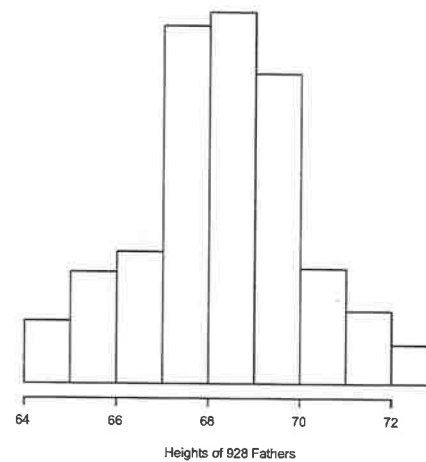


The normal curve

Many data have histograms that look bell-shaped, e.g. heights, weights, IQ scores:



‘The data follow the normal curve.’

But remember that some data have histograms that look quite different, e.g. incomes, house prices.

The empirical rule

If the data follow the normal curve, then

- ▶ about 2/3 (68%) of the data fall within one standard deviation of the mean
- ▶ about 95% fall within 2 standard deviations of the mean
- ▶ about 99.7% fall within 3 standard deviations of the mean

Galton's measurements of heights of fathers have $\bar{x} = 68.3$ in and $s = 1.8$ in.

Therefore about 95% of all heights are between 68.3 in -2×1.8 in $= 64.7$ in and 68.3 in $+2 \times 1.8$ in $= 71.9$ in.

The empirical rule

Recall that in a histogram, percentages are given by areas:

Standardizing data

A normal curve is determined by \bar{x} and s : If the data follow the normal curve, then knowing \bar{x} and s means knowing the whole histogram.

To compute areas under the normal curve, we first **standardize** the data by subtracting off \bar{x} and then dividing by s :

$$z = \frac{\text{height} - \bar{x}}{s}$$

z is called the **standardized value** or **z-score**.

z has no unit (height, \bar{x} and s all have the unit 'inches')

For example, $z = 2$ means the height is 2 standard deviations above average.

$z = -1.5$ means the height is 1.5 standard deviations *below* average.

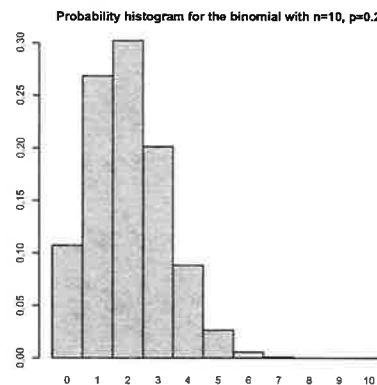
Random variables

You play an online game 10 times. Each time there is a 20% chance to win. The outcomes of the 10 games are due to chance, so the number of wins is random: One set of 10 games might result in 4 wins, another set might result in 7 wins.

X = 'number of successes' is called a **random variable**.

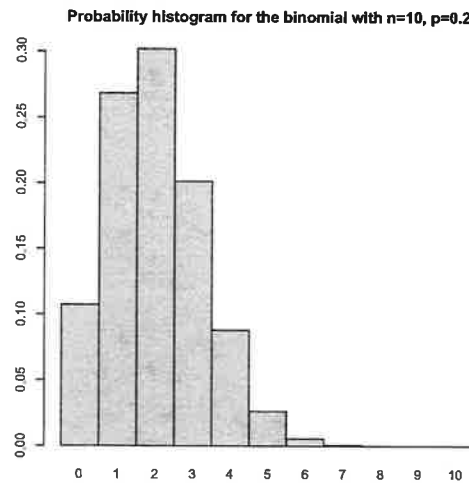
One can calculate $P(X = 2) = 30.2\%$.

We can visualize the probabilities of the various outcomes of X with a **probability histogram**:



The probability histogram

We can visualize the probabilities of the various outcomes of X with a **probability histogram**:



A histogram of data gives percentages for observed data. In contrast, a probability histogram is a theoretical construct: it visualizes probabilities rather than data that have been empirically observed.

Parameter and statistic

What is the average height of adult men in the US?

This is difficult to evaluate as there are 120 million adult men, but it can be *estimated* quite well with a relatively small sample.

The **population** consists of all adult men in the US (about 120 million).

A **parameter** is a quantity of interest about the population: the population average μ , or the population standard deviation σ .

A **statistic (estimate)** is the quantity of interest as measured in the sample: the sample average \bar{x} , or the sample standard deviation s .

The expected value

If we sample an adult male at random, then we expect his height to be around the population average μ , give or take about one standard deviation σ .

The **expected value** of one random draw is the population average μ .

How about \bar{x}_n , the average of n draws?

The **expected value of the sample average**, $E(\bar{x}_n)$, is the population average μ .

But remember that \bar{x}_n is a random variable because sampling is a random process.

So \bar{x}_n won't be exactly equal to $\mu = 69.3$ in: We might get, say, $\bar{x}_n = 70.1$ in. Taking another sample of size n might result in $\bar{x}_n = 69.1$ in.

How far off from μ will \bar{x}_n be?

The **standard error (SE)** of a statistic tells roughly how far off the statistic will be from its expected value.

The standard error for the sample average

The **standard error (SE)** of a statistic tells roughly how far off the statistic will be from its expected value.

So the SE for a statistic plays the same role that the standard deviation σ plays for one observation drawn at random.

The **square root law** is key for statistical inference:

$$\text{SE}(\bar{x}_n) = \frac{\sigma}{\sqrt{n}}$$

The importance of the square root law is twofold:

- ▶ It shows that the SE becomes smaller if we use a larger sample size n . We can use the formula to determine what sample size is required for a desired accuracy.
- ▶ The formula for the standard error **does not depend on the size of the population**, only on the size of the sample.

Expected value and standard error for percentages

What percentage of likely voters approve of the way the US President is handling his job?

The 'percentage of likely voters' is an average. This becomes clear by using the **framework for counting and classifying**:

- ▶ The population consists of all likely voters (about 140 million).
- ▶ Each likely voter falls into one of two categories: approve or not approve.
- ▶ Put the label '1' on each likely voter who approves, and '0' on each who doesn't.
- ▶ Then the number of likely voters who approve equals the sum of all 140 million labels.



- ▶ The percentage of likely voters who approve is the percentage of 1s among the labels.

Expected value and standard error for percentages

In a sample of n likely voters

- ▶ the number of voters in the sample who are approving is the sum of the draws
- ▶ the percentage of voters approving is the percentage of 1s, which is $\frac{\text{sum}}{n} \times 100\% = \bar{x}_n \times 100\%$

Therefore

$$E(\text{percentage of 1s}) = \mu \times 100\% \qquad SE(\text{percentage of 1s}) = \frac{\sigma}{\sqrt{n}} \times 100\%$$

where μ is the population average (=proportion of 1s) and σ is the standard deviation of the population of 0s and 1s.

All of the above formulas are for sampling with replacement. They are still approximately true when sampling without replacement if the sample size is much smaller than the size of the population.

The sampling distribution

Toss a coin 100 times. The number of tails has the following possible outcomes:
 $0, 1, 2, \dots, 100$.

How likely is each outcome?

The number of tails has the binomial distribution with $n = 100$ and $p = 0.5$.
(‘success’ = coin lands tails)

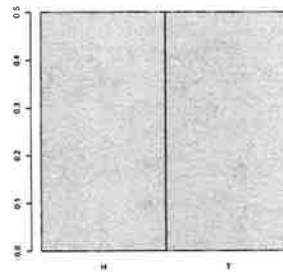
So if the statistic of interest is S_n = ‘number of tails’, then S_n is a random variable whose probability histogram is given by the binomial distribution. This is called the **sampling distribution** of the statistic S_n .

The sampling distribution of S_n provides more detailed information about the chance properties of S_n than the summary numbers given by the expected value and the standard error.

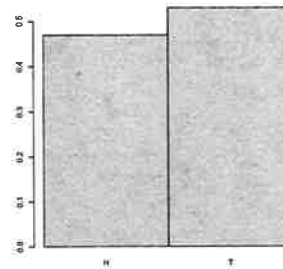
There are three histograms

The chance process of tossing a coin 100 times comes with three different histograms:

1. The probability histogram for producing the data:

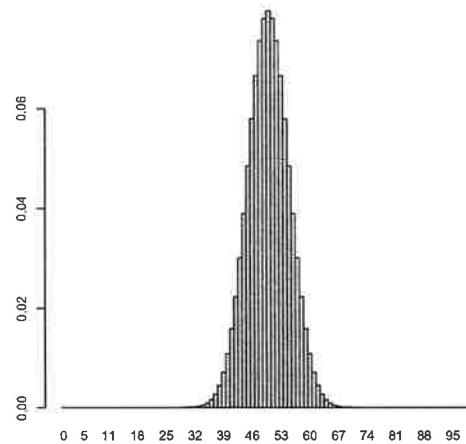


2. The histogram of the 100 observed tosses. This is an empirical histogram of real data:



There are three histograms

3. The probability histogram of the statistic $S_{100} = \text{'number of tails'}$, which shows the sampling distribution of S_{100} :



When doing statistical inference it is important to carefully distinguish these three histograms.

The law of large numbers

The square root law says that $SE(\bar{x}_n)$, the standard error of the sample mean, goes to zero as the sample size increases.

Therefore the sample mean \bar{x}_n will likely be close to its expected value μ if the sample size is large. This is the **law of large numbers**.

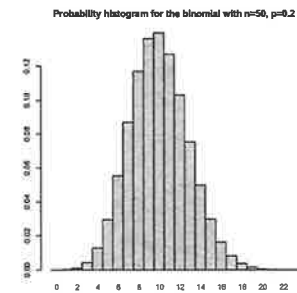
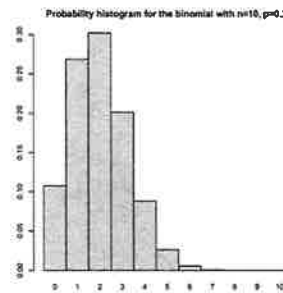
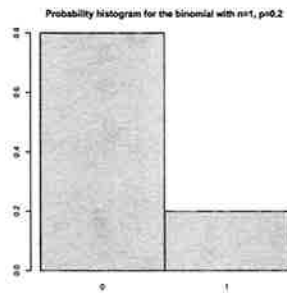
Keep in mind that the law of large numbers applies

- ▶ for averages and therefore also for percentages, but not for sums as their SE *increases*
- ▶ for sampling with replacement from a population, or for simulating data from a probability histogram

More advanced versions of the law of large numbers state that the empirical histogram of the data (the histogram in 2. in the previous section) will be close to the probability histogram in 1. if the sample size is large.

The central limit theorem

Recall the online game where you win with probability 0.2. We looked at the random variable X = 'number of wins' in n gambles and found that X has the binomial distribution with that n and $p = 0.2$.

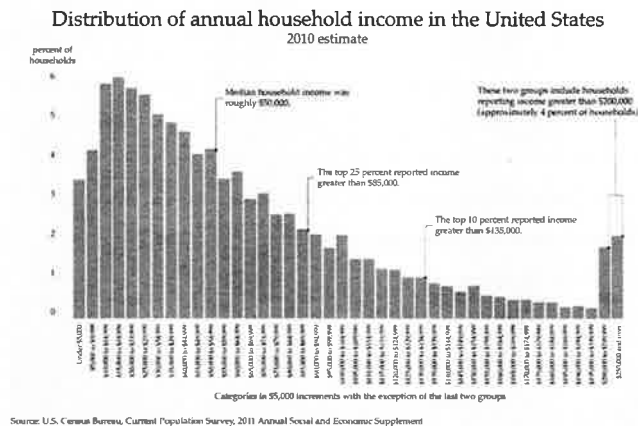


As n gets large, the probability histogram looks more and more similar to the normal curve. This is an example of the **central limit theorem**:

When sampling with replacement and n is large, then the sampling distribution of the sample average (or sum or percentage) approximately follows the normal curve. To standardize, subtract off the expected value of the statistic, then divide by its SE.

The central limit theorem

The key point of the theorem is that we know that the sampling distribution of the statistic is normal *no matter what the population histogram is*:



$$\mu = \$67,000$$

$$\sigma = \$38,000$$

If we sample n incomes at random, then the sample average \bar{x}_n follows the normal curve centered at $E(\bar{x}_n) = \mu = \$67,000$ and with its spread given by

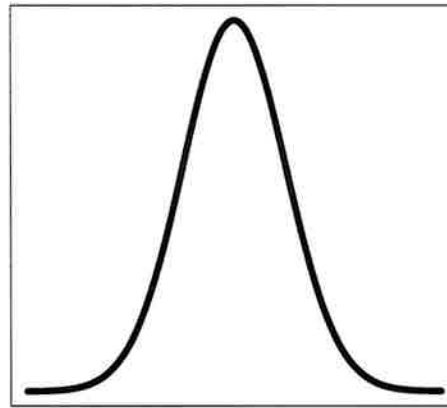
$$SE(\bar{x}_n) = \frac{\sigma}{\sqrt{n}} = \frac{\$38,000}{\sqrt{n}}$$

The central limit theorem

If we sample n incomes at random, then the sample average \bar{x}_n follows the normal curve centered at $E(\bar{x}_n) = \mu = \$67,000$ and with its spread given by

$$SE(\bar{x}_n) = \frac{\sigma}{\sqrt{n}} = \frac{\$38,000}{\sqrt{n}}.$$

For example, if we sample 100 incomes, then by the empirical rule there is about a 16% chance that \bar{x}_n is larger than \$ 70,800:



When does the central limit theorem apply?

For the normal approximation to work, the key requirements are:

- ▶ We sample with replacement, or we simulate independent random variables from the same distribution.
- ▶ The statistic of interest is a sum (averages and percentages are sums in disguise).
- ▶ The sample size is large enough: the more skewed the population histogram is, the larger the required sample size n .
(if there is no strong skewness then $n \geq 15$ is sufficient)

Mini quiz

1. There are two candidates running for governor in CA and they are said to have roughly equal support from the voters. To get a better idea who is ahead, a company polls 400 of the 20 million registered voters in California. Likewise, there are two candidates running for mayor in Palo Alto who are said to have roughly equal support, and the company polls 400 out of the 20,000 registered voters in Palo Alto. Will the first poll be more/equal/less accurate than the second?
2. The average taxable income reported on tax returns for the year 2016 is \$ 45,000, and the standard deviation of the taxable incomes is \$ 23,000. For each of the following two statements, state whether it is true or false:
 - a. The percentage of taxable incomes that fall below \$ 30,000 can be computed from the above information using normal approximation.
 - b. The chances that the sum of 100 randomly selected taxable incomes exceeds \$ 4 million can be computed from the above information using normal approximation.