

# **AP Stats Review**

**DRAFT – Not Final**

Rafael A. Irizarry

April 30, 2025

# Table of contents

<b>Preface</b>	<b>4</b>
<b>1 Introduction</b>	<b>5</b>
1.1 Population parameters . . . . .	5
1.2 Random variables . . . . .	5
1.3 Sample estimates . . . . .	6
<b>2 One variable data</b>	<b>8</b>
2.1 The normal distribution . . . . .	8
<b>3 Random Variables and Probability Distributions</b>	<b>10</b>
3.1 Discrete random variables . . . . .	10
3.2 Mean and standard deviation of a random variable . . . . .	12
3.3 Bernoulli trials . . . . .	13
3.4 Combining, shifting, and scaling random variables . . . . .	14
3.5 Binomial distribution . . . . .	15
3.6 Geometric distribution . . . . .	16
3.7 Continuous random variables . . . . .	17
3.8 Approximation to Normal . . . . .	18
<b>4 Sampling distributions</b>	<b>20</b>
4.1 Population and parameters . . . . .	20
4.2 The sample average . . . . .	21
4.3 Central Limit Theorem . . . . .	22
4.4 Proportions . . . . .	22
4.5 Means . . . . .	24
<b>5 Proportions</b>	<b>26</b>
5.1 Confidence interval . . . . .	26
5.2 Critical values . . . . .	27
5.3 p-values . . . . .	27
5.4 Difference of two proportion . . . . .	29
5.4.1 Confidence interval . . . . .	29
5.4.2 p-value . . . . .	29
5.4.3 Confidence interval and p-value connection . . . . .	30

<b>6</b>	<b>Means</b>	<b>31</b>
6.1	Confidence interval . . . . .	31
6.2	t-test . . . . .	31
6.3	Difference of two means . . . . .	32
6.3.1	Confidence interval . . . . .	32
6.3.2	p-value . . . . .	33
<b>7</b>	<b>Goodness of fit</b>	<b>34</b>
<b>8</b>	<b>Slopes</b>	<b>37</b>
8.1	Motivation to help understand . . . . .	37
8.1.1	Example is height hereditary? . . . . .	37
8.1.2	The population correlation coefficient . . . . .	37
8.1.3	Sample correlation coefficient . . . . .	38
8.1.4	Correlation is not always a useful summary . . . . .	39
8.1.5	The regression line . . . . .	39
8.1.6	Linear model representation . . . . .	40
8.2	Notes for the exam . . . . .	41

# Preface

These are draft notes intended to support review for AP Statistics. They are a work in progress and may change frequently as material is added or revised. They have not been fully proofread and are provided *as is*, without any guarantees regarding completeness or accuracy.

# 1 Introduction

Main notation, definitions, and properties to remember.

## 1.1 Population parameters

For problems about proportions

$p$  is the population proportion of 1s,  $1 - p$  is the proportion of 0s.

$\sqrt{p(1-p)}$  is the standard deviation of a population with a proportion  $p$  of 1s.

For problems about continuous variables

$\mu$  is the population average.

$\sigma$  is the population standard deviation

For categorical variables

$p_1, p_2, \dots, p_k$  are the proportion of categories  $1, 2, \dots, k$  respectively.

## 1.2 Random variables

All random variables  $X$  have a distribution.

- $\mu_X$  is the mean of this distribution
- $\sigma_X$  is the standard deviation of this distribution
- $\mu_{X+Y} = \mu_X + \mu_Y$
- $\sigma_{X+Y} = \sqrt{\sigma_X^2 + \sigma_Y^2}$  **IF  $X$  and  $Y$  are independent.**

## 1.3 Sample estimates

For problems about proportions

- $\hat{p}$  is the sample population
- $\hat{p}$  estimates  $p$
- $N$  is the sample size
- $\mu_{\hat{p}} = p$
- $\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{N}}$
- If  $N \geq 30$

$$\frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{N}}}$$

follows standard normal distribution

For problems about continuous variables

- $\bar{X}$  is the sample average
- $\bar{X}$  estimates  $\mu$
- $s_X$  estimates  $\sigma$
- $N$  is the sample size
- $\mu_{\bar{X}} = \mu$
- $\sigma_{\bar{X}} = \frac{\sqrt{\sigma}}{N}$
- If  $N \geq 30$

$$\frac{\bar{X} - \mu}{\sqrt{\frac{s_X}{N}}}$$

follows standard normal distribution

- If  $N \leq 30$  and **population data** is approximately normal then it follows a  $t$  distribution with  $N - 1$  degrees of freedom.

For categorical variables

- We define  $E_i = Np_i$  as the expected counts for each category if we take a sample of size  $N$ .
- If we observe counts  $O_1, \dots, O_k$  then

$$\sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}$$

follow  $\chi^2$  with  $k - 1$  degrees of freedom.

Regression line

- $b$  is the slope of the least squares regression line. It is a random variable!
- $\sigma_b$  is the standard deviation of the regression line. It usually given in a table under **SE Coef** column.
- Under the null that there is no relationship and that residuals are approximately normal  $b/\sigma_b$  follows t distribution with  $N - 2$  degrees of freedom.
- The residuals have average 0.

**i** Note

Note that all the test statistics have the form

$$\frac{\text{random variable} - \text{mean of random variable}}{\text{standard deviation of random variable}}$$

In the case of the  $\chi^2$  it is squared because we both negative and positive deviations from expected to make the statistic bigger.

$$\frac{(\text{random variable} - \text{mean of random variable})^2}{\text{variance of random variable}}$$

## 2 One variable data

### 2.1 The normal distribution

The normal distribution, also known as the bell curve and as the Gaussian distribution, is one of the most famous mathematical concepts in history.

One reason for this is that approximately normal distributions occur in many situations, including gambling winnings, heights, weights, blood pressure, standardized test scores, and experimental measurement errors.

This is not for discrete random variables but **continuous random variables**. If  $X$  has a normal pdf it can take on any continuous value. Therefore  $Pr(X = x)$  does not make sense anymore.

The normal distribution is defined with a mathematical formula. For any interval  $(a, b)$ , the proportion of values in that interval can be computed using this formula:

$$Pr(a < x \leq b) = \int_a^b \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx$$

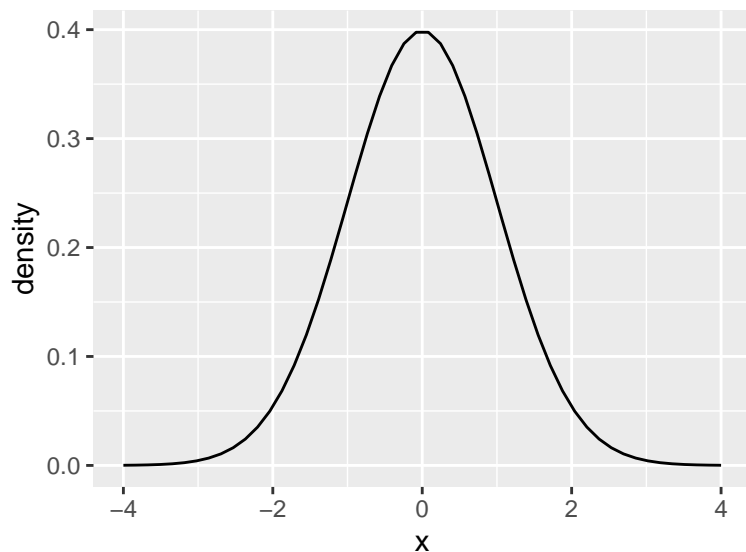
You don't need to memorize the formula.

The most important characteristics is that it is completely defined by just two parameters:  $\mu$  and  $\sigma$ . The rest of the symbols in the formula represent the interval ends,  $a$  and  $b$ , and known mathematical constants  $\pi$  and  $e$ .

These two parameters,  $\mu$  and  $\sigma$ , are referred to as the *mean* and the *standard deviation* (SD) of the distribution, respectively.

The distribution is symmetric, centered at  $\mu$ , and most values (about 95%) are within  $2\sigma$  from  $\mu$ . Here is what the normal distribution looks like when the  $\mu = 0$  and  $\sigma = 1$ :





If  $\mu = 0$  and  $\sigma = 1$  it is called **standard**.

Memorize the following: For a standard normal random variable  $Z$ :

- 68% are between -1 and 1
- 95% are between -2 and 2 (-1.96 and 1.96 to be precise)
- 99.7% between -3 and 3

You can use the `normalcdf` function to obtain these probabilities

- $P(a < Z < b)$  is `normalcdf(a,b,0,1)`, 0 and 1 are mean and SD, respectively
- $P(Z < a)$  is `normalcdf(-1E99,a,0,1)`. -1E99 is  $-\infty$

Also memorize:

- If  $Z$  is normal and  $a$  is a constant  $aZ$  is normal.
- If  $Z$  is normal and  $b$  is a constant  $Z + b$  is normal.
- If  $Z$  is normal and  $a, b$  are a constant  $aZ + b$  is normal.
- If  $Z$  and  $W$  are normal and  $Z + W$  is normal.

## 3 Random Variables and Probability Distributions

### 3.1 Discrete random variables

Discrete random variable

$X$  can have different outcomes, each one with a probability. Discrete means the possible outcomes are finite.

Probability density function (pdf)

Defines the probability  $P(X = k)$  for each outcome  $k$   
The exam often uses short hand  $p_k = P(X = k)$

Note

Later we will need the *cumulative distribution function* which is simply defined by

$$\begin{aligned} F(a) &= P(X \leq a) \\ &= \sum_{x_i \leq a} P(X = x_i) \end{aligned}$$

Interpretation  $F(a)$  tells us the probability of  $X$  being less than  $a$  for any  $a$ .

Example 1: Fair coin

Define  $X = 0$  for tails and  $X = 1$  for heads

$$P(X = 0) = 1/2$$

$$P(X = 1) = 1/2$$

### Example 2: A die

$$P(X = 1) = 1/6$$

$$P(X = 2) = 1/6$$

$$\vdots$$

$$P(X = 6) = 1/6$$

### Example 3: Sum of two dice

$$P(X = 2) = 1/36,$$

$$P(X = 3) = 2/36$$

$$\vdots$$

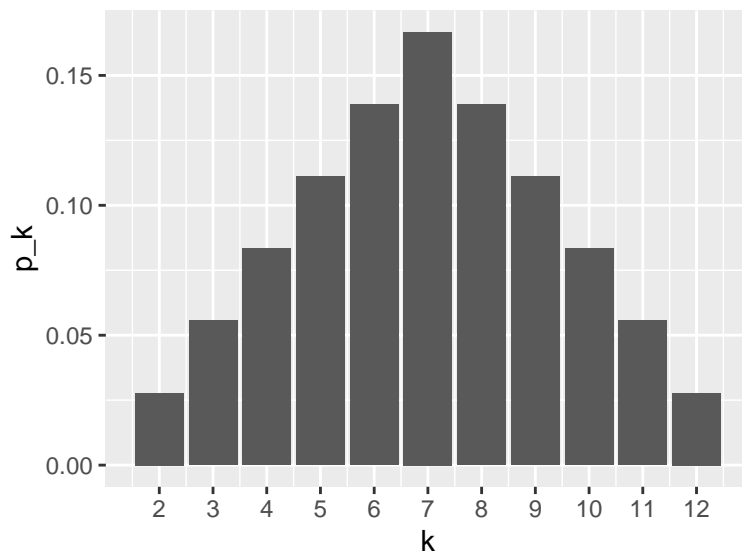
$$P(X = 7) = 1/6$$

$$\vdots$$

$$P(X = 12) = 1/36$$

The pdf is often shown as a table or a graph. For the graph we simply plot a bar for each  $k$  going up to  $p_k$ .

Here is the pdf for the sum of two dice:



## 3.2 Mean and standard deviation of a random variable

The mean and standard deviation of the the distribution of a random variable  $X$  are referred as the mean and standard deviation of  $X$ .

Mean

$$\mu_X = \sum_{i=1}^n x_i P(X = x_i)$$

With

$$x_1, \dots, x_n$$

all the possible outcomes.

For any random variable  $*$  we use the symbol  $\mu_*$  to represent it's mean. It can be any random variable and we don't always use the name  $X$ .

The standard deviation (SD)

The standard de deviation of the distribution of a random variable  $X$  is defined as:

$$\sigma_X = \sqrt{\sum_{i=1}^n (x_i - \mu_X)^2 P(X = x_i)}$$

We use  $\sigma_*$  just like  $\mu_*$

You can think of this as the typical distance you see  $X$  from the mean  $\mu_X$ . For the heads or tails this is 1/2 since both 0 and 1 are 1/2 from the mean 1/2.

**i** Note

We sometimes say *the standard error of  $X$*  to mean *the standard deviation of the distribution of a random variable  $X$* .

The variance

The variance is simply the standard deviation squared  $\sigma_X^2$ .

- We define the variance because mathematical calculations are easier without the square root. Other than that we always use standard deviation.
- The standard deviation has the same units as  $X$ . The variance has units squared which has no interpretation. What is kilograms squared or dollars squared?

### 3.3 Bernoulli trials

- A super common example of a useful random variable are Bernoulli trials.
- These are either a 0 (failure) or 1 (success) and each trial is independent of others.

#### Bernoulli trial definition

- $X$  is either 1 (success) or failure (0).
- Completely defined by the probability of success:  $P(X = 1) = p$ .
- The probability of failure is simply  $1 - p$ , sometimes called  $q$ .

Bernoulli trials are popular because we can use them to count random things: number of heads when we toss coins, number of lottery winners, number of defective light bulbs made in a day by a factory, number of patients that got cured by a drug, number of COVID-19 hospitalizations in a day, and so on.

#### Bernoulli trial mean and SD

If  $X$  is a Bernoulli trial:

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

You need to memorize this but here is the derivation

$$\begin{aligned}\mu_X &= 0 \times P(X = 0) + 1 \times P(X = 1) \\ &= p\end{aligned}$$

and

$$\begin{aligned}\sigma_X^2 &= (0 - p)^2(1 - p) + (1 - p)^2p \\ &= (1 - p)p(p + 1 - p) \\ &= p(1 - p)\end{aligned}$$

#### Examples

- Tossing coins,  $p = 0.5$
- Steph Curry free throws,  $p = 0.9$
- Lottery winners,  $p < 10^{-6}$
- Celtics win a game in NBA finals  $p = ?$

### 3.4 Combining, shifting, and scaling random variables

#### Mean of linear combinations

Need to memorize these (they are intuitive). If  $X$  and  $Y$  random variables and  $a$  is a constant:

- $\mu_{X+Y} = \mu_X + \mu_Y$
- $\mu_{X+a} = \mu_X + a$
- $\mu_{aX} = a\mu_X$

#### Example

If  $X$  and  $Y$  are two random variables, what is  $\mu_{X-Y}$ ?

$$\begin{aligned}\mu_{X-Y} &= \mu_X + \mu_{-Y} \\ &= \mu_X + -1\mu_Y \\ &= \mu_X - \mu_Y\end{aligned}$$

#### SD of linear combinations

For these we use the variance. But you can take square root at the end.

- $\sigma_{X+a}^2 = \sigma_X^2$ : shifting does not change variability.
- $\sigma_{aX}^2 = a^2\sigma_X^2 \implies \sigma_{aX} = |a|\sigma_X$ : change of scale also scales measure of variability.
- **If  $X$  and  $Y$  are independent**,  $\sigma_{X+Y}^2 = \sigma_X^2 + \sigma_Y^2 \implies \sigma_{X+Y} = \sqrt{\sigma_X^2 + \sigma_Y^2}$ : Adding two things that vary, varies more.

#### Example

If  $X$  and  $Y$  are two **independent** random variables, what is  $\sigma_{X-Y}$ ?

$$\begin{aligned}\sigma_{X-Y}^2 &= \sigma_X^2 + \sigma_{-Y}^2 \\ &= \sigma_X^2 + (-1)^2\sigma_Y^2 \\ &= \sigma_X^2 + \sigma_Y^2\end{aligned}$$

Which implies

$$\sigma_{X-Y} = \sqrt{\sigma_X^2 + \sigma_Y^2}$$

Interpretation: Subtracting two variables that vary independently has more variability than each.

### 3.5 Binomial distribution

Another popular random variable is the sum of Bernoulli trials.

$$S = \sum_{i=1}^n X_i$$

It tells us the number of successes and it is also a random variable.

Examples:

- Number of heads if I toss coins
- Number of free throws curry makes

Example

What is  $\mu_S$ ?

$$\begin{aligned}\mu_S &= \mu_{X_1 + \dots + X_n} \\ &= \mu_{X_1} + \dots + \mu_{X_n} \\ &= np\end{aligned}$$

What is  $\sigma_S$ ?

$$\begin{aligned}\sigma_S^2 &= \sigma_{X_1 + \dots + X_n}^2 \\ &= \sigma_{X_1}^2 + \dots + \sigma_{X_n}^2 \\ &= np(1-p)\end{aligned}$$

This implies

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

Binomial pdf

We can compute the pdf for the sum of  $n$  trials:

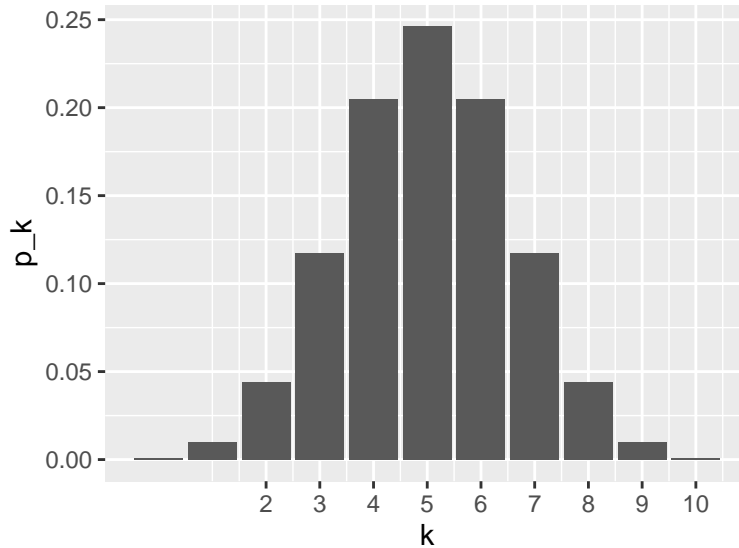
$$P(S = k) = \binom{n}{k} p^k (1-p)^{n-k}$$

This is called the binomial distribution and can be computed in AP-test with `binompdf(n, p, k)` and the CDF with `binomcdf(n, p, k)`

The CDF is useful for answering questions such as “what is the chance that we see 3 heads or less?” or “what is the chance we see 4,5,6 heads?”

#### Example

pdf of the number of heads when tossing 10 coins:



### 3.6 Geometric distribution

It is also common to ask how many trials do I need to see a success. For example, how many free throws will Curry take until he misses.

#### Geometric distribution

Define a random variable as  $X$ =number of trials if we stop after the first success.  
It is not hard to see that this is:

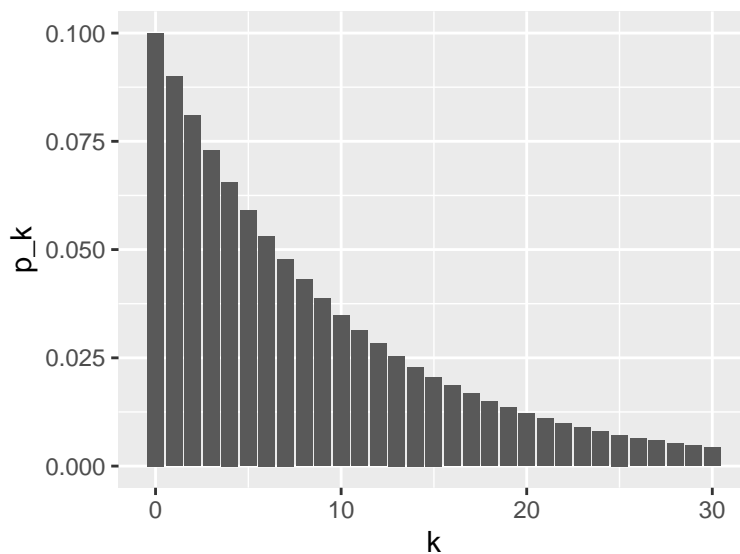
$$P(X = k) = (1 - p)^{k-1}p$$

This is called the Geometric distribution, defined for  $k = 1, 2, \dots, \infty$ .

Example: number of free throws before Curry misses.

Here miss is the success we are waiting for so  $p = 0.1$





We use this to calculate, for example, that the chance of seeing 10 or more free throws in a row to start the game is  $1 - \text{geomcdf}(10, .1) = 0.3138106$

### 3.7 Continuous random variables

#### Continuous random variables

Some random variables are continuous. Height, weight, and temperature are examples.

#### CDF

A continuous random variable  $X$  can take an infinity number of values  $x$  so it does not make sense to write:

$$P(X = x)$$

Instead we define the cumulative distribution function as

$$F(a) = P(X < a)$$

We can then define the *probability density function*  $f(x)$  so that

$$F(a) = \int_{-\infty}^a f(x) dx$$

**!!!** We use continuous distribution to approximate discrete ones. We will use

- normal distribution
- t distribution
- Chi-square distribution

In the test you either use a function on the calculator or they provide a look up table for the cdf  $F(a)$ .

:::

## 3.8 Approximation to Normal

When the number of trials is large binomial is very well approximated by the normal distribution.

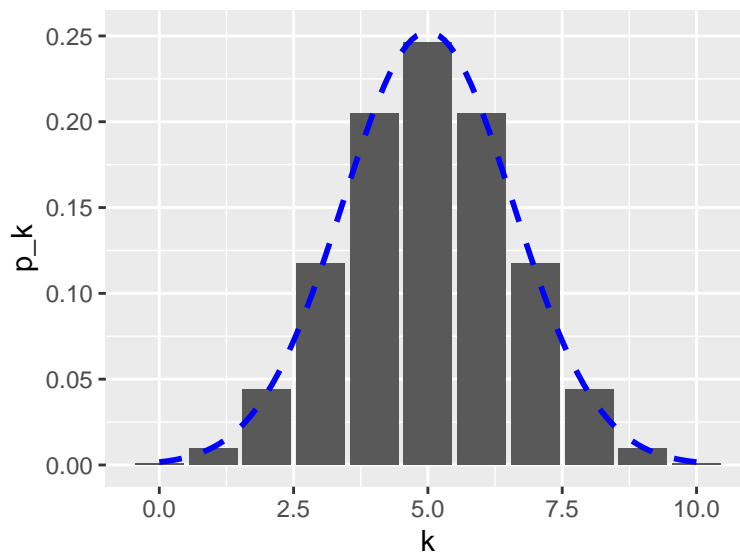
Define

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

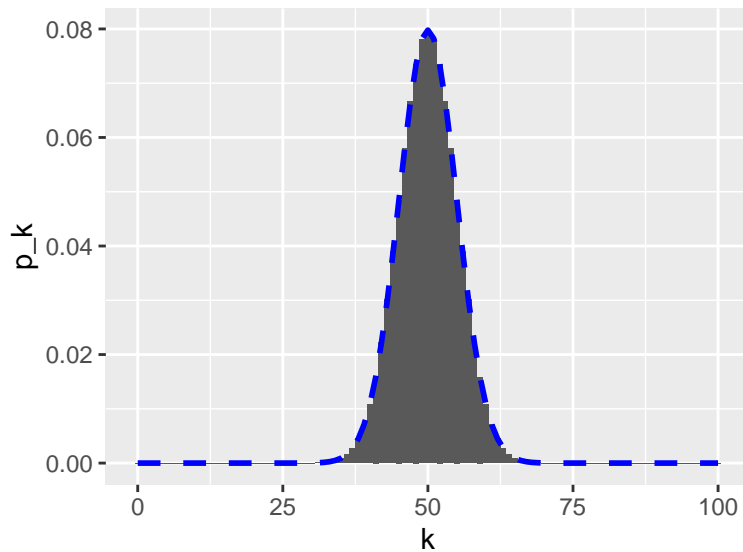
Then  $Z$  is approximated by standard normal

Here is a  $n = 10, p = 0.5$  binomial with a normal with mean  $np = 5$  and standard deviation  $\sqrt{np(1-p)} \approx 1.6$  added in blue

Warning: Using `size` aesthetic for lines was deprecated in ggplot2 3.4.0.  
i Please use `linewidth` instead.



Here it is for 100. In this case  $np = 50$  and  $\sqrt{np(1-p)} = 5$



#### Example

If I toss 100 coins, what is the probability that I see between 45 and 55 heads?

We can use the binomial to answer this exactly `binomcdf(100, 0.5, 55) - binomcdf(100, 0.5, 44)` which is 0.728747.

But we can also use the normal distribution:

$$\begin{aligned}
 P(45 \leq S \leq 55) &= P(44.5 < S < 55.5) \\
 &= P(44.5 - 50 < S - 50 < 55.5 - 50) \\
 &= P\left(\frac{44.5 - 50}{\sqrt{100 \times 0.5 \times 0.5}} < \frac{S - 50}{\sqrt{100 \times 0.5 \times 0.5}} < \frac{55.5 - 50}{\sqrt{100 \times 0.5 \times 0.5}}\right) \\
 &= P(-1.1 < Z < 1.1)
 \end{aligned}$$

We use `normalcdf(-1.1 1.1, 0, 1) = 0.7286679`

Which is almost identical to the binomial result.

#### **i** Note

Important to understand why we do the first  $P(45 \leq S \leq 55) = P(44.5 < S < 55.5)$ . The normal distribution is continuous so it can't be equal to anything. So we do the adjustment to make sure we include 45 and 55 in the approximation.

## 4 Sampling distributions

We want to learn about populations from samples.

### 4.1 Population and parameters

#### Population

The population is defined by the list of numbers  $x_1, x_2, \dots, x_n$ .

The  $x$ s are **not random**.

We can't see the entire population but we want to learn about it.

#### Example 1: Trump voters

The population is the people who will vote on election day. Trump voters get a 1 and others get a 0. So all the  $x$ s are either 0 or 1.

#### Example 2: High school SAT scores

The population are all the students that took SAT. The  $x$ s are the scores for each student.

#### Population parameters

The population parameters are summaries of the  $x_1, x_2, \dots, x_n$  we are interested in.

In the AP test we almost always care about the **population mean**:

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

Most of this chapter is about *estimating* the population mean  $\mu$ .

Another parameter we will need is the **population standard deviation**:

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}$$

### Population proportion

When the  $x$ s are 0s or 1s, then the population mean is equivalent to the proportion of 1s.

In this case we use the symbol  $p$  instead of  $\mu$ .

And the standard deviation can be shown to be  $\sqrt{p(1-p)}$

## 4.2 The sample average

The strategy to estimate the population parameter is to take a random sample: we can't examine all the sample so we examine a much smaller subset.

We learn that we can learn a lot about population parameters from samples.

By far the most common example is using a *sample average* to estimate a *population mean*.

### A sample

- A sample are the resulting observed values we obtain when picking individuals at random from the population.
- We represent them with capital letters because they are **random variables**:

$$X_1, \dots, X_N$$

### The sample size

$N$  is called the sample size.

Do not confuse it with the number of individuals in the population  $n$ .

In the election poll example  $n$  is over 100 million while a typical sample size  $N$  is 1,000 or less.

### The sample average

- A sample are the resulting observed values we obtain when picking individuals at random from the population.
- We represent them with capital letters because they are **random variables**:

$$X_1, \dots, X_N$$

\*  $N$  is called the sample size.

## 4.3 Central Limit Theorem

tldr: The distribution of the sample average is approximated by a normal distribution when the sample size is large.

### Central Limit Theorem (CLT)

- If  $X_1, \dots, X_N$  are random variables that are independent and have the same distribution, the sum  $\sum_{i=1}^N X_i$  gets closer and closer to being normally distributed when  $N$  gets very large.
- Because dividing a normal random variable by a constant is still normal, the CLT applies to the average  $\frac{1}{N} \sum_{i=1}^N X_i$  as well.
- **Rule of thumb**  $N \geq 30$  is considered large enough.

## 4.4 Proportions

A very common application of statistics is estimating a population proportion

Examples:

- Proportion of voters voting for trump.
- Proportion of patients that a drug cures.
- Proportion of adults with a job.

We want to estimate  $p$ , the population parameter.

Note that:

- Each  $X$  in the sample is a Bernoulli trial because  $P(X = 1) = p$ .
- This implies that for all  $i$ ,  $\mu_{X_i} = p$  and  $\sigma_{X_i} = \sqrt{p(1-p)}$
- Because we sample with replacement the  $X$ s are independent.

### Mean and SD of sample proportion

The sample proportion is

$$\hat{p} = \frac{1}{N} \sum_{i=1}^N X_i$$

Using what we have learned about mean and SD of combinations and re-scaling we have:

- $\mu_{\hat{p}} = p$
- $\sigma_{\hat{p}} = \frac{\sigma_X}{\sqrt{N}} = \frac{\sqrt{p(1-p)}}{\sqrt{N}}$

#### Distribution of sample proportion

- $\hat{p}$  is a sum of Bernoulli trials divided by a constant. So we could use the Binomial distribution to compute  $P(\hat{p} = k/N)$
- However, in the exam they want you to use the CLT.
- $\hat{p}$  is approximated by normal distribution with mean  $p$  and SD  $\sqrt{p(1-p)/n}$

#### Example

If I take a poll of 1000 people to get an idea of how many people are voting for Trump, what is the chance that my sample proportion  $\hat{p} = 0.45$  is within 1% of the actual proportion?

We are asking  $P(|\hat{p} - p| < 0.01)$

Let's figure it out:

$$\begin{aligned} P(|\hat{p} - p| < 0.01) &= P(-0.01 < \hat{p} - p < 0.01) \\ &= P\left(\frac{-0.01}{\sqrt{\frac{p(1-p)}{N}}} < \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{N}}} < \frac{0.01}{\sqrt{\frac{p(1-p)}{N}}}\right) \\ &= P\left(\sqrt{1000} \frac{-0.01}{\sqrt{p(1-p)}} < Z < \sqrt{1000} \frac{0.01}{\sqrt{p(1-p)}}\right) \end{aligned}$$

I don't know  $p$  but in the exam they want you to stick in  $\hat{p}$  for the SD calculation. So  $\sqrt{1000}/(0.45 \times .55) \approx 63.5$

so we have  $P(0.645 < Z < 0.645)$  or `normcdf(-0.645, 0.645, 0, 1)` which is `pnorm(0.645)-pnorm(-0.645)`

#### 💡 Tip

In the exam compute the standard deviation  $\sqrt{p(1-p)/N}$  first and stick that in the calculations instead of the formula.

## 4.5 Means

Another common application of statistics is estimating a population mean

Examples:

- What is the average SAT score in a high school?
- What is the average blood pressure for people taking a drug?

We want to estimate  $\mu$ , the population parameter.

Note that

- Each  $X$  in the sample has the same distribution  $P(X = x_i) = 1/n$  for all  $i$ .
- This implies that for all  $i$ ,  $\mu_{X_i} = \mu$  and  $\sigma_{X_i} = \sigma$
- Because we sample with replacement the  $X$ s are independent.

### Mean and SD of sample average

The sample average is

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$$

Using what we have learned about mean and SD of combinations and re-scaling we have:

- $\mu_{\bar{X}} = \mu$
- $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{N}}$

### Distribution of sample average

CLT tells us that  $\bar{X}$  is approximated by a normal distribution with mean  $\mu$  and SD  $\sqrt{\sigma/n}$

### Sample standard deviation

If I want to make probability calculations I need to know  $\sigma_{\bar{X}}$ , but I don't know *sigma*.  
For proportions we used  $\sqrt{\hat{p}(1-\hat{p})}$  as an approximation of the standard deviation  $\sqrt{p(1-p)}$

But when sample means are not based on Bernoulli trials, we can't do that.  
Instead we use the sample standard deviation.



$$s = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (X - \bar{X})^2}$$

# 5 Proportions

## 5.1 Confidence interval

Task:

- We take a sample of 1,000 voters.
- 45% of our respondents say they will vote for trump.
- Provide an interval with 95% of containing the true proportion  $p$

Solution:

Our proportion  $\hat{p}$  follows a normal distribution with mean  $p$  and standard deviation  $\sqrt{p(1-p)/N} \approx \sqrt{0.45 \times 0.55/1000} \approx 0.016$

Lets consider symmetric intervals with:  $[\hat{p} - B, \hat{p} + B]$

We want to find a *margin of error* MOE such that:

$$P(p \in [\hat{p} - \text{MOE}, \hat{p} + \text{MOE}]) = 0.95$$

We accomplish it by setting  $\text{MOE} = 1.96\sigma_{\hat{p}} \approx 1.96\sqrt{\hat{p}(1-\hat{p})}/\sqrt{N} \approx 0.03$

$$\begin{aligned} P(p \in [\hat{p} - 2\sigma_{\hat{p}}, \hat{p} + 2\sigma_{\hat{p}}]) &= P(\hat{p} - 2\sigma_{\hat{p}} < p < \hat{p} + 2\sigma_{\hat{p}}) \\ &= P(-2\sigma_{\hat{p}} < \hat{p} - p < 2\sigma_{\hat{p}}) \\ &= P\left(-2 < \frac{\hat{p} - p}{\sigma_{\hat{p}}} < 2\right) \\ &= P(-2 < Z < 2) \\ &= 0.95 \end{aligned}$$

So our interval is  $0.45 \pm 0.03$ .

- If we want to be 99.7% sure we can use 3 instead of 2.
- If we want to be 68% sure we can use 1 instead of 2.

We refer to the 0.03 as the *margin of error* (MOE).

## 5.2 Critical values

We already knew that using 2 would give us a 95% confidence interval.

But what if we didn't know? Or if we wanted a 99% confidence interval?

The function `invNorm` will do this for us. The impute is the area to the left.

So to obtain 95% we need 0.5% to the left and 0.5% to the right.

We use `invNorm(0.995)` which gives us `r qnorm(0.995)`

So we multiply by 2.57 not 2 to get a 99% confidence interval.

Note that to get exactly 95% we actually use `invNorm(0.975)` which is `r qnorm(0.975)`, a little bit less than 2. In some books you will see 1.96 instead of 2.

## 5.3 p-values

- We want to know if a coin is biased.
- We toss it 100 times and observe 60% heads.

Is it biased or can this happen by chance?

Let's compute the probability of seeing  $\hat{p} = 0.6$  or more extreme.

Note that 0.4 is as extreme: we usually permit both directions.

Null hypothesis: It is fair or  $p = 0.5$

We will reject if the *p-value* is 0.05 or smaller.

The p-value is the probability observing something as extreme as we did when the null hypothesis holds

$$\begin{aligned} P(|\hat{p} - p| \geq 0.1) &= 1 - P(|\hat{p} - p| < 0.1) \\ &= 1 - P\left(\left|\frac{\hat{p} - p}{\sigma_{\hat{p}}}\right| < \frac{0.1}{\sigma_{\hat{p}}}\right) \\ &= 1 - P\left(|Z| < \frac{0.1}{\sigma_{\hat{p}}}\right) \end{aligned}$$

When the null hypothesis holds,  $\sigma_{\hat{p}} = \sqrt{0.5 \times 0.5/100} = 0.05$

So the p-value is  $1 - P(|Z| < 0.1/0.05 = 1.96)$  which is a bit less than 0.05

We reject the null hypothesis.

### Type of errors

- Type I error is rejecting the null hypothesis when it is true. Example say the coin is biased when it was fair.
- Type II error is failing to reject the null hypothesis when it is not true. Example saying the coin is fair when it was biased.
- Power is the 1 - probability of Type II error.

To help remember:



	Null hypothesis is TRUE	Null hypothesis is FALSE
Reject null hypothesis	Type I Error (False positive)	Correct outcome! (True positive)
Fail to reject null hypothesis	Correct outcome! (True negative)	Type II Error (False negative)

## 5.4 Difference of two proportion

### 5.4.1 Confidence interval

- Does drug work better than placebo?
- The proportion of the populations are  $p_1$  and  $p_2$ .
- For both placebo and drugged populations we obtain sample means  $\hat{p}_1$  and  $\hat{p}_2$
- The sample sizes are  $N_1$  and  $N_2$
- Provide a 95% confidence interval

We know the difference  $\hat{p}_1 - \hat{p}_2$  had the following mean and SD:

- $\mu_{\hat{p}_1 - \hat{p}_2} = p_1 - p_2$
- $\sigma_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{p_1(1-p_1)}{N_1} + \frac{p_2(1-p_2)}{N_2}}$

To construct a 95% confidence interval we use  $\hat{p}_1 - \hat{p}_2 \pm \sigma_{\hat{p}_1 - \hat{p}_2}$

As before we estimate

$$\sigma_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{N_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{N_2}}$$

### 5.4.2 p-value

- Suppose we have sample sizes of 25 and 100 for the the drug and placebo group respectively and
- we observe  $\hat{p}_1 = 0.25$  and  $\hat{p}_2 = 0.15$
- The null is that there is no difference so  $p_1 - p_2$  or  $p_1 = p_2 = p$

Under the null hypothesis we have

- $\mu_{\hat{p}_1 - \hat{p}_2} = 0$
- $\sigma_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{p(1-p)}{N_1} + \frac{p(1-p)}{N_2}}$

To compute the p-value we need an estimate for  $\sigma_{\hat{p}_1 - \hat{p}_2}$  which depends on  $p$

We estimate  $p$  with pooled data:

$$\frac{N_1 \hat{p}_1 + N_2 \hat{p}_2}{N_1 + N_2} = 0.17$$

which means  $\sigma_{\hat{p}_1 - \hat{p}_2} \approx \sqrt{(0.17 \times 0.83)(1/25 + 1/100)} \approx 0.08$

With this we can compute

$$\begin{aligned} P(|\hat{p}_1 - \hat{p}_2| \geq 0.1) &= 1 - P(|Z| < 0.1/\sigma_{\hat{p}_1 - \hat{p}_2}) \\ &= 1 - P(|Z| < 0.1/0.08) \\ &\approx 0.20 \end{aligned}$$

We do not reject.

### 5.4.3 Confidence interval and p-value connection

You can do the math and see that **if a 95% confidence interval does not include the null hypothesis mean, then a p-value will be less than 0.05**

The math:

If the null hypothesis says the mean is  $p$  and the observed  $\hat{p}$  resulted in a p-value less than 0.05, we know:

$$\left| \frac{\hat{p} - p}{\sigma_{\hat{p}}} \right| > 2$$

This implies that either

$$p > \hat{p} + 2\sigma_{\hat{p}} \text{ or } p < \hat{p} - 2\sigma_{\hat{p}}$$

## 6 Means

### 6.1 Confidence interval

Task:

- We take a sample of 36 student SAT scores.
- We observe a sample average of  $\bar{X} = 1100$  and a sample standard deviation  $s = 204$
- Provide an interval with 95% of containing the high school population average  $\mu$ .

Solution:

The sample average  $\bar{X}$  follows a normal distribution with mean  $\mu$  and standard deviation  $\sigma/\sqrt{N} \approx 204/6 = 34$

As with proportions we have

$$\begin{aligned} P(\mu \in [\bar{X} - 1.96\sigma_{\bar{X}}, \bar{X} + 1.96\sigma_{\bar{X}}]) &= P(\bar{X} - 1.96\sigma_{\bar{X}} < \mu < \bar{X} + 1.96\sigma_{\bar{X}}) \\ &= P(-1.96\sigma_{\bar{X}} < \bar{X} - \mu < 1.96\sigma_{\bar{X}}) \\ &= P\left(-1.96 < \frac{\bar{X} - \mu}{\sigma_{\bar{X}}} < 1.96\right) \\ &= P(-1.96 < Z < 1.96) \\ &= 0.95 \end{aligned}$$

So our interval is  $1100 \pm 68$

### 6.2 t-test

When  $N < 30$  we can't use CLT.

So what is the distribution of  $\bar{X}$ ?

If the population values are also approximately normal, as they are for SAT scores, then

$$t = \frac{\bar{X} - \mu}{s/\sqrt{N}}$$

Follows a t-distribution with  $N - 1$  degrees of freedom.

#### Example

Let's repeat the above example but  $N = 15$

All we have to do now is use the cutoff that gives us 0.95 for a t-distribution with 14 degrees of freedom.

You can use `invT` with area to the left and degrees of freedom.

So instead of 1.96 we use `invT(0.975, 14)` which is `qt(0.975, 14)`, a little bit bigger than 1.96.

We make our confidence interval

$$1100 \pm 2.14 \times 34$$

or

$$1100 \pm 73$$

## 6.3 Difference of two means

### 6.3.1 Confidence interval

- Are the mean SAT scores in two high schools different?
- The sample averages are  $\bar{X}_1 = 1200$  and  $\bar{X}_2 = 1100$  and the sample standard deviations are  $s_1 = 200$  and  $s_2 = 180$
- The sample sizes are  $N_1 = 30$  and  $N_2 = 35$
- Provide a 95% confidence interval

We know the difference  $\bar{X}_1 - \bar{X}_2$  had the following mean and SD:

- $\mu_{\bar{X}_1 - \bar{X}_2} = \mu_1 - \mu_2$
- $\sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{\sigma_1^2}{N_1} + \frac{\sigma_2^2}{N_2}}$

To construct a 95% confidence interval we use  $\bar{X}_1 - \bar{X}_2 \pm \sigma_{\bar{X}_1 - \bar{X}_2}$

We approximate  $\sigma_1$  and  $\sigma_2$  with  $s_1$  and  $s_2$



$$\sigma_{\bar{X}_1 - \bar{X}_2} \approx \sqrt{\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2}}$$

The sample sizes are large enough that we can use CLT so the confidence interval is

$$\bar{X}_1 - \bar{X}_2 \pm 2\sqrt{\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2}}$$

$$\bar{X}_1 - \bar{X}_2 = 100$$

and

$$2\sqrt{\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2}} = 2\sqrt{200^2/30 + 180^2/35} \approx 48$$

So the confidence interval is

$$100 \pm 48$$

### 6.3.2 p-value

Is the difference we saw significant?

We already computed  $\sigma_{\bar{X}_1 - \bar{X}_2} \approx 24$

$$P(|\bar{X}_1 - \bar{X}_2| > 100) = P(|Z| > 100/47) \approx 0.03$$

We reject.

## 7 Goodness of fit

Example:

Are all color Skittles equally likely?  
Here is the data:

Color	Observed Count
Red	20
Yellow	25
Green	15
Purple	18
Orange	22
<b>Total</b>	<b>100</b>

We see more yellow and less green?  
Can this happen by chance?

Chi-square test

This is a **goodness of fit** test.

If we have  $k$  categories each with  $p_i$ ,  $i = 1, \dots, k$  and observe  $N$  outcomes, we expect to see  $E_i = Np_i$  of each.

If we observe  $O_i$  for categories  $i = 1, \dots, k$  then the distribution of

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}$$

follows what is called a chi-square distribution with  $k - 1$  degrees of freedom.

Example:

Getting back to our example, if equally likely we have

$$E_i = 100/5 \approx 20$$

for all categories.

We can compute the  $\chi^2$  stat by

$$((20 - 20)^2 + (25 - 20)^2 + (15 - 20)^2 + (18 - 20)^2 + (22 - 20)^2) / 20 = 2.8$$

We can look up this probability for a  $\chi^2$  with 4 degrees of freedom and see that the p-value is 0.6.

So can easily see this by chance.

### Example

Is promotion status **independent** of gender, or is there evidence of **gender bias**?

Observed Data:

Gender	Promoted	Not Promoted	Total
Men	45	55	100
Women	30	70	100
<b>Total</b>	<b>75</b>	<b>125</b>	<b>200</b>

Calculate expected counts under the assumption of independence:

$$E_{ij} = \frac{(\text{row total}) \times (\text{column total})}{\text{grand total}}$$

Gender	Promoted	Not Promoted
Men	$\frac{100 \times 75}{200} = 37.5$	$\frac{100 \times 125}{200} = 62.5$
Women	37.5	62.5

Compute the chi-square test statistic:

$$\begin{aligned} \chi^2 &= \sum \frac{(O - E)^2}{E} = \frac{(45 - 37.5)^2}{37.5} + \frac{(55 - 62.5)^2}{62.5} + \frac{(30 - 37.5)^2}{37.5} + \frac{(70 - 62.5)^2}{62.5} \\ &= \frac{56.25}{37.5} + \frac{56.25}{62.5} + \frac{56.25}{37.5} + \frac{56.25}{62.5} = 1.5 + 0.9 + 1.5 + 0.9 = \mathbf{4.8} \end{aligned}$$

This has 1 degree of freedom.

Use the chi-square cumulative distribution function:

$$P(\chi^2 \geq 4.8) \approx 0.028$$

Since the p-value is approximately 0.028, which is less than the typical significance level of 0.05, we **reject**  $H_0$ .

### ! Important

In the test there will be two types of problems related to goodness of fit:

- Is the distribution of categorical data as expected (for example, Skittles)?
  - Several categories  $k$
  - probabilities of each category given
  - degrees of freedom is  $k - 1$
- Are two variables independent (for example, gender bias)?
  - Two categories
  - proportions calculated are from data.
  - degrees of freedom is 1.

## 8 Slopes

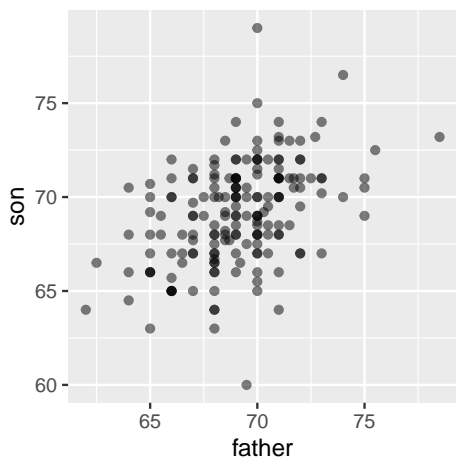
### 8.1 Motivation to help understand

This section gives some background and motivation to understand the material you are tested on. You can skip to the next section if you just want to see notes on the materials on the test.

#### 8.1.1 Example is height hereditary?

How well can we predict a child's height based on the parents' height?

We can summarize the data with the two averages and two standard deviations. However, this summary fails to describe an important characteristic of the data: the trend that the taller the father, the taller the son.



#### 8.1.2 The population correlation coefficient

The correlation coefficient is defined for a list of pairs  $(x_1, y_1), \dots, (x_n, y_n)$  as the average of the product of the standardized values:

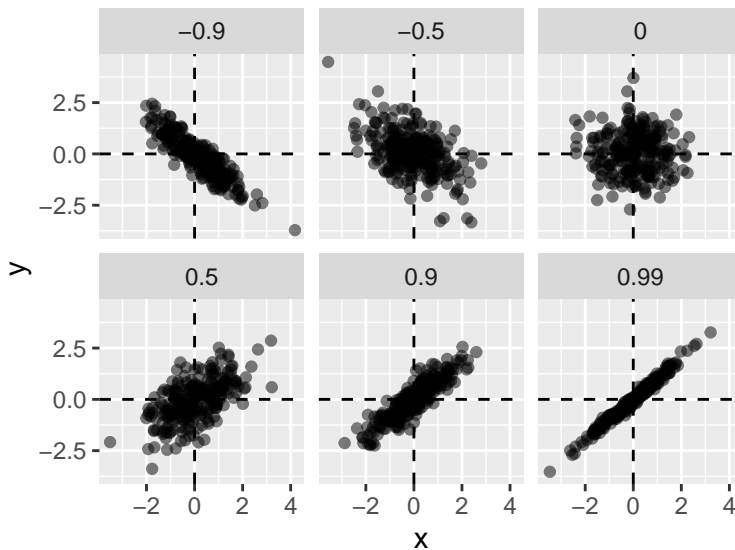
$$\rho = \frac{1}{n} \sum_{i=1}^n \left( \frac{x_i - \mu_x}{\sigma_x} \right) \left( \frac{y_i - \mu_y}{\sigma_y} \right)$$

with  $\mu_x, \mu_y$  the averages of  $x_1, \dots, x_n$  and  $y_1, \dots, y_n$ , respectively, and  $\sigma_x, \sigma_y$  the standard deviations.

To understand why this equation does in fact summarize how two variables move together, consider the  $i$ -th entry of  $x$  is  $\left( \frac{x_i - \mu_x}{\sigma_x} \right)$  SDs away from the average. Similarly, the  $y_i$  that is paired with  $x_i$ , is  $\left( \frac{y_i - \mu_y}{\sigma_y} \right)$  SDs away from the average  $y$ . If  $x$  and  $y$  are unrelated, the product  $\left( \frac{x_i - \mu_x}{\sigma_x} \right) \left( \frac{y_i - \mu_y}{\sigma_y} \right)$  will be positive ( $+\times+$  and  $-\times-$ ) as often as negative ( $+\times-$  and  $-\times+$ ) and will average out to about 0. This correlation is the average and therefore unrelated variables will have 0 correlation. If instead the quantities vary together, then we are averaging mostly positive products ( $+\times+$  and  $-\times-$ ) and we get a positive correlation. If they vary in opposite directions, we get a negative correlation.

The correlation coefficient is always between -1 and 1.

To see what data looks like for different values of  $\rho$ , here are six examples of pairs with correlations ranging from -0.9 to 0.99:



### 8.1.3 Sample correlation coefficient

The  $\rho$  defined above is for a population  $(x_1, y_1), \dots, (x_n, y_n)$

If we have a sample  $(X_1, Y_1), \dots, (X_N, Y_N)$  we can estimate with the sample correlation:

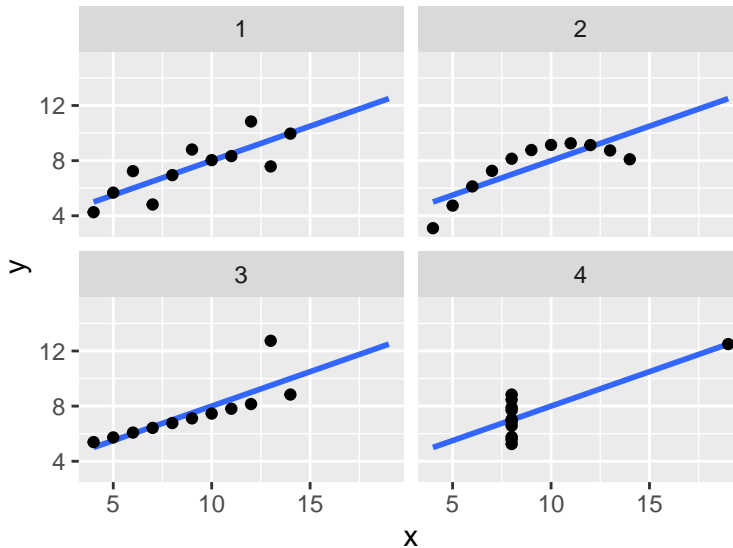
$$r = \frac{1}{N-1} \sum_{i=1}^N \left( \frac{X - \bar{X}}{s_X} \right) \left( \frac{Y - \bar{Y}}{s_Y} \right)$$

! Important

$r$  is a random variable!

### 8.1.4 Correlation is not always a useful summary

Correlation is not always a good summary of the relationship between two variables. The following four artificial datasets, referred to as Anscombe's quartet, famously illustrate this point. All these pairs have a correlation of 0.82:



### 8.1.5 The regression line

If we are predicting a random variable  $Y$  knowing the value of another  $X = x$  using a regression line, then we predict that for every standard deviation,  $\sigma_X$ , that  $x$  increases above the average  $\mu_X$ , our prediction  $\hat{Y}$  increase  $\rho$  standard deviations  $\sigma_Y$  above the average  $\mu_Y$  with  $\rho$  the correlation between  $X$  and  $Y$ . The formula for the regression is therefore:

$$\left( \frac{\hat{Y} - \mu_Y}{\sigma_Y} \right) = \rho \left( \frac{x - \mu_X}{\sigma_X} \right)$$

We can rewrite it like this:

$$\hat{Y} = \mu_Y + \rho \frac{\sigma_Y}{\sigma_X} (x - \mu_X)$$

- The  $\rho \frac{\sigma_Y}{\sigma_X}$  of the regression line is proportional to  $r$ .
- You can think of  $\frac{\sigma_Y}{\sigma_X}$  as needed for unit conversion.

**Interpretation:** If there is perfect correlation, the regression line predicts an increase that is the same number of SDs. If there is 0 correlation, then we don't use  $x$  at all for the prediction and simply predict the average  $\mu_Y$ . For values between 0 and 1, the prediction is somewhere in between. If the correlation is negative, we predict a reduction instead of an increase.

**i** Note

Note that if the correlation is positive and lower than 1, our prediction is closer, in standard units, to the average height than the value used to predict,  $x$ , is to the average of the  $x$ s. This is why we call it *regression*: the son regresses to the average height. In fact, the title of Galton's paper was: *Regression toward mediocrity in hereditary stature*.

### 8.1.6 Linear model representation

Note you can write the equations above in the form lines are usually written:  $y = a + bx$

In the case of the regression line  $Y$  is random, not exactly a line so we can write the relationship like this

$$Y = a + bx + \varepsilon$$

Here  $a + bx$  is the regression line and  $\varepsilon$  are the *errors* representing the distance between observed points and the line.

Note that using the equations above we can write

$$b = \rho \frac{\sigma_y}{\sigma_x}$$

which as the correlation after unit conversion.



## 8.2 Notes for the exam

### The regression line

In the AP test they write the regression line like this:

$$\hat{y} = a + bx$$

The little hat means on top of  $y$  signifies that it's a prediction.

The prediction  $\hat{y}$  is on a line, but the observed  $y$  is not.

The slope  $b$  is what we care about because it summarizes the association between  $y$  and  $x$ .

### The slope is a random

Because we estimate  $b$  from data, it is a random variable with a mean  $\mu_b$  and standard error  $\sigma_b$ .

The standard error can be computed to be (we don't explain why)

$$\sigma_b = \frac{\sigma}{\sigma_x \sqrt{n-1}}$$

We estimate this with the sample standard deviations

$$s_b = \frac{s}{s_x \sqrt{n-1}}$$

Here  $s$  is the standard deviation of the residuals  $y - \hat{y}$ .

### ! Important

In the exam you don't need to compute  $s_b$ . It is given to you as computer output. Looks something like this:

Predictor	Coef	SE Coef	T	P
Constant	40.93831	4.40460	9.294	< 2e-16
son	0.40713	0.06363	6.398	1.36e-09

s = 2.301 R-sq (adj) = 18.33%

- $b$  is calculated for you, it's the Coef for son.
- $s_b$  is calculated as well, it's the SE Coef for son.
- T is simply  $b/s_b$

- $P$  is the p-value for the null the slope being 0.

### Confidence intervals and p-values

If you think of  $b$  as a random variable, you can apply what you have learned before to construct confidence intervals and p-values.

This is because, if the data is normally distributed,  $b/s_b$  follows a t-distribution with  $n - 2$  degrees of freedom with  $n$  the number of observations.

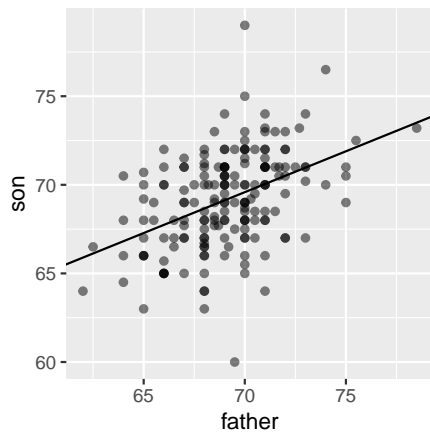
You can use this to compute p-values for the null that there is no relationship.

You calculate

$$P(|b/s_b| > \text{observed value})$$

### Example: heights

Here we add the regression line to the original data:



### Note

The regression formula implies that if we first standardize the variables, that is subtract the average and divide by the standard deviation, then the regression line has intercept 0 and slope equal to the correlation  $\rho$ .

You can make same plot, but using standard units like this:

