

STAT 305: Chapter 5

Part IV

Amin Shirazi

ashirazist.github.io/stat305.github.io

Chapter 5 4. Joint Distributions and Independence

Working with Multiple Random Variables

Joint Distributions

Joint Distributions

We often need to consider two random variables together - for instance, we may consider

- the length and weight of a squirrel,
- the loudness and clarity of a speaker,
- the blood concentration of Protein A, B, and C and so on.

This means that we need a way to describe the probability of two variables *jointly*. We call the way the probability is simultaneously assigned the "joint distribution".

Joint Distributions

Discrete RVs

Joint distribution of discrete random variables

For several discrete random variable, the device typically used to specify probabilities is a *joint probability function*. The two-variable version of this is defined.

A **joint probability function (joint pmf)** for discrete random variables X and Y is a nonnegative function $f(x, y)$, giving the probability that (simultaneously) X takes the values x and Y takes the values y . That is,

$$f(x, y) = P[X = x \text{ and } Y = y]$$

Properties of a valid joint pmf:

- $f(x, y) \in [0, 1]$ for all x, y $f(x, y) \geq 0$ $\sum_{x,y} f(x, y)$
- $\sum_{x,y} f(x, y) = 1$

Joint Distributions

Discrete RVs

Joint distribution of discrete random variables

So we have probability functions for X , probability functions for Y and now a probability function for X and Y together - that's a lot of f s floating around though! In order to be clear which function we refer to when we refer to " f ", we also add some subscripts

Suppose X and Y are two discrete random variables.

- we may need to identify the *joint probability function* using $\underline{f_{XY}(x, y)}$,
- we may need to identify the probability function of X by itself (aka the *marginal probability function* for X) using $\underline{f_X(x)}$,
- we may need to identify the probability function of Y by itself (aka the *marginal probability function* for Y) using $\underline{f_Y(y)}$

Joint Distributions

Discrete RVs

Joint pmf

For the discrete case, it is useful to give $f(x, y)$ in a **table**.

Two bolt torques, cont'd

Recall the example of measure the bolt torques on the face plates of a heavy equipment component to the nearest integer. With

X = the next torque recorded for bolt 3

Y = the next torque recorded for bolt 4



Joint Distributions

Discrete RVs

Joint pmf

the joint probability function, $f(x, y)$, is

$x \rightarrow$	11	12	13	14	15	16	17	18	19	20
$y \backslash x$	20	0	0	0	0	0	0	2/34	2/34	1/34
20	0	0	0	0	0	0	2/34	0	0	0
19	0	0	0	0	0	0	2/34	0	0	0
18	0	0	1/34	1/34	0	0	1/34	1/34	1/34	0
17	0	0	0	0	2/34	1/34	1/34	2/34	0	0
16	0	0	0	1/34	2/34	2/34	0	0	2/34	0
15	1/34	1/34	0	0	3/34	0	0	0	0	0
14	0	0	0	0	1/34	0	0	2/34	0	0
13	0	0	0	0	1/34	0	0	0	0	0

$$P(x=12, y=15) = \frac{1}{34}$$

$$P(x=19, y=16) = \frac{2}{34}$$

Joint Distributions

Calculate:

- $P[X = 14 \text{ and } Y = 19]$

Discrete RVs

- $P[X = 18 \text{ and } Y = 17]$

Joint Distributions

Discrete RVs

By summing up certain values of $f(x, y)$, probabilities associated with X and Y with patterns of interest can be obtained.

Consider: $P(X \geq Y) = P(X=13 \& Y=13) + P(X=14, Y=13)$

$$+ \dots + P(X=20 \& Y=20)$$

$$= \frac{3}{34} + \frac{4}{34} + \dots + \frac{1}{34} = \frac{17}{34}$$

$y \setminus x$	11	12	13	14	15	16	17	18	19	20
20										α
19									α	α
18									α	α
17								α	α	α
16						α	α	α	α	α
15					α	α	α	α	α	α
14			α							
13	α									

$f(13, 13)$ $f(14, 13)$

Joint Distributions

Discrete RVs

Pl "the torque recorded for the 3rd & 4th bolt are 0 or n Lb from each other ")

$$P(|X - Y| \leq 1) = P(X=12 \& Y=13) + P(X=13 \& Y=12) + \dots + P(X=20 \& Y=20) = \frac{2}{34} + \frac{3}{34} + \dots + \frac{1}{34} = \frac{18}{34}$$

y \ x	11	12	13	14	15	16	17	18	19	20
20									x	x
19								x	x	x
18							x	x	x	
17						x	x	x		
16					x	x	x			
15			x	x	x					
14		x	x	x						
13	x	x	x							

Joint Distributions

Discrete RVs

$P(X = 17)$

$f_X(x=17)$

$y \downarrow$

$x \rightarrow$

$y \setminus x$	11	12	13	14	15	16	17	18	19	20
20	X						X			
19							X			
18							X			
17							X			
16							X			
15							X			
14							X			
13							X			

$P_Y(19)$

$$P(X=17) = P(X=17 \& Y=13) + P(X=17 \& Y=14) + \dots + P(X=17, Y=20)$$

$$= \frac{1}{34} + \frac{1}{34} + \dots + \frac{1}{34} = \frac{4}{34}$$

Marginal Distribution

Joint Distributions

Discrete RVs

Marginal distributions

In a bivariate problem, one can add down columns in the (two-way) table of $f(x, y)$ to get values for the probability function of X , $f_X(x)$ and across rows in the same table to get values for the probability distribution of Y , $f_Y(y)$.

The individual probability functions for discrete random variables X and Y with joint probability function $f(x, y)$ are called **marginal probability functions**. They are obtained by summing $f(x, y)$ values over all possible values of the other variable.

Joint Distributions

Discrete RVs

Connecting Joint and Marginal Distributions

In continuous joint

dist. of

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

Use: Joint to Marginal for Discrete RVs

Let X and Y be discrete random variables with joint probability function Then the marginal probability function for X can be found by:

$$\underbrace{f_X(x)}_{\text{a function of } x} = \sum_y f_{XY}(x,y)$$

and the marginal probability function for Y can be found by:

$$\downarrow \quad f_Y(y) = \sum_x f_{XY}(x,y)$$

a function of y

Joint Distributions

Discrete RVs

y

Example: [Torques, cont'd]

Find the marginal probability functions for X and Y from the following joint pmf.

$y \setminus x$	11	12	13	14	15	16	17	18	19	20	$P_Y(y)$
20	0	0	0	0	0	0	0	2/34	2/34	1/34	5/34
19	0	0	0	0	0	0	2/34	0	0	0	2/34
18	0	0	1/34	1/34	0	0	1/34	1/34	1/34	0	5/34
17	0	0	0	0	2/34	1/34	1/34	2/34	0	0	6/34
16	0	0	0	1/34	2/34	2/34	0	0	2/34	0	7/34
15	1/34	1/34	0	0	3/34	0	0	0	0	0	8/34
14	0	0	0	0	1/34	0	0	2/34	0	0	9/34
13	0	0	0	0	1/34	0	0	0	0	0	10/34

$$P_X(x=x) : \underbrace{y_{34} \ y_{34} \ y_{34}}_{\text{1st row}} \ \underbrace{2y_{34} \ 2y_{34}}_{\text{2nd row}} \ \underbrace{\frac{9}{34} \ \frac{3}{34}}_{\text{3rd row}} \ \underbrace{4y_{34} \ 7y_{34}}_{\text{4th row}} \ \underbrace{\frac{5}{34} \ y_{34}}_{\text{5th row}}$$

(1)

marginal dist. of X and/or Y

So,

X	$P_X(x) = \sum_{y=13}^{20} P(x,y)$
11	1/34
12	1/34
13	1/34
14	2/34
15	9/34
16	3/34
17	4/34
18	7/34
19	5/34
20	1/34

Y	$P_Y(y) = \sum_{x=11}^{20} P(x,y)$
13	1/34
14	3/34
15	5/34
16	7/34
17	6/34
18	5/34
19	2/34
20	5/34

$$E X = \sum_{x=11}^{20} x \cdot P_x(x) = 11 \left(\frac{1}{34} \right) + 12 \left(\frac{1}{34} \right) \\ \dots + 20 \left(\frac{1}{34} \right) = \dots$$

Joint Distributions

Getting marginal probability functions from joint probability functions begs the question whether the process can be reversed.

Discrete RVs

Can we find joint probability functions from marginal probability functions?

No! \rightarrow we need more information (later!)

	x	1	2	3	$f_y(y)$
y					
1	0.4	0	0		0.4
2	0	0.4	0		0.4
3	0	0	0.2		0.2
		0.4	0.4	0.2	$f_x(x)$

	x	1	2	3	$f_y(y)$
y					
1	0.16	0.16	0.08		0.4
2	0.16	0.16	0.08		0.4
3	0.08	0.08	0.04		0.2
		0.4	0.4	0.2	$f_x(x)$

Note: $P(x=1, y=1) = 0.4$

\neq

$P(x=1, y=1) = 0.16$

Conditional Distribution

Joint Distributions

Discrete RVs

Conditional Distribution

Conditional Distribution of Discrete Random Variables

When working with several random variables, it is often useful to think about what is expected of one of the variables, given the values assumed by all others.

For discrete random variables X and Y with joint probability function $f(x, y)$, the **conditional probability function of X given $Y = y$** is a function of x

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

joint dist. of x, y
marginal f_Y

and the **conditional probability function of Y given $X = x$** is a function of y

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

Joint Distributions

Discrete RVs

Conditional Distribution

Example: [Torque, cont'd]

y \ x	11	12	13	14	15	16	17	18	19	20
20	0	0	0	0	0	0	0	2/34	2/34	1/34
19	0	0	0	0	0	0	2/34	0	0	0
18	0	0	1/34	1/34	0	0	1/34	1/34	1/34	0
17	0	0	0	0	2/34	1/34	1/34	2/34	0	0
16	0	0	0	1/34	2/34	2/34	0	0	2/34	0
15	1/34	1/34	0	0	3/34	0	0	0	0	0
14	0	0	0	0	1/34	0	0	2/34	0	0
13	0	0	0	0	1/34	0	0	0	0	0

9/34

$$f_{Y|X}(18) = \frac{9}{34}$$

Find the following probabilities:

$$\begin{aligned} \bullet f_{Y|X}(20|18) &= \frac{P(x=18, y=20)}{P_x(18)} \\ &= \frac{P(x=18, y=20)}{P_x(20)} \\ &= \frac{2/34}{7/34} = \frac{2}{7} \end{aligned}$$

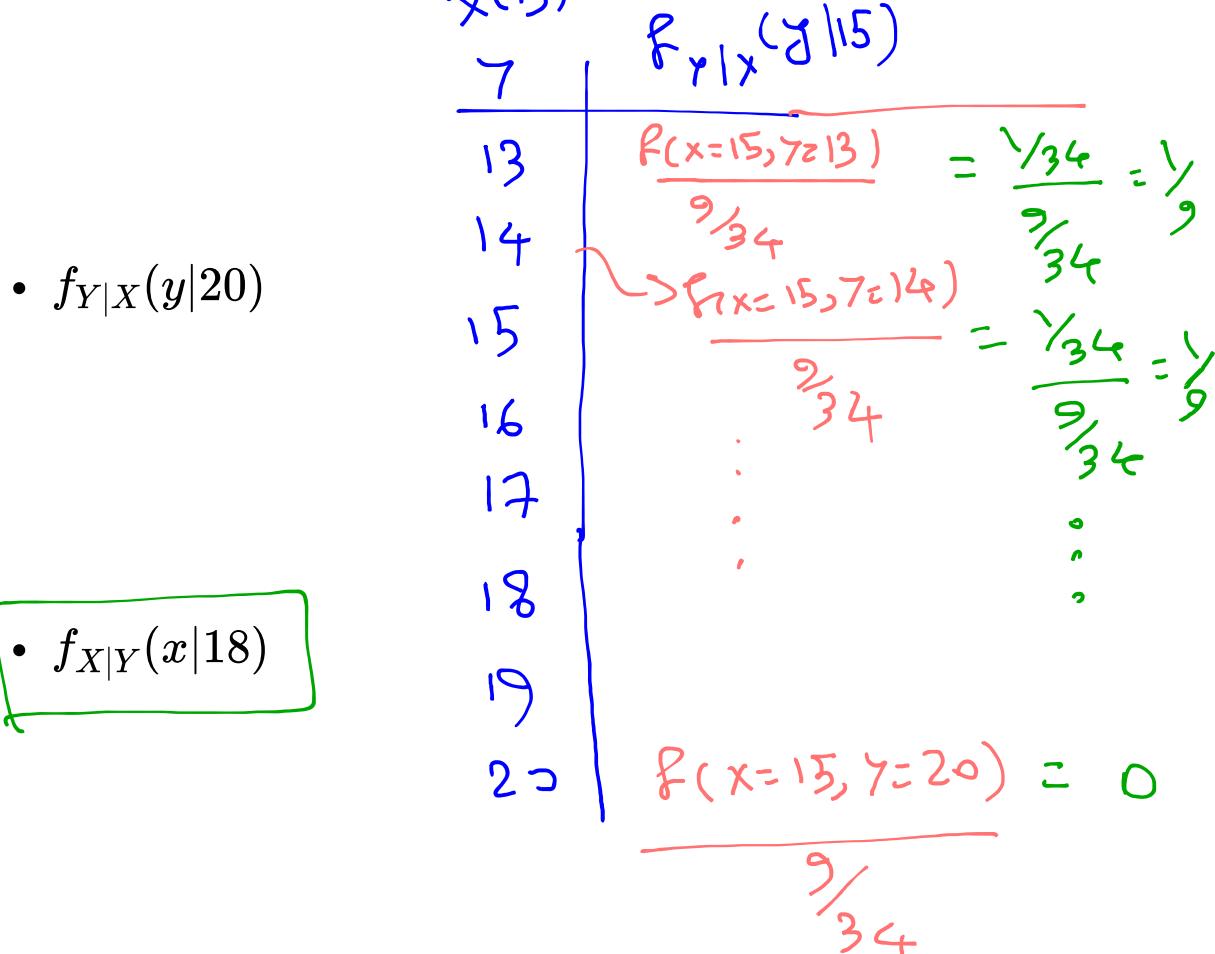
Joint Distributions

Discrete RVs

Conditional Distribution

Example: [Torque, cont'd]

$$\bullet f_{Y|X}(y|15) = \frac{P(x=15, y=y)}{P_x(15)} = \frac{P(x=15, y=5)}{9/34}$$



$$\bullet f_{X|Y}(x|18)$$

$$f_{X|Y}(x|y=18) = \frac{P(x, y=18)}{P(y=18)}$$

x	$f_{X Y}(x y=18)$
11	0
12	0
13	$\frac{1}{5}$
14	$\frac{1}{5}$
15	0
16	0
17	$\frac{1}{5}$
18	$\frac{1}{5}$
19	$\frac{1}{5}$
20	0

Notes Conditional
distributions

are Valid
dist.

$$E(X|Y=15) = \sum_{x=11}^{20} x \cdot f_{X|Y}(x|y=15)$$

$$= 11(0) + 12(0) + 13(\frac{1}{5}) + \\ 14(\frac{1}{5}) + 15(0) + 16(0) + 17(\frac{1}{5}) + \\ 18(\frac{1}{5}) + 19(\frac{1}{5}) + 20(0) = \dots$$

Independence

Joint Distributions

Discrete RVs

Conditional Distribution

Independence

Let's start with an example. Look at the following joint probability distribution and the associated marginal probabilities.

$y \setminus x$	1	2	3	$f_Y(y)$
3	0.08	0.08	0.04	0.20
2	0.16	0.16	0.08	0.40
1	0.16	0.16	0.08	0.40
$f_X(x)$	0.40	0.40	0.20	1.00

What do you notice?

①

$$P_{x,y}(x,y) = P_x(x) \cdot P_y(y)$$

②

$$P_{y|x}(y|3) = \frac{P(x=3, y)}{P_x(3)} = \frac{f(x=3, y=1)}{0.2}$$

$$\begin{array}{c|ccc}
 Y & f_{Y|x}(y|x=3) = & \frac{P(x=3,y)}{P_x(3)} = & \frac{f(x=3,y)}{0.2} \\
 \hline
 1 & \frac{0.08}{0.2} = 0.4 & = f_Y(y=1) \\
 2 & \frac{0.08}{0.2} = 0.4 & = f_Y(y=2) \\
 3 & \frac{0.04}{0.2} = 0.2 & = f_Y(y=3)
 \end{array}$$

$f_{Y|x}(y|x=3) = f_Y(y)$. Actually this is true for all values of x , i.e knowing what value x takes, doesn't matter

in the questions about Y .

$\Rightarrow X$ and Y are independent.

Joint Distributions

Discrete RVs

Conditional Distribution

Independence

Discrete random variables X and Y are **independent** if their joint distribution function $f(x, y)$ is the product of their respective marginal probability functions. This is,

independence means that

$$f(x, y) = f_X(x)f_Y(y)$$



for all x, y .

If this does not hold, then X and Y are **dependent**

Alternatively, discrete random variables X and Y are independent if for all x and y ,

If X and Y are not only independent but also have the same marginal distribution, then they are **independent and identically distributed (iid)**.

$$F_{Y|X}(y|x) = F_Y(y) \quad \& \quad F_{X|Y}(x|y) = F_X(x)$$

Chapter 5.5: Functions of Random Variables

Results and Theorems

Functions of RVs

Functions of Random Variables

A random variable can be thought of as a function whose input is an outcome and whose output is a real number. When we take a function of the value the random variable takes, the resulting value is still depends on the outcome of a random experiment - in other words: functions of random variables are random variables.

This means that a function of a random variable will have probabilities attached to the value it takes, based on the value taken by the random variable. It also means functions of random variables will have:

- probability functions (if discrete) or probability density functions (if continuous)
- cumulative probability functions (if discrete) or cumulative density functions (if continuous)
- expected values and variances ...

Functions of RVs

Linear Combinations

Linear combinations

For engineering purposes, it often suffices to know the mean and variance for a function of several random variables, $U = g(X_1, X_2, \dots, X_n)$ (as opposed to knowing the whole distribution of U). When g is **linear**, there are explicit functions.

Proposition: If X_1, X_2, \dots, X_n are n **independent** random variables and a_0, a_1, \dots, a_n are $n + 1$ constants, then consider

Recall: X is a C.V

$\Rightarrow ax + b$
is a C.V

generalize
the idea

$$U = a_0 + a_1 X_1 + a_2 X_2 + \dots + a_n X_n$$

U is itself a random variable as it is a linear combination of n *independent* random variables

Functions of RVs

Linear Combinations

Linear combinations [cont'd]

Recall: $E(ax + b) = aE(x) + b$ & $\text{Var}(ax + b) = a^2 \text{Var}(x)$

U, as a random variable has mean

$$\rightarrow EU = a_0 + a_1 E(X_1) + a_2 E(X_2) + \cdots + a_n E(X_3)$$

and variance

$$\rightarrow \text{Var}U = a_1^2 \text{Var}X_1 + a_2^2 \text{Var}X_2 + \cdots + a_n^2 \text{Var}X_3$$

These hold true when X_1, \dots, X_n are

$\underbrace{\text{iid}}$

Note: In general $\text{Var}(X_1 + X_2) \neq \text{Var}(X_1) + \text{Var}(X_2)$
(just in iid form it's true)

Functions of RVs

Linear Combinations

Example:

Say we have two independent random variables X and Y with $\text{E}X = 3.3$, $\text{Var}X = 1.91$, $\text{E}Y = 25$, and $\text{Var}Y = 65$. Find the mean and variance for

- $\underline{U} = 3 + 2X - 3Y$

$$\begin{aligned}\text{E}U &= \text{E}(3 + 2X - 3Y) = 3 + 2\text{E}X - 3\text{E}Y \\ &= 3 + 2(3.3) - 3(25) = -65.4\end{aligned}$$

$$\begin{aligned}\text{V}(U) &= \text{V}(3 + 2X - 3Y) \\ &= \text{V}(3) + 2^2 \text{V}(X) + (-3)^2 \text{V}(Y) \\ &= 0 + 4(1.91) + 9(65) = 592.64\end{aligned}$$

$$\begin{aligned}\bullet \quad \underline{V} &= -4X + 3Y \quad \xrightarrow{\text{indep.}} \\ \Rightarrow \text{E}V &= \text{E}(-4\underbrace{X}_{\text{indep.}} + 3Y) = -4\text{E}X + 3\text{E}Y\end{aligned}$$

$$\left. \begin{aligned} &\text{& } \text{V}(V) = \text{Var}(-4X + 3Y) = \text{Var}(-4X) + \text{V}(3Y) \\ &= (-4)^2 \text{Var}(X) + 3^2 \text{Var}(Y) \end{aligned} \right\} \begin{aligned} &= -4(3.3) + 3(25) \\ &= 61.8 \end{aligned}$$

Functions of RVs

Linear Combinations

$$= 16(1.91) + 9(65) = 615.56$$

Example:

Say $X \sim \text{Binomial}(n = 10, p = 0.5)$ and $Y \sim \text{Poisson}(\lambda = 3)$. Calculate the mean and variance of $Z = \underline{\underline{5 + 2X - 7Y}}$.

$\underline{\underline{X, Y \text{ are independent.}}}$

Note: $X \sim \text{Binomial}(10, 0.5) \rightarrow E(X) = 10(0.5) = 5$

$$\rightarrow V(X) = 10(0.5)(1-0.5) = 2.5$$

$$Y \sim \text{Poisson}(\lambda = 3) \rightarrow E(Y) = \text{Var } Y = 3$$

$$\begin{aligned} E(Z) &= E(5 + 2X - 7Y) = E(5) + E(2X) + E(-7Y) \\ &= 5 + 2E(X) - 7E(Y) \\ &= 5 + 2(5) - 7(3) = -6 \end{aligned}$$

$$\text{Var}(z) = \text{Var}(5 + 2x - 7y) =$$

$$= \underbrace{\text{Var}(5)}_{0} + \text{Var}(2x) + \text{Var}(-7y)$$

$$= 0 + 4\text{Var}(x) + 49\text{Var}(y)$$

$$= 4(2.5) + 49(3) = \underbrace{157}_{157}$$

Functions of RVs

Linear Combinations

Sample Mean

A particularly important use of functions of random variables concerns n iid random variables where each $a_i = \frac{1}{n}$ for $i = 1, 2, \dots, n$. Then we can define the random variable \bar{X} as follows

$$\bar{X} = \frac{1}{n} X_1 + \dots + \frac{1}{n} X_n = \frac{1}{n} \sum_{i=1}^n X_i$$

Note that \bar{X} is a random variable

→ We can then find the mean and variance of this random variable.

Functions of RVs

Linear Combinations

Sample Mean

$$\bar{X} = \frac{1}{n} X_1 + \cdots + \frac{1}{n} X_n = \frac{1}{n} \sum_{i=1}^n X_i$$

as they relate to the population parameters $\mu = E X_i$ and $\sigma^2 = \text{Var } X_i$.

For **independent** variables X_1, \dots, X_n with common mean μ and variance σ^2 ,

$$E(\bar{X}) : E\left(\frac{1}{n} x_1 + \frac{1}{n} x_2 + \cdots + \frac{1}{n} x_n\right)$$

$$= \frac{1}{n} E(x_1) + \cdots + \frac{1}{n} E x_n$$

$$x_1, \dots, x_n \stackrel{\text{iid}}{\sim} = \frac{1}{n} \underbrace{\mu + \mu + \cdots + \mu}_{n \text{ times}} = n \cdot \frac{1}{n} \mu = \mu$$

$$\text{Var}(\bar{X}) : \text{Var}\left(\frac{1}{n} x_1 + \cdots + \frac{1}{n} x_n\right)$$

$$\text{indep.} = \text{Var}\left(\frac{1}{n} x_1\right) + \cdots + \text{Var}\left(\frac{1}{n} x_n\right)$$

$$= \frac{1}{n^2} \text{Var}(x_1) + \cdots + \frac{1}{n^2} \text{Var}(x_n)$$

Functions of RVs

Linear Combinations

Sample Mean

$$E x_1 = E x_{1,0} = \underline{\mu}$$

$$\text{var}(x_3) = \text{var}(x_{n-1}) = \underline{\sigma^2}$$

$$E x_i = \underline{\mu + i}$$

X

not identically distributed.

$$\text{var}(x_n) = \underline{\sigma^2 / n} \quad X$$

$$= \frac{1}{n^2} \sigma^2 + \dots + \frac{1}{n^2} \sigma^2 = n \cdot (\frac{1}{n^2} \sigma^2) = \underline{\frac{\sigma^2}{n}}$$

What is the point? ^{n times}

It does not matter if we are working with discrete or continuous random variables, as long as we have an independent and identically distributed (iid) sample of size n with the same mean μ and the same variance σ^2 , the random variable \bar{X} has

$$\underline{E(\bar{X}) : \mu} \quad \checkmark$$

and

$$\underline{V(\bar{X}) : \frac{\sigma^2}{n}} \quad \checkmark$$

The point is that the variance of a sample mean of size n is the population variance devided by the sample size n which makes it smaller

i.e. as the sample size increases, the variability of the sample mean decreases.

$$E x_i = \underline{n p}^{\text{fixed}} \quad \text{in binomial.}$$

Functions of RVs

Linear Combinations

Sample Mean

Example:[Seed lengths]

One botanist measured the length of 10 seeds from the same plant. The seed lengths measurements are

X_1, X_2, \dots, X_{10} . Suppose it is known that the seed lengths are iid with mean $\mu = 5$ mm and variance $\sigma^2 = 2$ mm.

Calculate the mean and variance of the average of 10 seed measurements.

$$E \bar{X} = E X_1 = \mu = 5$$

$$\text{Var}(\bar{X}) = \frac{\text{Var}(X_1)}{n} = \frac{\sigma^2}{10} = \frac{2}{10} = 0.2$$

Central Limit Theorem

The Most Important Result in Statistics

Functions of RVs

Linear Combinations

Sample Mean

CLT $n \geq 25$

Central limit theorem

One of the most frequently used statistics in engineering applications is the sample mean. We can relate the mean and variance of the probability distribution of the sample mean to those of a single observation when an iid model is appropriate.

In the case of the sample mean, if the sample size (n) is large enough, we can also approximate the shape of the probability distribution function of the sample mean!

Functions of RVs

Linear Combinations

Sample Mean

CLT

(CLT)

Central limit theorem

If X_1, \dots, X_n are **independent** and **identically** distributed (iid) random variable (with mean μ and variance σ^2), then for large n , the variable \bar{X} is approximately normally distributed. That is,

Sample mean $\bar{X} \sim \text{Normal} \left(\mu, \frac{\sigma^2}{n} \right)$

* This is one of the **most important** results in statistics.

Functions of RVs

Linear Combinations

Sample Mean

CLT

W_i discrete r.v.
 $i=1, 2$

Example: [Tool serial numbers]

Consider selecting the last digit of randomly selected serial numbers of pneumatic tools. Let

- * W_1 = the last digit of the serial number observed next Monday at 9am
- * W_2 = the last digit of the serial number observed the following Monday at 9am

A plausible model for the pair of random variables W_1, W_2 is that they are independent, each with the marginal probability function

$$f(w) = \begin{cases} .1 & w = 0, 1, 2, \dots, 9 \\ 0 & \text{otherwise} \end{cases}$$

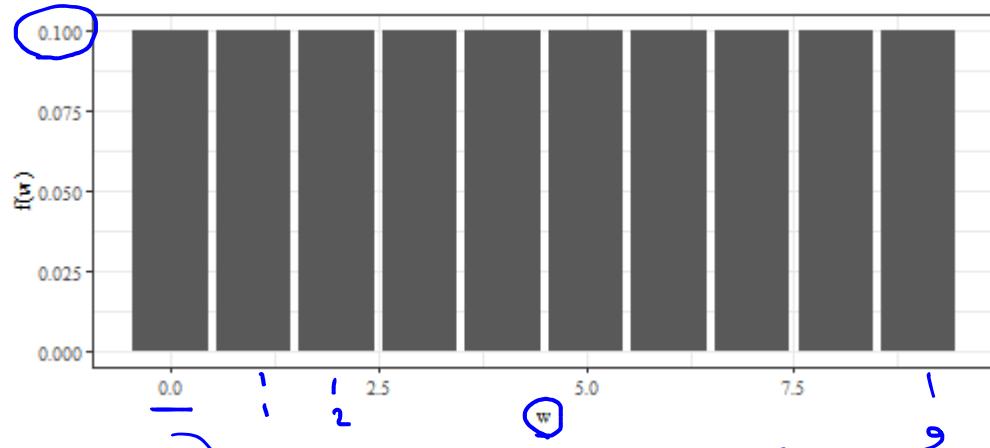
Functions of RVs

Linear Combinations

Sample Mean

CLT

Example: [Tool serial numbers]



With $\underline{EW} = 4.5$ and $\underline{\text{Var}}W = 8.25$.

Using such a distribution, it is possible to see that

→ $\overline{W} = \frac{1}{2}(W_1 + W_2)$ has probability distribution

$\begin{matrix} w_1 & \leftarrow w_2 \\ (0,0) & \rightarrow \end{matrix}$

\overline{w}	$f(\overline{w})$								
0.00	0.01	2.00	0.05	4.00	0.09	6.00	0.07	8	0.03
0.50	0.02	2.50	0.06	4.50	0.10	6.50	0.06	8.5	0.02
1.00	0.03	3.00	0.07	5.00	0.09	7.00	0.05	9	0.01
1.50	0.04	3.50	0.08	5.50	0.08	7.50	0.04		

$(w_1=1, w_2=2)$

$(\overline{w}_1=4.5, \overline{w}_2=4.5)$

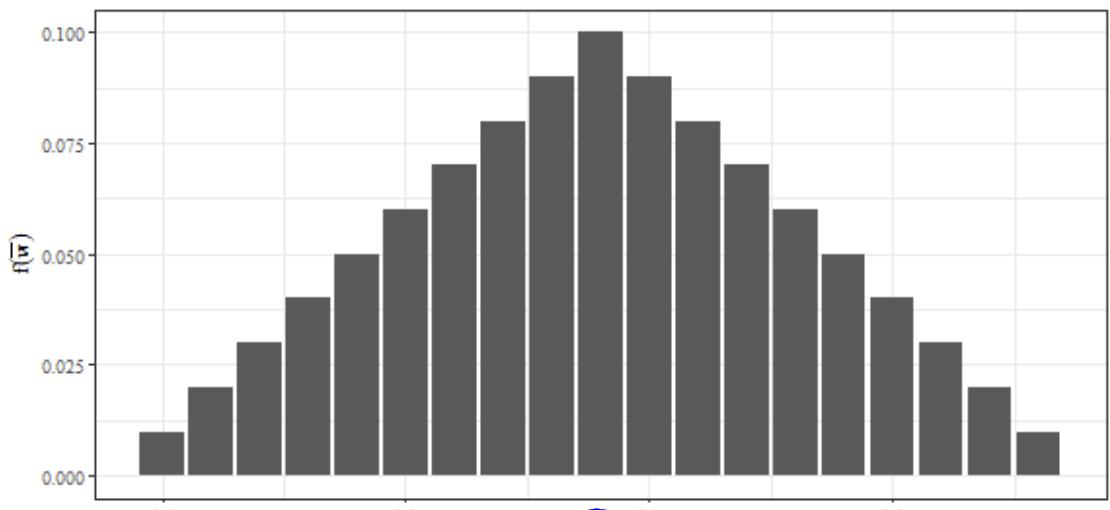
Functions of RVs

Linear Combinations

Sample Mean

CLT

Example: [Tool serial numbers]



$$E \bar{w} = E w_1 = 4.5$$

$$\text{Var}(\bar{w}) = \frac{\text{Var}(w_1)}{n} = 4.125$$

Comparing the two distributions, it is clear that even for a completely flat/uniform distribution of W and a small sample size of $n = 2$, the probability distribution of \bar{w} looks more bell-shaped than the underlying distribution.

Functions of RVs

Linear Combinations

Sample Mean

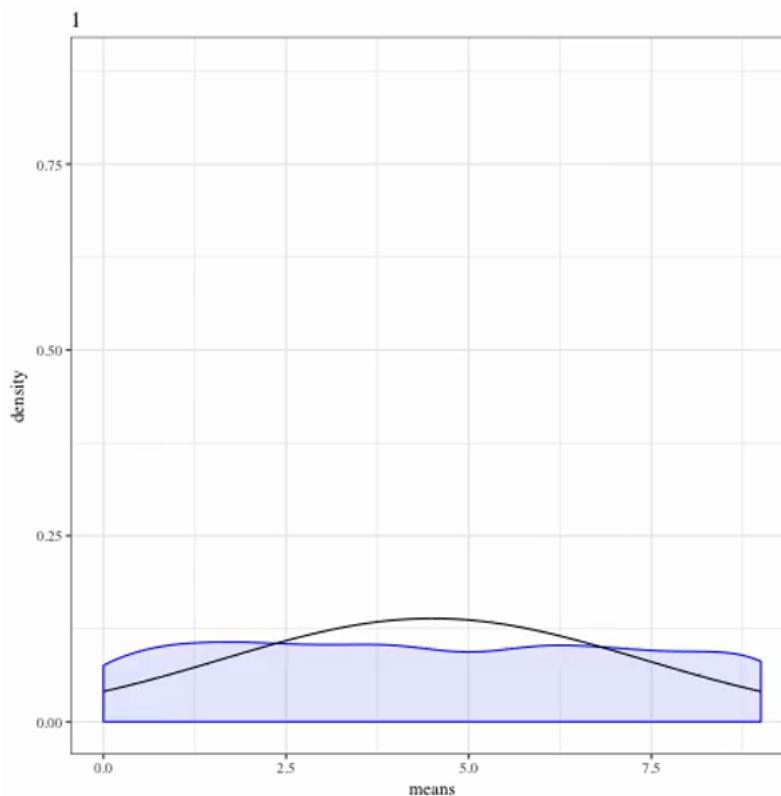
CLT

\bar{w} will always
have $E\bar{w} = Ew_i = 4.5$,
but $\text{var } \bar{w} = \frac{\text{var}(w_i)}{n}$
decreases as $n \rightarrow \infty$

and \bar{w} will approximately have Normal shape
(distribution)

Now consider larger and larger sample sizes,
 $n = 1, \dots, 40$:

Watch how CLT works [here](#)



Functions of RVs

Linear Combinations

Sample Mean

CLT

Example: [Stamp sale time]

Imagine you are a stamp salesperson (on eBay). Consider the time required to complete a stamp sale as S , and let

assume iid

\bar{S} = the sample mean time required to complete the next 100 sales

$n \geq 25$

Each individual sale time should have an $Exp(\alpha = 16.5s)$ distribution. We want to consider approximating $P[\bar{S} > 17]$.

$$S_i \stackrel{iid}{\sim} Exp(\alpha = 16.5) \rightarrow E S_i = \alpha = 16.5$$

$$\text{var } S_i = \alpha^2 = 16.5^2 = 272.25$$

one of them

↓

$$\text{Now: } E \bar{S} = E S_i = 16.5$$

$$\text{var}(\bar{S}) = \frac{\text{var}(S_i)}{n} = \frac{272.25}{100} = 2.72225$$

Since $n=100 \geq 25$,

$$\xrightarrow{1.65^2}$$

using CLT: $\bar{S} \sim N(\mu=16.5, \sigma^2 = 2.72225)$

$$\rightarrow P(\bar{S} > 17) = P\left(\frac{\bar{S} - 16.5}{\sqrt{2.72225}} > \frac{17 - 16.5}{\sqrt{2.72225}}\right)$$

$$= P(Z > \frac{17 - 16.5}{1.65})$$

0.303

$$= 1 - P(Z \leq 0.303)$$

$$= 1 - \Phi(0.303)$$

$$\text{table } \approx 1 - 0.6217 = 0.3783$$

Functions of RVs

Linear Combinations

Sample Mean

CLT

Example: [Cars]

Suppose a bunch of cars pass through certain stretch of road. Whenever a car comes, you look at your watch and record the time. Let X_i be the time (in minutes) between when the i^{th} car comes and the $(i + 1)^{th}$ car comes for $i = 1, \dots, 44$. Suppose you know the average time between cars is 1 minute.

Find the probability that the average time gap between cars for the next 44 cars exceeds 1.05 minutes.

x_i : time (minute) between i^{th} car & $(i+1)^{th}$ car \leftarrow

$$x_i \stackrel{\text{iid}}{\sim} \text{Exp}(\alpha=1)$$

$\bar{x} = \frac{1}{44} \sum_{i=1}^{44} x_i$: the average time gap between cars for 44 cars.

$$P(\bar{x} > 1.05) = ?$$

we have iid sample + $n = 44 \geq 25$

use CLT.

$$\bar{X} \sim N \left(E x_i = \alpha = 1, \frac{\text{Var}(x_i)}{n} = \frac{\alpha^2}{n} = \frac{1}{44} \right)$$

$$P(\bar{X} > 1.05) = P\left(\frac{\bar{X} - 1}{\sqrt{\frac{1}{44}}} > \frac{1.05 - 1}{\sqrt{\frac{1}{44}}}\right)$$

$$= P(Z > 0.332) \quad , Z \sim N(0,1)$$

$$= 1 - P(Z \leq 0.332)$$

$$= 1 - \Phi(0.332)$$

$$= 1 - 0.633 = 0.3669.$$

with 36.69% probability, the average time gap between 44

Cars is > 1.05 minute.

Functions of RVs

Linear Combinations

Sample Mean

CLT (assume iid)

Example: [Baby food jars, cont'd]

The process of filling food containers appears to have an inherent standard deviation of measured fill weights on the order of $1.6g$. Suppose we want to calibrate the filling machine by setting an adjustment knob and filling a run of n jars. Their sample mean net contents will serve as an indication of the process mean fill level corresponding to that knob setting.

we're looking for \underline{n} .

You want to choose a sample size, n , large enough that there is an 80% chance the sample mean is within $.3g$ of the actual process mean.



$$* P(\mu - 0.3 < \bar{x} < \mu + 0.3) = 0.8$$

Note that, μ & n are not given!

we're looking for \underline{n}

x_i : the weight of one jar.

\bar{x} : The sample mean weight of n jars.

$$E\bar{x} = E x_1 = E x_i = \mu$$

$$\text{Var}(\bar{x}) = \frac{\text{Var}(x_1)}{n} = \frac{1.6^2}{n}$$

For n large enough, by CLT: $\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$
 $\Rightarrow \bar{x} \sim N(\mu, \frac{1.6^2}{n})$

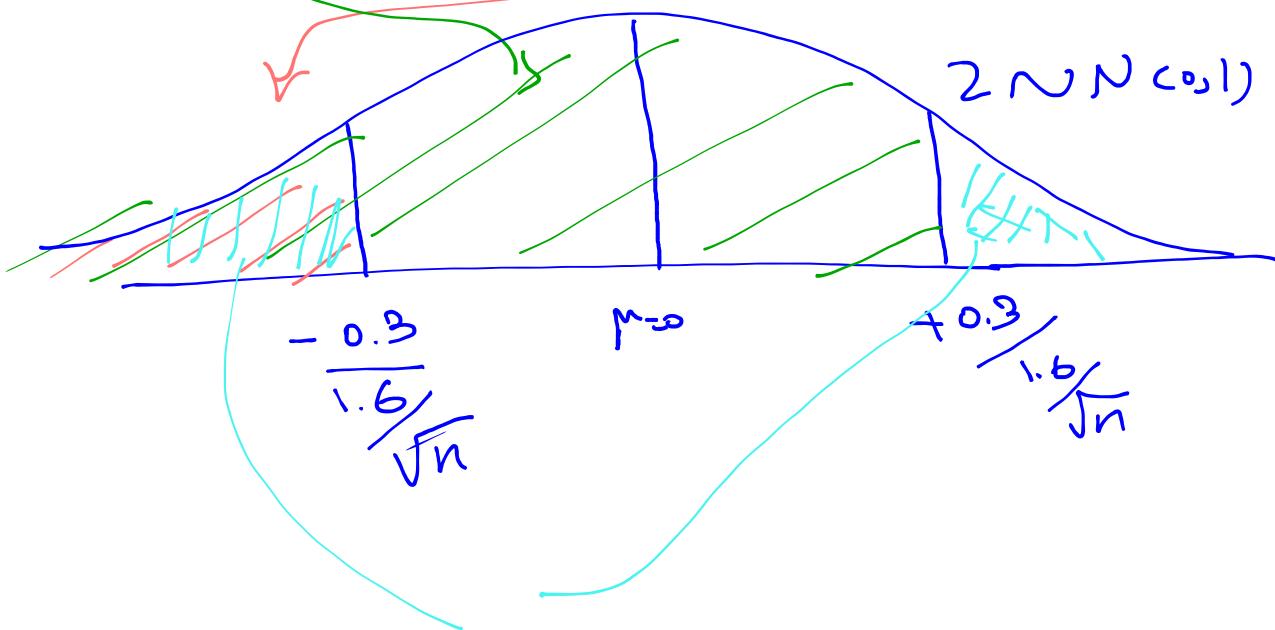
$$P(\mu - 0.3 < \bar{x} < \mu + 0.3) = 0.8$$

$$= P\left(\frac{\mu - 0.3 - \mu}{1.6/\sqrt{n}} < \frac{\bar{x} - \mu}{1.6/\sqrt{n}} < \frac{\mu + 0.3 - \mu}{1.6/\sqrt{n}} \right) = 0.8$$

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

$$= P\left(-\frac{0.3}{1.6/\sqrt{n}} < Z < \frac{+0.3}{1.6/\sqrt{n}}\right) = 0.8$$

$$= \Phi\left(\frac{0.3}{1.6/\sqrt{n}}\right) - \Phi\left(-\frac{0.3}{1.6/\sqrt{n}}\right) = 0.8$$



The same area

$$= \Phi\left(\frac{0.3}{1.6\sqrt{n}}\right) - \left(1 - \Phi\left(\frac{0.3}{1.6\sqrt{n}}\right)\right) = 0.8$$

$$= 2\Phi\left(\frac{0.3}{1.6\sqrt{n}}\right) - 1 = 0.8$$

$$\Rightarrow \Phi\left(\frac{0.3}{1.6\sqrt{n}}\right) = \frac{0.8}{2} = 0.9$$

by the

table $\Rightarrow \frac{0.3}{1.6\sqrt{n}} = 1.29 \Rightarrow \sqrt{n} = \frac{1.29 \times 1.6}{0.3}$

$$\Rightarrow n = \left(\frac{1.29 \times 1.6}{0.3}\right)^2 = 47.3344$$

choose : $\underbrace{n = 48}$

Functions of RVs

Linear Combinations

Sample Mean

CLT

$$\sigma^2 \rightarrow$$

Example: [Printing mistakes]

Suppose the number of printing mistakes on a page follows some unknown distribution with a mean of 4 and a variance of 9. Assume that number of printing mistakes on a printed page are iid.

- What is the approximate probability distribution of the average number of printing mistakes on 50 pages?

$$\bar{x}$$

$$n > 25$$

$$\bar{x} \sim N(4, \frac{9}{50}) \text{ by CLT.}$$

- Can you find the probability that the number of printing mistakes on a single page is less than 3.8?

$$x$$

No, because the probability dist. of # of printing mistakes is unknown on a single page.

Functions of RVs

Linear Combinations

Sample Mean

CLT

Example: [Printing mistakes]

$$\bar{x}$$

- Can you find the probability that the average number of printing mistakes on 10 pages is less than 3.8?

No, because $n=10 < 25$ & we cannot use CLT. Thus, the dist. of \bar{x} is unknown.

- Can you find the probability that the average number of printing mistakes on 50 pages is less than 3.8?

$$\bar{x}$$

Yes, because $n=50 > 25$ & x_i are iid.

So, by CLT $\bar{x} \sim N(\mu=4, \frac{9}{50})$