

Stat 202-2015XD-W4 - Thursday

## Review

The mean and standard deviation of a discrete random variable are calculated from the probability table

Values of $X$	$x_1$	$x_2$	$\dots$	$x_k$
Probability	$p_1$	$p_2$	$\dots$	$p_k$

$$\mu = \text{mean} = \sum x_i p_i = x_1 p_1 + \dots + x_k p_k$$

$$\sigma^2 = \text{Variance} = \sum (x_i - \mu)^2 p_i \\ = (x_1 - \mu)^2 p_1 + \dots + (x_k - \mu)^2 p_k$$

$$\sigma_x = \text{std dev} = \sqrt{\sigma^2}$$

## Formulas

All:  $\mu_{x+y} = \mu_x + \mu_y$  Independent:  $\sigma^2_{x+y} = \sigma^2_x + \sigma^2_y$

$$\mu_{x-y} = \mu_x - \mu_y \quad \sigma^2_{x-y} = \sigma^2_x + \sigma^2_y$$

## P Correlation

All  $\sigma^2_{x+y} = \sigma^2_x + \sigma^2_y + 2\rho\sigma_x\sigma_y$  <sup>↑ not a typo</sup>  
 $\sigma^2_{x-y} = \sigma^2_x + \sigma^2_y - 2\rho\sigma_x\sigma_y$

(pg 2)

## Linear transformation (all)

$$\mu_{ax+b} = a\mu_x + b$$

$$\sigma^2_{ax+b} = b^2 \sigma^2_x$$

Mistake in Wednesday's notes

Values of $x$	0	1	2	3
Probability	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{1}{8}$

$$\sigma^2_x = \frac{3}{4} \text{ not } \frac{9}{16}$$

$$\sigma_x = \sqrt{\frac{3}{4}}$$

~~BBM101~~

$$\mu_x = 0 \cdot \frac{1}{8} + 1 \cdot \frac{3}{8} + 2 \cdot \frac{3}{8} + 3 \cdot \frac{1}{8} = \frac{3}{2}$$

$$\sigma^2_x = (0 - \frac{3}{2})^2 \cdot \frac{1}{8} + (1 - \frac{3}{2})^2 \cdot \frac{3}{8}$$

$$+ (2 - \frac{3}{2})^2 \cdot \frac{3}{8} + (3 - \frac{3}{2})^2 \cdot \frac{1}{8}$$

$$= \frac{3}{4}$$

$$\sigma_x = \sqrt{\frac{3}{4}}$$

(Pg 3)

## Why are these formulas useful?

Well in Chapter 5 (and elsewhere)  
we are going to need the mean  
and variance of a Binomial  
Random Variable

remember this was like tossing an unfair coin where the probability of "success" (heads) was  $P$  and probability of failure (tails) was  $1-p$ , many times.

(Remember Binomial calculator from exam.)

The Binomial RV is the count of ~~that~~ successes. Well I could just give you the formulas for mean and variance.

But they would be very bizarre. It wouldn't make sense why they were true.

With the formulas shown today I can actually explain why these formulas are true.

They will be similarly helpful in other situations.

Lets derive the mean and variance of the random variable "number of heads in three tosses of a coin" that we have seen earlier

We have already derived them

$\mu_X = 3/2$   $\sigma_X^2 = 3/4$  But it was hard and our method will get harder. Now we will use a method with large values of  $n$ .

Consider that  $X$  is the sum of three independent RVs each one the number of heads in the toss of one coin,

$$X = Y_1 + Y_2 + Y_3$$

Values of $Y_i$	0	1
Probability	$\frac{1}{2}$	$\frac{1}{2}$

Mean  $0 \cdot \frac{1}{2} + 1 \cdot \frac{1}{2} = \frac{1}{2}$

$$\begin{aligned}\sigma_{Y_i}^2 &= (0 - \frac{1}{2})^2 \cdot \frac{1}{2} + (1 - \frac{1}{2})^2 \cdot \frac{1}{2} \\ &= \frac{1}{4} \cdot \frac{1}{2} + \frac{1}{4} \cdot \frac{1}{2} = \frac{1}{4}\end{aligned}$$

Now we apply our rules

$$\mu_X = \cancel{\mu_{Y_1}} + \cancel{\mu_{Y_2}} + \cancel{\mu_{Y_3}} = \frac{1}{2} + \frac{1}{2} + \frac{1}{2} = \frac{3}{2}$$

check ✓

Because independent

$$\sigma_X^2 = \sigma_{Y_1}^2 + \sigma_{Y_2}^2 + \sigma_{Y_3}^2 = \frac{1}{4} + \frac{1}{4} + \frac{1}{4} = \frac{3}{4}$$

check ✓

Now in general

Values of $Y_i$	0	1
Probability	$1-p$	$p$

$$\mu_Y = 0 \cdot (1-p) + 1 \cdot p = p$$

$$\sigma_{Y_i}^2 = \cancel{(0-p)^2(1-p)} + (1-p)^2 \cdot p$$

$$= p^2(1-p) + (1-p)^2 \cdot p$$

$$= p(1-p)[p + (1-p)]$$

$$= \cancel{p(1-p)}$$

Because:  $\mu_X = \mu_{Y_1} + \dots + \mu_{Y_n} = np$

Independent  $\sigma_X^2 = \sigma_{Y_1}^2 + \dots + \sigma_{Y_n}^2 = np(1-p)$

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

6

Notes about Correlation between  $x$  and  $y$

If  $\rho = 1$   $x$  and  $y$  are related linearly with positive slope

$$Y = mx + b \quad m > 0$$



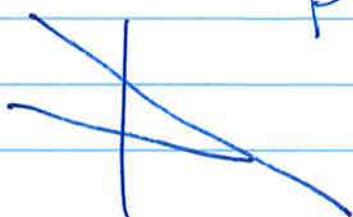
Data  
look like  
this

$$\rho = 1$$

$$\begin{aligned} 3Y &= 4 + 5X \\ X - Y &= 2 \end{aligned} \quad \left. \begin{array}{l} \text{both have} \\ \text{correlation} \\ \rho = 1 \end{array} \right\}$$

Likewise  $Y = mx + b \quad m < 0$

$$\rho = -1$$

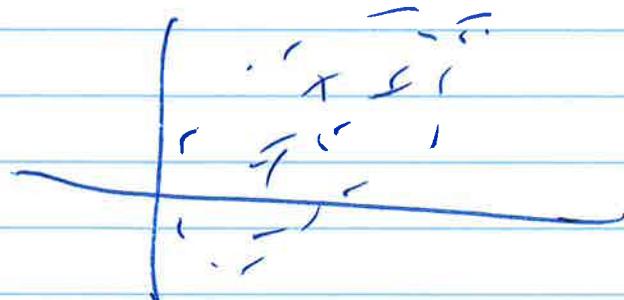


$$\text{e.g., } Y + X = 3$$

7

IF  $0 < |\rho| < 1$  there is scatter between  $X$  and  $Y$

When values are simulated get a scatter plot with



$r = \rho$  approximately

If  $\rho = 0$  there is no linear relationship between  $X$  and  $Y$

If  $X$  and  $Y$  are independent

$$\rho = 0$$

If  $\rho = 0$   $X$  and  $Y$  ~~are not necessarily~~ could be but are not necessarily independent.

Independent mean knowing  $X$  tells you nothing about  $Y$  and vice-versa

## New Sampling Distributions

Statistics such as means and proportions summarize data

- mean height of 30 people drawn from a population
- the proportion of 30 people drawn from a population who prefer red over blue,

A statistic from a random sample or randomized experiment is a random variable

The probability distribution of such a statistic is its sampling distribution

The population distribution of a variable is the distribution of its values over all members of the population. Also the prob distrib of the variable when we choose one individual at random from the population

## 55.1 The Sampling distribution of the sample mean

Random phenomenon: Draw a sample  
of size  $n$  from a population

Random Variable: mean of  $n$  observations  
(e.g. heights) from sample,

called Sample mean

Facts about Sample means

1. Sample means are less variable than individual observations
2. Sample means are more normal than individual observations

In other words the sampling distribution  
of a sample mean has less spread  
than the population distribution and  
is closer to Normal — Q-Q Plot  
looks more like a line ~~less like a bell~~  
~~unless it's a bell~~ (unless it's a line)

(Pg 5)

The mean and standard deviation of  $\bar{X}$

The sample mean is approximately the mean of the underlying distribution

But every time you draw a sample you get a different value (maybe only slightly different for  $\bar{X}$ )

In other words, there is spread in the sampling distribution of  $\bar{X}$

$\bar{X}$  is an estimate of the population mean  $\mu$ .

"Estimate" is a technical term.

What is the mean of  $\bar{X}$  and what is the variance / Standard deviation?

Let  $x_i$  be the observation from the  $i^{\text{th}}$  element of the sample. Then  $\bar{X} = \frac{1}{n}(x_1 + x_2 + \dots + x_n)$

From Rules  $\mu_{\bar{X}} = \frac{1}{n}(\mu_{x_1} + \mu_{x_2} + \dots + \mu_{x_n})$

But  $\bar{X}_i$  (mean of one observation from one individual in sample)

Must be the same as  ~~$\mu_{\text{population}}$~~   $\mu$   
 mean of population — they ~~are~~ two names for the same thing.

$$\begin{aligned}\text{Thus } \bar{X} &= \frac{1}{n} (\bar{X}_1 + \dots + \bar{X}_n) \\ &= \frac{1}{n} (\mu + \mu + \mu + \dots + \mu) \\ &= \frac{1}{n} \cdot n \mu \\ &= \mu\end{aligned}$$

The mean of the sampling distribution  $\bar{X}$  is the same as the mean of the population.

PJ7

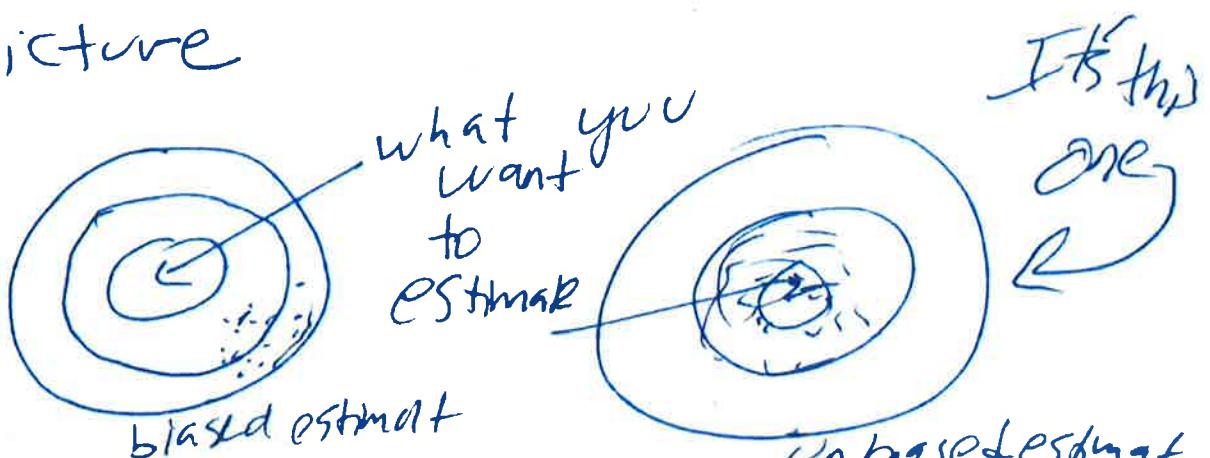
Because  $\mu_{\bar{X}} = \mu$  and

~~Average~~  $\bar{X}$  estimates  $\mu$

We say  $\bar{X}$  is an unbiased estimate of  $\mu$  (the mean of the estimate happens to be what you are trying to estimate)

Picture

TWO alternatives:



Bull's Eye is what you want to estimate. Cloud of points are different sample means each calculated from  $n$  observations (independent) of  $X$ .

When you say an estimate is unbiased (that's good) but it doesn't tell you how close to the bull's eye the points lie — it tells you that the center (mean) of the sampling distribution is spot on but it doesn't tell you about the spread in distribution — whether all darts fell close to bulls eye or whether they were all over the board,

The variance tells you this. Because the observations are independent we can use the rules for variances

$$\sigma_x^2 = \left(\frac{1}{n}\right)^2 (\sigma_{x_1}^2 + \sigma_{x_2}^2 + \dots + \sigma_{x_n}^2)$$

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

$\sigma^2$   
population variance

$$= \frac{\sigma^2}{n}$$

Standard deviation  $\frac{\sigma}{\sqrt{n}}$

## In Summary

Let  $\bar{X}$  be the mean of an SRS of size  $n$  from a population having mean  $\mu$  and standard deviation  $\sigma$ .

The mean and standard deviation of  $\bar{X}$  are

$$\mu_{\bar{X}} = \mu$$

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

$\uparrow$   
 $n > 1$  ~~smaller~~

so the standard deviation of the sampling distribution of the sample means is smaller than the standard deviation of the population

Remember we said sample means are less variable than individual observations

This shows why

If the population distribution is normal  $N(\mu, \sigma)$  then the sampling distribution of  $\bar{X}$  is also normal  $N(\mu, \frac{\sigma}{\sqrt{n}})$ .

Even if the population distribution is not normal, the sampling distribution of  $\bar{X}$  is approximately normal  $N(\mu, \frac{\sigma}{\sqrt{n}})$  for large  $n$ .

(Technically this is only true if  $\sigma$  is finite. I'll bet you didn't know  $\sigma$  could be infinite. The book mentions this.)

→ This is called The Central Limit Theorem. It is one of the two biggest results in statistics along with the ~~contradicting~~ law of large numbers.

Note that even if RV  $X$  is discrete (1 if heads, 0 if tails)

$\bar{X}$  is approximately normal  
~~(proportion of heads in n tosses)~~