}[Y]$
                    ‖
                    0

$= \int_{-1}^{1} x |x| \cdot \frac{1}{2} \, dx = 0$

$= \int_{-t}^{t} \frac{1}{2} \, dx$

$= \frac{2t}{2} = t$

$$\mathbb{E}[\, a_1 X_1 + a_2 X_2 + \cdots + a_n X_n \,]$$
$$= a_1 \mathbb{E}[X_1] + a_2 \mathbb{E}[X_2] + \cdots + a_n \mathbb{E}[X_n]$$

$$\mathrm{Var}(\, a_1 X_1 + a_2 X_2 + \cdots + a_n X_n \,)$$
$$= a_1^2 \mathrm{Var}(X_1) + a_2^2 \mathrm{Var}(X_2) + \cdots + a_n^2 \mathrm{Var}(X_n)$$
$$+ 2\Big( \mathrm{Cov}(X_1, X_2) + \mathrm{Cov}(X_1, X_3) + \cdots + \mathrm{Cov}(X_1, X_n)$$
$$+ \mathrm{Cov}(X_2, X_3) + \mathrm{Cov}(X_2, X_4) + \cdots \Big)$$

## Linear Combination

### Example

A certain automobile manufacturer equips a particular model with either a six-cylinder engine or a four-cylinder engine.

Let $X_1$ and $X_2$ be fuel efficiencies for independently and randomly selected six-cylinder and four-cylinder cars, respectively, with

$$\mu_1 = 22, \qquad \mu_2 = 26, \qquad \sigma_1 = 1.2, \qquad , \sigma_2 = 1.5.$$

Find $\mathbb{E}[X_1 - X_2]$ and $\mathrm{Var}(X_1 - X_2)$.

$$\mathbb{E}[X_1 - X_2] = \mathbb{E}[X_1] - \mathbb{E}[X_2] = 22 - 26 = -4$$
$$\mathrm{Var}(X_1 - X_2) = 1^2 \mathrm{Var}(X_1) + (-1)^2 \mathrm{Var}(X_2)$$
$$= (1.2)^2 + (1.5)^2$$
$$= 1.44 + 2.25$$
$$= 3.69 .$$

# Linear Combination

## Proposition

If $X_1, X_2, \cdots, X_n$ are independent, normally distributed RVs (with possibly different means and/or variances), then any linear combination also has a normal distribution.

In particular, the difference $X_1 - X_2$ between two independent, normally distributed variables is itself normally distributed.

$$X_i \sim N(\mu_i, \sigma_i^2) \qquad i = 1, \cdots, n$$

$$\Rightarrow a_1 X_1 + a_2 X_2 + \cdots + a_n X_n \sim N(\mu, \sigma^2)$$

$$\mu = a_1 \mu_1 + a_2 \mu_2 + \cdots + a_n \mu_n$$

$$\sigma^2 = a_1^2 \sigma_1^2 + a_2^2 \sigma_2^2 + \cdots + a_n^2 \sigma_n^2.$$

Example
- $X_1, \cdots, X_n \quad i.i.d. \quad N(\mu, \sigma^2)$

$$\bar{X} = \frac{1}{n}(X_1 + X_2 + \cdots + X_n) \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

- $X_1, \cdots, X_n \quad i.i.d. \quad \mu < \infty, \sigma^2 < \infty$

$$\bar{X} \approx N\left(\mu, \frac{\sigma^2}{n}\right) \qquad \text{by CLT}$$

## Exercise

(5.5-62) Manufacture of a certain component requires three different machining operations.

Machining time for each operation has a normal distribution, and the three times are independent of one another.

The mean values are 15, 30, and 20 min, respectively, and the standard deviations are 1, 2, and 1.5 min, respectively.

What is the probability that it takes at most 1 hour of machining time to produce a randomly selected component?