

# *Introduction to Machine Learning*

## **A quick refresher course in probability theory**

Third lecture, 19.01.2022

Phuc Loi Luu, PhD  
[p.luu@garvan.org.au](mailto:p.luu@garvan.org.au)  
[luu.p.loi@googlemail.com](mailto:luu.p.loi@googlemail.com)

## *Roadmap for today*

---

### 1. Hypothesis testing:

- I toss a coin ten times and get nine heads. How unlikely is that? Can we continue to believe that the coin is *fair* when it produces nine heads out of ten tosses?

### 2. Likelihood and estimation:

- Suppose we know that our random variable is (say)  $\text{Binomial}(10, p)$ , for some  $p$ , but we don't know the value of  $p$ . We will see how to *estimate* the value of  $p$  using maximum likelihood estimation.

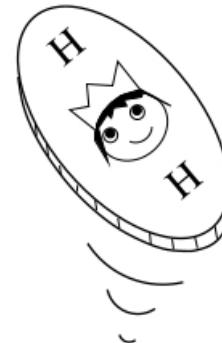
# *Hypothesis testing*

You have probably come across the idea of hypothesis tests,  $p$ -values, and significance in other courses. Common hypothesis tests include  $t$ -tests and chi-squared tests. However, hypothesis tests can be conducted in much simpler circumstances than these. The concept of the hypothesis test is at its easiest to understand with the Binomial distribution in the following example. All other hypothesis tests throughout statistics are based on the same idea.

## *Example: Weird Coin?*

I toss a coin 10 times and get 9 heads. How weird is that?

### What is ‘weird’?



- Getting 9 heads out of 10 tosses: we’ll call this *weird*.
- Getting 10 heads out of 10 tosses: *even more weird!*
- Getting 8 heads out of 10 tosses: *less weird*.
- Getting 1 head out of 10 tosses: *same as getting 9 tails out of 10 tosses: just as weird as 9 heads if the coin is fair.*
- Getting 0 heads out of 10 tosses: *same as getting 10 tails: more weird than 9 heads if the coin is fair.*

# *Hypothesis testing*

## Set of weird outcomes

*If* our coin is fair, the outcomes that are *as weird or weirder* than 9 heads are:

*9 heads, 10 heads, 1 head, 0 heads.*

## So how weird is 9 heads or worse, if the coin is fair?

Define  $X = \#\text{heads out of } 10 \text{ tosses}$ .

Distribution of  $X$ , if the coin is fair:  $X \sim \text{Binomial}(n = 10, p = 0.5)$ .

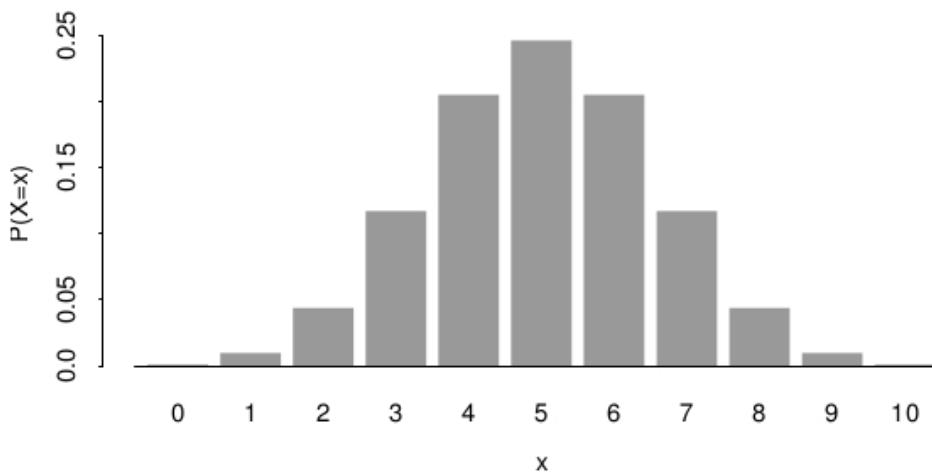
## Probability of observing something at least as weird as 9 heads, if the coin is fair:

We can add the probabilities of all the outcomes that are *at least as weird* as 9 heads out of 10 tosses, assuming that the coin is fair.

$$\mathbb{P}(X = 9) + \mathbb{P}(X = 10) + \mathbb{P}(X = 1) + \mathbb{P}(X = 0) \quad \text{where} \quad X \sim \text{Binomial}(10, 0.5).$$

# *Hypothesis testing*

Probabilities for Binomial( $n = 10, p = 0.5$ )



For  $X \sim \text{Binomial}(10, 0.5)$ , we have:

$$\begin{aligned}\mathbb{P}(X = 9) + \mathbb{P}(X = 10) + \mathbb{P}(X = 1) + \mathbb{P}(X = 0) &= \\ \binom{10}{9}(0.5)^9(0.5)^1 + \binom{10}{10}(0.5)^{10}(0.5)^0 + \\ \binom{10}{1}(0.5)^1(0.5)^9 + \binom{10}{0}(0.5)^0(0.5)^{10} &= \\ 0.00977 + 0.00098 + 0.00977 + 0.00098 &= \\ 0.021.\end{aligned}$$

# *Hypothesis testing*

## **Is this weird?**

---

Yes, it is quite weird. If we had a fair coin and tossed it 10 times, we would only expect to see something as extreme as 9 heads on about *2.1% of occasions*.

## **Is the coin fair?**

---

Obviously, we can't say. It might be: after all, on 2.1% of occasions that you toss a fair coin 10 times, you do get something as weird as 9 heads or more.

However, 2.1% is a small probability, so it is still very unusual for a fair coin to produce something as weird as what we've seen. If the coin really was fair, it would be very unusual to get 9 heads or more.

We can deduce that, *EITHER we have observed a very unusual event with a fair coin, OR the coin is not fair.*

In fact, this gives us *some evidence that the coin is not fair*.

p-value

The value 2.1% *measures the strength of our evidence. The smaller this probability, the more evidence we have.*

# *Hypothesis testing*

## Formal hypothesis test

We now formalize the procedure above. Think of the steps:

- We have a question that we want to answer: *Is the coin fair?*
- There are two alternatives:
  1. *The coin is fair.*
  2. *The coin is not fair.*
- Our observed information is  $X$ , the number of heads out of 10 tosses. We write down the distribution of  $X$  *if the coin is fair*:  
$$X \sim \text{Binomial}(10, 0.5).$$
- We calculate the probability of observing something *AT LEAST AS EXTREME as our observation,  $X = 9$ , if the coin is fair*:  $\text{prob}=0.021$ .
- The probability is small (2.1%). We conclude that this is unlikely with a fair coin, so *we have observed some evidence that the coin is NOT fair.*

# *Null hypothesis and alternative hypothesis*

We express the steps above as two competing hypotheses.

Null hypothesis: *the first alternative, that the coin IS fair.*

*We expect to believe the null hypothesis unless we see convincing evidence that it is wrong.*

Alternative hypothesis: *the second alternative, that the coin is NOT fair.*

In hypothesis testing, we often use this same formulation.

- The null hypothesis is *specific*.

It specifies an exact distribution for our observation:  $X \sim \text{Binomial}(10, 0.5)$ .

- The alternative hypothesis is *general*.

It simply states that the null hypothesis is wrong. It does not say what the *right* answer is.

We use  $H_0$  and  $H_1$  to denote the null and alternative hypotheses respectively.

## *Null hypothesis and alternative hypothesis*

The null hypothesis is  $H_0$  : *the coin is fair.*

The alternative hypothesis is  $H_1$  : *the coin is NOT fair.*

To set up the test, we write:

*Number of heads,  $X \sim \text{Binomial}(10, p)$ ,*

*and*

$$H_0 : p = 0.5$$

$$H_1 : p \neq 0.5.$$

Think of ‘null hypothesis’ as meaning the ‘default’: the hypothesis we will accept unless we have a good reason not to.

## *p-values*

In the hypothesis-testing framework above, we always *measure evidence AGAINST the null hypothesis.*

That is, we believe that our coin is fair unless we see convincing evidence otherwise.

We measure the strength of evidence against  $H_0$  using the *p-value*.

In the example above, the *p-value* was  $p = 0.021$ .

A *p-value* of 0.021 represents *quite strong evidence against the null hypothesis.*

It states that, if the null hypothesis is TRUE, we would only have *a 2.1% chance of observing something as extreme as 9 heads or tails.*

Some people might even see this as strong enough evidence to decide that the null hypothesis is not true, but this is generally an over-simplistic interpretation.

In general, the *p-value* is *the probability of observing something AT LEAST AS EXTREME AS OUR OBSERVATION, if  $H_0$  is TRUE.*

This means that *SMALL p-values represent STRONG evidence against  $H_0$ .*

Small *p*-values mean Strong evidence.  
Large *p*-values mean Little evidence.

## *p-values*

**Note:** Be careful not to confuse the term *p*-value, which is 0.021 in our example, with the Binomial probability  $p$ . Our hypothesis test is designed to test whether the Binomial probability is  $p = 0.5$ . To test this, we calculate the *p*-value of 0.021 as a measure of the strength of evidence ***against*** the hypothesis that  $p = 0.5$ .

## *Interpreting the hypothesis test*

There are different schools of thought about how a  $p$ -value should be interpreted.

- Most people agree that the  $p$ -value is a useful measure of the ***strength of evidence against the null hypothesis***. The smaller the  $p$ -value, the stronger the evidence against  $H_0$ .
- Some people go further and use an ***accept/reject framework***. Under this framework, the null hypothesis  $H_0$  should be *rejected* if the  $p$ -value is less than 0.05 (say), and *accepted* if the  $p$ -value is greater than 0.05.
- In this course we use the ***strength of evidence*** interpretation. The  $p$ -value measures how far out our observation lies in the tails of the distribution specified by  $H_0$ . We do not talk about accepting or rejecting  $H_0$ . This decision should usually be taken in the context of other scientific information.

However, as a rule of thumb, we consider that  $p$ -values of 0.05 and less start to suggest that the null hypothesis is doubtful.

## *Statistical significance*

You have probably encountered the idea of *statistical significance* in other courses.

*Statistical significance refers to the p-value.*

The result of a hypothesis test is *significant at the 5% level* if the *p-value is less than 0.05*.

This means that *the chance of seeing what we did see (9 heads), or more, is less than 5% if the null hypothesis is true.*

Saying the test is *significant* is a quick way of saying that there is evidence against the null hypothesis, usually at the 5% level.

## *Statistical significance*

In the coin example, we can say that our test of  $H_0 : p = 0.5$  against  $H_1 : p \neq 0.5$  is significant at the 5% level, because the *p*-value is 0.021 which is < 0.05.

This means:

- we have some evidence that  $p \neq 0.5$ .

It does **not** mean:

- the difference between  $p$  and 0.5 is *large*, or
- the difference between  $p$  and 0.5 is *important in practical terms*.

Statistically significant means that we have evidence, in OUR sample, that  $p$  is different from 0.5. It says NOTHING about the SIZE, or the IMPORTANCE, of the difference.

“Substantial evidence of a difference”, not “Evidence of a substantial difference.”

# *Statistical significance*

## **Beware!**

---

The *p*-value gives the *probability of seeing something as weird as what we did see, if  $H_0$  is true.*

This means that *5% of the time, we will get a p-value < 0.05 WHEN  $H_0$  IS TRUE!!*

Similarly, about once in every thousand tests, we will get a *p*-value < 0.001, when  $H_0$  is true!

*A small p-value does NOT mean that  $H_0$  is definitely wrong.*

## *One-sided and two-sided tests*

The test above is a *two-sided test*. This means that we considered it *just as weird to get 9 tails as 9 heads*.

If we had a good reason, *before* tossing the coin, to believe that the binomial probability could *only* be = 0.5 or > 0.5, i.e. that it would be *impossible* to have  $p < 0.5$ , then we could conduct a one-sided test:  $H_0 : p = 0.5$  versus  $H_1 : p > 0.5$ .

This would have the effect of halving the resultant  $p$ -value.

## *Example: Presidents and deep-sea divers*

Men in the class: would you like to have daughters? Then become a deep-sea diver, a fighter pilot, or a heavy smoker.

Would you prefer sons? Easy!  
Just become a US president.

Numbers suggest that men in different professions tend to have more sons than daughters, or the reverse. Presidents have sons, fighter pilots have daughters. But is it real, or just chance? We can use hypothesis tests to decide.



### The facts

- The 44 US presidents from George Washington to Barack Obama have had a total of 153 children, comprising 88 sons and only 65 daughters: a sex ratio of 1.4 sons for every daughter.
- Two studies of deep-sea divers revealed that the men had a total of 190 children, comprising 65 sons and 125 daughters: a sex ratio of 1.9 daughters for every son.

## *Example: Presidents and deep-sea divers*

Could this happen by chance?

---

Is it possible that the men in each group *really had a 50-50 chance of producing sons and daughters?*

This is the same as the question in Section 2.2.

For the presidents: *If I tossed a coin 153 times and got only 65 heads, could I continue to believe that the coin was fair?*

For the divers: If I tossed a coin 190 times and got only 65 heads, could I continue to believe that the coin was fair?

## *Example: Presidents*

### Hypothesis test for the presidents

We set up the competing hypotheses as follows.

*Let  $X$  be the number of daughters out of 153 presidential children.*

*Then  $X \sim \text{Binomial}(153, p)$ , where  $p$  is the probability that each child is a daughter.*

Null hypothesis:  $H_0 : p = 0.5$ .

Alternative hypothesis:  $H_1 : p \neq 0.5$ .

$p$ -value: We need the probability of getting a result AT LEAST AS EXTREME as  $X = 65$  daughters, if  $H_0$  is true and  $p$  really is 0.5.

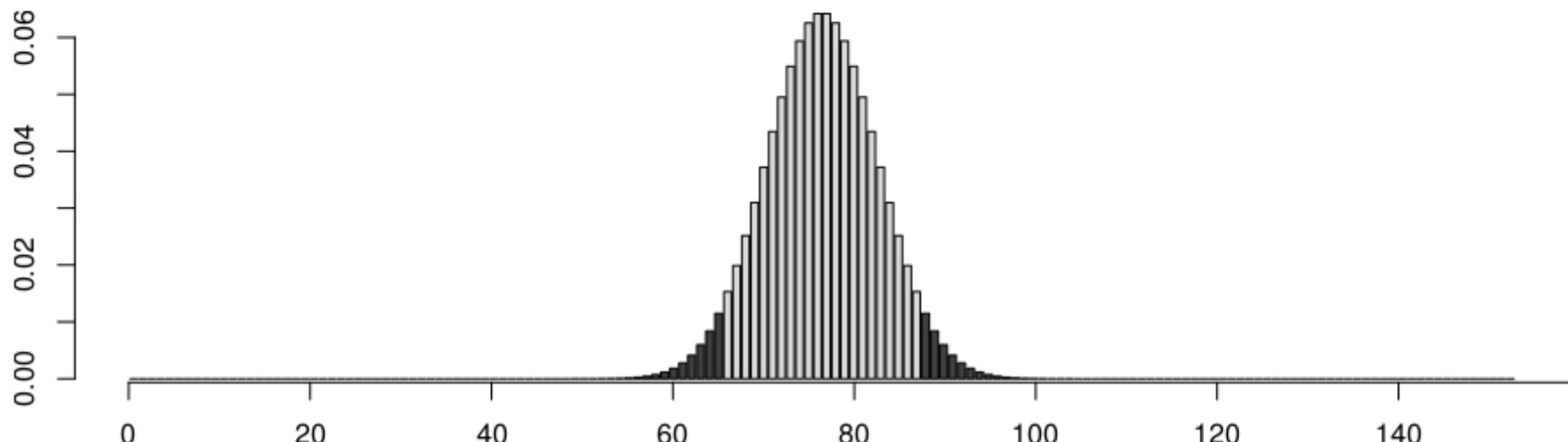
## *Example: Presidents*

Which results are at least as extreme as  $X = 65$ ?

$X = 0, 1, 2, \dots, 65$ , for even fewer daughters.

$X = (153 - 65), \dots, 153$ , for too many daughters, because we would be just as surprised if we saw  $\leq 65$  sons, i.e.  $\geq (153 - 65) = 88$  daughters.

Probabilities for  $X \sim \text{Binomial}(n = 153, p = 0.5)$



# *Example: Presidents*

## Calculating the $p$ -value

The  $p$ -value for the president problem is given by

$$\mathbb{P}(X \leq 65) + \mathbb{P}(X \geq 88) \text{ where } X \sim \text{Binomial}(153, 0.5).$$

In principle, we could calculate this as

$$\mathbb{P}(X = 0) + \mathbb{P}(X = 1) + \dots + \mathbb{P}(X = 65) + \mathbb{P}(X = 88) + \dots + \mathbb{P}(X = 153)$$

$$= \binom{153}{0} (0.5)^0 (0.5)^{153} + \binom{153}{1} (0.5)^1 (0.5)^{152} + \dots$$

This would take a lot of calculator time! Instead, we use a computer with a package such as  $R$ .

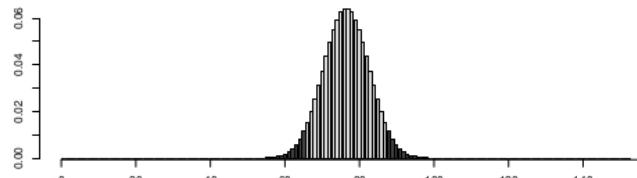
## $R$ command for the $p$ -value

The  $R$  command for calculating the *lower-tail p-value for the Binomial( $n = 153, p = 0.5$ ) distribution is*

`pbinom(65, 153, 0.5).`

Typing this in  $R$  gives:

```
> pbinom(65, 153, 0.5)
[1] 0.03748079
```



This gives us the *lower-tail p-value only*:

$$\mathbb{P}(X \leq 65) = 0.0375.$$

## *Example: Presidents*

To get the overall  $p$ -value:

*Multiply the lower-tail  $p$ -value by 2:*

$$2 \times 0.0375 = 0.0750.$$

In  $R$ :

```
> 2 * pbinom(65, 153, 0.5)
[1] 0.07496158
```

This works because the upper-tail  $p$ -value, by definition, is always going to be the same as the lower-tail  $p$ -value. The upper tail gives us the probability of finding something *equally surprising* at the opposite end of the distribution.

## *Example: Presidents*

**Note:** The *R* command `pbinom` is equivalent to the *cumulative distribution function* for the Binomial distribution:

$$\begin{aligned}\text{pbinom}(65, 153, 0.5) &= \mathbb{P}(X \leq 65) \quad \text{where } X \sim \text{Binomial}(153, 0.5) \\ &= F_X(65) \quad \text{for } X \sim \text{Binomial}(153, 0.5).\end{aligned}$$

The overall *p*-value in this example is  $2 \times F_X(65)$ .

**Note:** In the *R* command `pbinom(65, 153, 0.5)`, the order that you enter the numbers 65, 153, and 0.5 is important. If you enter them in a different order, you will get an error. An alternative is to use the longhand command `pbinom(q=65, size=153, prob=0.5)`, in which case you can enter the terms in any order.

## *Summary: are presidents more likely to have sons?*

Back to our hypothesis test. Recall that  $X$  was the number of daughters out of 153 presidential children, and  $X \sim \text{Binomial}(153, p)$ , where  $p$  is the probability that each child is a daughter.

Null hypothesis:  $H_0 : p = 0.5$ .

Alternative hypothesis:  $H_1 : p \neq 0.5$ .

p-value:  $2 \times F_X(65) = 0.075$ .

### What does this mean?

The  $p$ -value of 0.075 means that, *if the presidents really were as likely to have daughters as sons, there would only be 7.5% chance of observing something as unusual as only 65 daughters out of the total 153 children.*

This is slightly unusual, but not very unusual.

## *Summary: are presidents more likely to have sons?*

We conclude that *there is no real evidence that presidents are more likely to have sons than daughters. The observations are compatible with the possibility that there is no difference.*

Does this mean presidents are equally likely to have sons and daughters? No: *the observations are also compatible with the possibility that there is a difference. We just don't have enough evidence either way.*

## *Hypothesis test for the deep-sea divers*

For the deep-sea divers, there were 190 children: 65 sons, and 125 daughters.

Let  $X$  be the *number of sons out of 190 diver children*.

Then  $X \sim \text{Binomial}(190, p)$ , where  $p$  is the probability that each child is a son.

**Note:** We could just as easily formulate our hypotheses in terms of daughters instead of sons. Because `pbinom` is defined as a lower-tail probability, however, it is usually easiest to formulate them in terms of the *low* result (sons).

Null hypothesis:  $H_0 : p = 0.5$ .

Alternative hypothesis:  $H_1 : p \neq 0.5$ .

p-value: Probability of getting a result AT LEAST AS EXTREME as  $X = 65$  sons, if  $H_0$  is true and  $p$  really is 0.5.

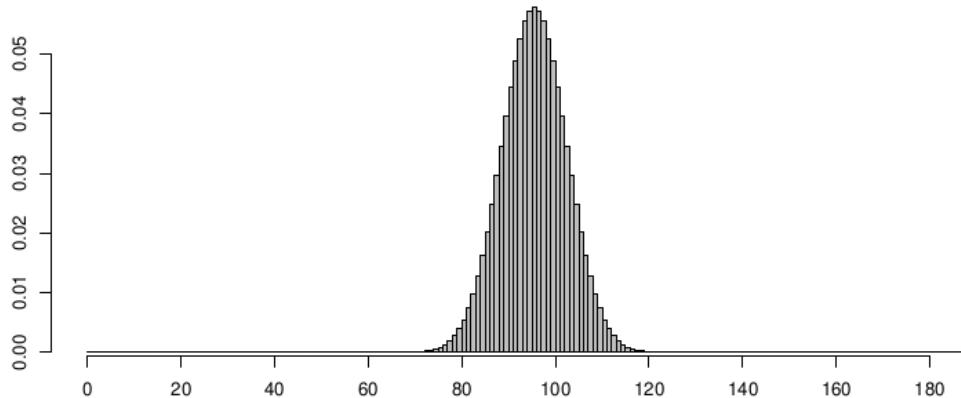
Results at least as extreme as  $X = 65$  are:

$X = 0, 1, 2, \dots, 65$ , for even fewer sons.

$X = (190 - 65), \dots, 190$ , for the equally surprising result in the opposite direction (too many sons).

# *Hypothesis test for the deep-sea divers*

Probabilities for  $X \sim \text{Binomial}(n = 190, p = 0.5)$



R command for the  $p$ -value

$p\text{-value} = 2 \times \text{pbinom}(65, 190, 0.5).$

Typing this in R gives:

```
> 2*pbinom(65, 190, 0.5)
[1] 1.603136e-05
```

This is 0.000016, or a little more than *one chance in 100 thousand*.

We conclude that *it is extremely unlikely that this observation could have occurred by chance, if the deep-sea divers had equal probabilities of having sons and daughters.*

We have *very strong evidence that deep-sea divers are more likely to have daughters than sons. The data are not really compatible with  $H_0$ .*

## *Likelihood and estimation*

So far, the hypothesis tests have only told us whether the Binomial probability  $p$  *might be*, or *probably isn't*, equal to the value specified in the null hypothesis. They have told us nothing about the size, or potential importance, of the departure from  $H_0$ .

For example, for the deep-sea divers, we found that *it would be very unlikely to observe as many as 125 daughters out of 190 children if the chance of having a daughter really was  $p = 0.5$* .

But what does this say about the *actual* value of  $p$ ?

Remember the  $p$ -value for the test was 0.000016. Do you think that:

1.  $p$  could be as big as 0.8?

*No idea! The p-value does not tell us.*

2.  $p$  could be as close to 0.5 as, say, 0.51?

*The test doesn't even tell us this much!*

*If there was a huge sample size (number of children), we COULD get a p-value as small as 0.000016 even if the true probability was 0.51.*

Common sense, however, gives us a hint. Because there were almost twice as many daughters as sons, my guess is that the probability of a having a daughter is something close to  $p = 2/3$ . We need some way of formalizing this.

# *Estimation*

The process of using observations to suggest a value for a parameter is called *estimation*.

The value suggested is called the *estimate* of the parameter.

In the case of the deep-sea divers, we wish to estimate the probability  $p$  that the child of a diver is a daughter. The common-sense estimate to use is

$$p = \frac{\text{number of daughters}}{\text{total number of children}} = \frac{125}{190} = 0.658.$$

However, there are many situations where our common sense fails us. For example, what would we do if we had a regression-model situation (see Section 3.8) and wished to specify an alternative form for  $p$ , such as

$$p = \alpha + \beta \times (\text{diver age}).$$

How would we estimate the unknown intercept  $\alpha$  and slope  $\beta$ , given known information on diver age and number of daughters and sons?

We need a general framework for estimation that can be applied to any situation. The most useful and general method of obtaining parameter estimates is the method of maximum likelihood estimation.

## *Likelihood*

Likelihood is one of the most important concepts in statistics.  
Return to the deep-sea diver example.

$X$  is the *number of daughters out of 190 children*.

We know that  $X \sim \text{Binomial}(190, p)$ ,

and we wish to estimate the value of  $p$ .

The available data is the observed value of  $X$ :  $X = 125$ .

Suppose for a moment that  $p = 0.5$ . What is the probability of observing  $X = 125$ ?

When  $X \sim \text{Binomial}(190, 0.5)$ ,

$$\begin{aligned}\mathbb{P}(X = 125) &= \binom{190}{125} (0.5)^{125} (1 - 0.5)^{190-125} \\ &= 3.97 \times 10^{-6}.\end{aligned}$$

*Not very likely!!*

# Likelihood

What about  $p = 0.6$ ? What would be the probability of observing  $X = 125$  if  $p = 0.6$ ?

*When  $X \sim \text{Binomial}(190, 0.6)$ ,*

$$\begin{aligned}\mathbb{P}(X = 125) &= \binom{190}{125} (0.6)^{125} (1 - 0.6)^{190-125} \\ &= 0.016.\end{aligned}$$

*This still looks quite unlikely, but it is almost 4000 times more likely than getting  $X = 125$  when  $p = 0.5$ .*

So far, we have discovered that *it would be thousands of times more likely to observe  $X = 125$  if  $p = 0.6$  than it would be if  $p = 0.5$ .*

This suggests that  $p = 0.6$  *is a better estimate than  $p = 0.5$ .*

You can probably see where this is heading. If  $p = 0.6$  is a better estimate than  $p = 0.5$ , what if we move  $p$  even closer to our common-sense estimate of 0.658?

*When  $X \sim \text{Binomial}(190, 0.658)$ ,*

$$\begin{aligned}\mathbb{P}(X = 125) &= \binom{190}{125} (0.658)^{125} (1 - 0.658)^{190-125} \\ &= 0.061.\end{aligned}$$

*This is even more likely than for  $p = 0.6$ . So  $p = 0.658$  is the best estimate yet.*

# Likelihood

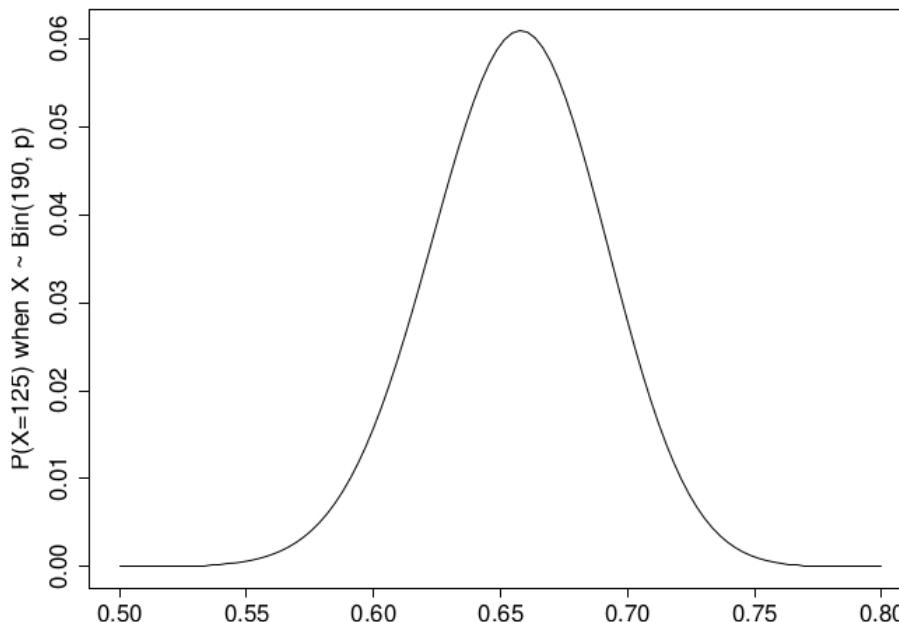
Can we do any better? What happens if we increase  $p$  a little more, say to  $p = 0.7$ ?

When  $X \sim \text{Binomial}(190, 0.7)$ ,

$$\begin{aligned}\mathbb{P}(X = 125) &= \binom{190}{125} (0.7)^{125} (1 - 0.7)^{190-125} \\ &= 0.028.\end{aligned}$$

This has decreased from the result for  $p = 0.658$ , so our observation of 125 is LESS likely under  $p = 0.7$  than under  $p = 0.658$ .

Overall, we can plot a graph showing **how likely** our observation of  $X = 125$  is under each different value of  $p$ .



## *Likelihood*

The graph reaches a *clear maximum*. This is a **value of  $p$**  at which the **observation  $X = 125$  is **MORE LIKELY than**** at any other value of  $p$ .

This ***maximum likelihood*** value of  $p$  is our ***maximum likelihood estimate***.

We can see that the maximum occurs somewhere close to our common-sense estimate of  $p = 0.658$ .

## *The likelihood function*

Look at the graph we plotted overleaf:

Horizontal axis: *The unknown parameter,  $p$ .*

Vertical axis: *The probability of our observation,  $X = 125$ , under this value of  $p$ .*

This function is called the *likelihood function*.

It is a function of *the unknown parameter  $p$* .

For our *fixed* observation  $X = 125$ , the likelihood function shows *how LIKELY the observation 125 is for every different value of  $p$* .

The likelihood function is:

$$L(p) = \mathbb{P}(X = 125) \text{ when } X \sim \text{Binomial}(190, p),$$

$$= \binom{190}{125} p^{125} (1-p)^{190-125}$$

$$= \binom{190}{125} p^{125} (1-p)^{65} \quad \text{for } 0 < p < 1.$$

## *The likelihood function*

This function of  $p$  is the curve shown on the graph on page 55.

In general, if our observation were  $X = x$  rather than  $X = 125$ , the likelihood function is *a function of  $p$  giving  $\mathbb{P}(X = x)$  when  $X \sim \text{Binomial}(190, p)$ .*

We write:

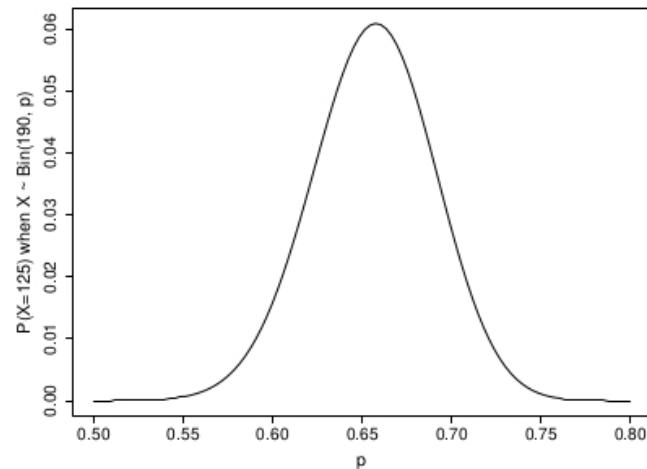
$$\begin{aligned} L(p; x) &= \mathbb{P}(X = x) \text{ when } X \sim \text{Binomial}(190, p), \\ &= \binom{190}{x} p^x (1-p)^{190-x}. \end{aligned}$$

# *Difference between the likelihood function and the probability function*

The likelihood function is *a probability of  $x$ , but it is a FUNCTION of  $p$ .*

The likelihood gives *the probability of a FIXED observation  $x$ , for every possible value of the parameter  $p$ .*

Compare this with the *probability function*, which is *the probability of every different value of  $x$ , for a FIXED value of  $p$ .*

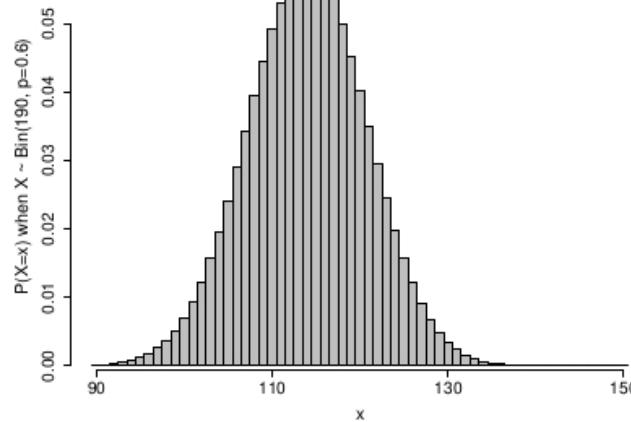


Likelihood function,  $L(p ; x)$ .

Function of  $p$  for fixed  $x$ .

Gives  $\mathbb{P}(X = x)$  as  $p$  changes.

( $x = 125$  here, but could be anything.)



Probability function,  $f_X(x)$ .

Function of  $x$  for fixed  $p$ .

Gives  $\mathbb{P}(X = x)$  as  $x$  changes.

( $p = 0.6$  here, but could be anything.)

## *Maximizing the likelihood*

We have decided that a sensible parameter estimate for  $p$  is the maximum likelihood estimate: *the value of  $p$  at which the observation  $X = 125$  is more likely than at any other value of  $p$ .*

We can find the maximum likelihood estimate using *calculus*.

The likelihood function is

$$L(p; 125) = \binom{190}{125} p^{125} (1-p)^{65}.$$

## *Maximizing the likelihood*

We wish to find the value of  $p$  that maximizes this expression.

To find the maximizing value of  $p$ , *differentiate the likelihood with respect to  $p$ :*

$$\begin{aligned}\frac{dL}{dp} &= \binom{190}{125} \times \left\{ 125 \times p^{124} \times (1-p)^{65} + p^{125} \times 65 \times (1-p)^{64} \times (-1) \right\} \\ &\quad (\text{Product Rule}) \quad (\mathbf{uv}') = u'v + uv' \\ &= \binom{190}{125} \times p^{124} \times (1-p)^{64} \left\{ 125(1-p) - 65p \right\} \\ &= \binom{190}{125} p^{124} (1-p)^{64} \left\{ 125 - 190p \right\}.\end{aligned}$$

The maximizing value of  $p$  occurs when

$$\frac{dL}{dp} = 0.$$

This gives:

$$\begin{aligned}\frac{dL}{dp} &= \binom{190}{125} p^{124} (1-p)^{64} \left\{ 125 - 190p \right\} = 0 \\ \Rightarrow \quad \left\{ 125 - 190p \right\} &= 0 \\ \Rightarrow \quad p &= \frac{125}{190} = 0.658.\end{aligned}$$

## *Maximizing the likelihood*

For the diver example, the maximum likelihood estimate of 125/190 is *the same as the common-sense estimate (page 53)*:

$$p = \frac{\text{number of daughters}}{\text{total number of children}} = \frac{125}{190}.$$

This gives us confidence that the method of maximum likelihood is sensible.

### **The ‘hat’ notation for an estimate**

---

It is conventional to write the estimated value of a parameter with a ‘hat’, like this:  $\hat{p}$ .

For example,

$$\hat{p} = \frac{125}{190}.$$

The correct notation for the maximization is:

$$\left. \frac{dL}{dp} \right|_{p=\hat{p}} = 0 \quad \Rightarrow \quad \hat{p} = \frac{125}{190}.$$

## *Summary of the maximum likelihood procedure*

1. Write down the distribution of  $X$  in terms of the unknown parameter:

$$X \sim \text{Binomial}(190, p).$$

2. Write down the observed value of  $X$ :

*Observed data:  $X = 125$ .*

3. Write down the likelihood function for this observed value:

$$L(p; 125) = \mathbb{P}(X = 125) \text{ when } X \sim \text{Binomial}(190, p)$$

$$= \binom{190}{125} p^{125} (1-p)^{65} \quad \text{for } 0 < p < 1.$$

4. Differentiate the likelihood with respect to the parameter, and set to 0 for the maximum:

$$\frac{dL}{dp} = \binom{190}{125} p^{124} (1-p)^{64} \left\{ 125 - 190p \right\} = 0, \text{ when } p = \hat{p}.$$

This is the *Likelihood Equation*.

5. Solve for  $\hat{p}$ : *From the graph, we can see that  $p = 0$  and  $p = 1$  are not maxima.*

$$\therefore \hat{p} = \frac{125}{190}.$$

This is the *maximum likelihood estimate (MLE)* of  $p$ .

## *Summary of the maximum likelihood procedure*

### **Verifying the maximum**

---

Strictly speaking, when we find the maximum likelihood estimate using

$$\frac{dL}{dp} \Big|_{p=\hat{p}} = 0,$$

we should verify that the result is a maximum (rather than a minimum) by showing that

$$\frac{d^2L}{dp^2} \Big|_{p=\hat{p}} < 0.$$

In Stats 210, we will be relaxed about this. You will usually be told to assume that the MLE occurs in the interior of the parameter range. Where possible, it is always best to *plot the likelihood function, as on page 55*.

This confirms that the maximum likelihood estimate *exists and is unique*.

In particular, *care must be taken when the parameter has a restricted range like  $0 < p < 1$  (see later)*.

## *Estimators*

For the example above, we had observation  $X = 125$ , and the maximum likelihood estimate of  $p$  was

$$\hat{p} = \frac{125}{190}.$$

It is clear that we could follow through the same working with *any* value of  $X$ , which we can write as  $X = x$ , and we would obtain

$$\hat{p} = \frac{x}{190}.$$

**Exercise:** Check this by maximizing the likelihood using  $x$  instead of 125.

This means that even *before* we have made our observation of  $X$ , we can provide a *RULE for calculating the maximum likelihood estimate once  $X$  is observed*:

**Rule:** Let

$$X \sim \text{Binomial}(190, p).$$

*Whatever value of  $X$  we observe, the maximum likelihood estimate of  $p$  will be*

$$\hat{p} = \frac{X}{190}.$$

## *Estimators*

Note that this expression is now a *random variable*: *it depends on the random value of  $X$ .*

A random variable specifying how an estimate is calculated from an observation is called *an estimator*.

In the example above, *the maximum likelihood estimator of  $p$  is*

$$\hat{p} = \frac{X}{190}.$$

*The maximum likelihood estimate of  $p$ , once we have observed that  $X = x$ , is*

$$\hat{p} = \frac{x}{190}.$$

## *General maximum likelihood estimator for $\text{Binomial}(n, p)$*

Take *any* situation in which our observation  $X$  has the distribution

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

where  $n$  is KNOWN and  $p$  is to be estimated.

We make a single observation  $X = x$ .

Follow the steps on page 59 to find the maximum likelihood estimator for  $p$ .

1. Write down the distribution of  $X$  in terms of the unknown parameter:

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

( $n$  is known.)

2. Write down the observed value of  $X$ :

*Observed data:  $X = x$ .*

3. Write down the likelihood function for this observed value:

$$L(p; x) = \mathbb{P}(X = x) \text{ when } X \sim \text{Binomial}(n, p)$$

$$= \binom{n}{x} p^x (1 - p)^{n-x} \quad \text{for } 0 < p < 1.$$

## *General maximum likelihood estimator for Binomial( $n, p$ )*

4. Differentiate the likelihood with respect to the parameter, and set to 0 for the maximum:

$$\frac{dL}{dp} = \binom{n}{x} p^{x-1} (1-p)^{n-x-1} \left\{ x - np \right\} = 0, \text{ when } p = \hat{p}.$$

(Exercise)

5. Solve for  $\hat{p}$ :

$$\hat{p} = \frac{x}{n}.$$

This is the *maximum likelihood estimate of  $p$* .

## *General maximum likelihood estimator for Binomial( $n, p$ )*

The maximum likelihood estimator of  $p$  is

$$\hat{p} = \frac{X}{n}.$$

(Just replace the  $x$  in the MLE with an  $X$ , to convert from the estimate to the estimator.)

By deriving the general maximum likelihood estimator for *any* problem of this sort, we can plug in values of  $n$  and  $x$  to get an instant MLE for any Binomial( $n, p$ ) problem in which  $n$  is known.

**Example:** Recall the president problem in Section 2.3. Out of 153 children, 65 were daughters. Let  $p$  be the probability that a presidential child is a daughter. What is the maximum likelihood estimate of  $p$ ?

**Solution:** Plug in the numbers  $n = 153$ ,  $x = 65$ :

the maximum likelihood estimate is

$$\hat{p} = \frac{x}{n} = \frac{65}{153} = 0.425.$$

**Note:** We showed in Section 2.3 that  $p$  was not significantly different from 0.5 in this example.

However, the MLE of  $p$  is definitely different from 0.5.

This comes back to the meaning of *significantly different* in the statistical sense.

Saying that  $p$  is not significantly different from 0.5 just means that we can't DISTINGUISH any difference between  $p$  and 0.5 from routine sampling variability.