

Introduction

Class 1, 18.05

Jeremy Orloff and Jonathan Bloom

1 Probability vs. Statistics

In this introduction we will preview what we will be studying in 18.05. Don't worry if many of the terms are unfamiliar, they will be explained as the course proceeds.

Probability and statistics are deeply connected because all statistical statements are at bottom statements about probability. Despite this the two sometimes feel like very different subjects. Probability is logically self-contained; there are a few rules and answers all follow logically from the rules, though computations can be tricky. In statistics we apply probability to draw conclusions from data. This can be messy and usually involves as much art as science.

Probability example

You have a fair coin (equal probability of heads or tails). You will toss it 100 times. What is the probability of 60 or more heads? There is only one answer (about 0.028444) and we will learn how to compute it.

Statistics example

You have a coin of unknown provenance. To investigate whether it is fair you toss it 100 times and count the number of heads. Let's say you count 60 heads. Your job as a statistician is to draw a conclusion (inference) from this data. There are many ways to proceed, both in terms of the form the conclusion takes and the probability computations used to justify the conclusion. In fact, different statisticians might draw different conclusions.

Note that in the first example the random process is fully known (probability of heads = .5). The objective is to find the probability of a certain outcome (at least 60 heads) arising from the random process. In the second example, the outcome is known (60 heads) and the objective is to illuminate the unknown random process (the probability of heads).

2 Frequentist vs. Bayesian Interpretations

There are two prominent and sometimes conflicting schools of statistics: [Bayesian](#) and [frequentist](#). Their approaches are rooted in differing interpretations of the meaning of probability.

Frequentists say that probability measures the [frequency of various outcomes of an experiment](#). For example, saying a fair coin has a 50% probability of heads means that if we toss it many times then we expect about half the tosses to land heads.

Bayesians say that probability is an abstract concept that [measures a state of knowledge or a degree of belief](#) in a given proposition. In practice Bayesians do not assign a single value for the probability of a coin coming up heads. Rather they consider a range of values each with its own probability of being true.

In 18.05 we will study and compare these approaches. The frequentist approach has long

been dominant in fields like biology, medicine, public health and social sciences. The Bayesian approach has enjoyed a resurgence in the era of powerful computers and big data. It is especially useful when incorporating new data into an existing statistical model, for example, when training a speech or face recognition system. Today, statisticians are creating powerful tools by using both approaches in complementary ways.

3 Applications, Toy Models, and Simulation

Probability and statistics are used widely in the physical sciences, engineering, medicine, the social sciences, the life sciences, economics and computer science. The list of applications is essentially endless: tests of one medical treatment against another (or a placebo), measures of genetic linkage, the search for elementary particles, machine learning for vision or speech, gambling probabilities and strategies, climate modeling, economic forecasting, epidemiology, marketing, googling... We will draw on examples from many of these fields during this course.

Given so many exciting applications, you may wonder why we will spend so much time thinking about [toy models](#) like coins and dice. By understanding these thoroughly we will develop a good feel for the simple essence inside many complex real-world problems. In fact, the modest coin is a realistic model for any situations with two possible outcomes: success or failure of a treatment, an airplane engine, a bet, or even a class.

Sometimes a problem is so complicated that the best way to understand it is through computer simulation. Here we use software to run *virtual* experiments many times in order to estimate probabilities. In this class we will use R for simulation as well as computation and visualization. Don't worry if you're new to R; we will teach you all you need to know.

Counting and Sets
Class 1, 18.05
Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Know the definitions and notation for sets, intersection, union, complement.
2. Be able to visualize set operations using Venn diagrams.
3. Understand how counting is used computing probabilities.
4. Be able to use the rule of product, inclusion-exclusion principle, permutations and combinations to count the elements in a set.

2 Counting

2.1 Motivating questions

Example 1. A coin is *fair* if it comes up heads or tails with equal probability. You flip a fair coin three times. What is the probability that exactly one of the flips results in a head?

answer: With three flips, we can easily list the eight possible **outcomes**:

$$\{TTT, TTH, THT, THH, HTT, HTH, HHT, HHH\}$$

Three of these outcomes have exactly one head:

$$\{TTH, THT, HTT\}$$

Since all outcomes are equally probable, we have

$$P(1 \text{ head in } 3 \text{ flips}) = \frac{\text{number of outcomes with 1 head}}{\text{total number of outcomes}} = \frac{3}{8}.$$

Think: Would listing the outcomes be practical with 10 flips?

A deck of 52 cards has 13 **ranks** (2, 3, ..., 9, 10, J, Q, K, A) and 4 **suits** (). A poker hand consists of 5 cards. A *one-pair* hand consists of two cards having one rank and three cards having three other ranks, e.g., {2, 2♠, 5, 8♣, K♦}

Test your intuition: the probability of a one-pair hand is:

- (a) less than 5%
- (b) between 5% and 10%
- (c) between 10% and 20%

- (d) between 20% and 40%
- (e) greater than 40%

At this point we can only guess the probability. One of our goals is to learn how to compute it exactly. To start, we note that since every set of five cards is **equally probable**, we can compute the probability of a one-pair hand as

$$P(\text{one-pair}) = \frac{\text{number of one-pair hands}}{\text{total number of hands}}$$

So, to find the exact probability, we need to **count** the number of elements in each of these sets. And we have to be clever about it, because there are too many elements to simply list them all. We will come back to this problem after we have learned some counting techniques.

Several times already we have noted that all the possible outcomes were equally probable and used this to find a probability by counting. Let's state this carefully in the following principle.

Principle: Suppose there are n possible outcomes for an experiment and each is equally probable. If there are k desirable outcomes then the probability of a desirable outcome is k/n . Of course we could replace the word desirable by any other descriptor: undesirable, funny, interesting, remunerative, ...

Concept question: Can you think of a scenario where the possible outcomes are not equally probable?

Here's one scenario: on an exam you can get any score from 0 to 100. That's 101 different possible outcomes. Is the probability you get less than 50 equal to 50/101?

2.2 Sets and notation

Our goal is to learn techniques for counting the number of elements of a set, so we start with a brief review of sets. (If this is new to you, please come to office hours).

2.2.1 Definitions

A **set** S is a collection of elements. We use the following notation.

Element: We write $x \in S$ to mean the element x is in the set S .

Subset: We say the set A is a subset of S if all of its elements are in S . We write this as $A \subset S$.

Complement: The complement of A in S is the set of elements of S that are **not** in A . We write this as A^c or $S - A$.

Union: The union of A and B is the set of all elements in A **or** B (or both). We write this as $A \cup B$.

Intersection: The intersection of A and B is the set of all elements in both A **and** B . We write this as $A \cap B$.

Empty set: The empty set is the set with no elements. We denote it \emptyset .

Disjoint: A and B are **disjoint** if they have no common elements. That is, if $A \cap B = \emptyset$.

Difference: The difference of A and B is the set of elements in A that are not in B . We write this as $A - B$.

Let's illustrate these operations with a simple example.

Example 2. Start with a set of 10 animals

$$S = \{\text{Antelope, Bee, Cat, Dog, Elephant, Frog, Gnat, Hyena, Iguana, Jaguar}\}.$$

Consider two subsets:

$$M = \text{the animal is a mammal} = \{\text{Antelope, Cat, Dog, Elephant, Hyena, Jaguar}\}$$

$$W = \text{the animal lives in the wild} = \{\text{Antelope, Bee, Elephant, Frog, Gnat, Hyena, Iguana, Jaguar}\}.$$

Our goal here is to look at different set operations.

Intersection: $M \cap W$ contains all wild mammals: $M \cap W = \{\text{Antelope, Elephant, Hyena, Jaguar}\}$.

Union: $M \cup W$ contains all animals that are mammals or wild (or both).

$$M \cup W = \{\text{Antelope, Bee, Cat, Dog, Elephant, Frog, Gnat, Hyena, Iguana, Jaguar}\}.$$

Complement: M^c means everything that is *not* in M , i.e. not a mammal. $M^c = \{\text{Bee, Frog, Gnat, Iguana}\}$.

Difference: $M - W$ means everything that's in M and not in W .

$$\text{So, } M - W = \{\text{Cat, Dog}\}.$$

There are often many ways to get the same set, e.g. $M^c = S - M$, $M - W = M \cap L^c$.

The relationship between union, intersection, and complement is given by [DeMorgan's laws](#):

$$(A \cup B)^c = A^c \cap B^c$$

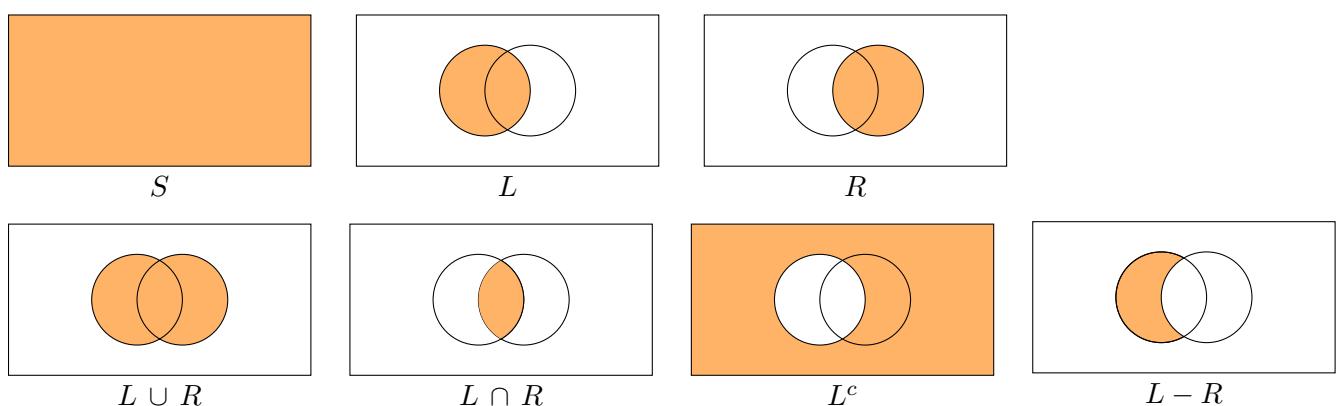
$$(A \cap B)^c = A^c \cup B^c$$

In words the first law says everything not in (A or B) is the same set as everything that's (not in A) and (not in B). The second law is similar.

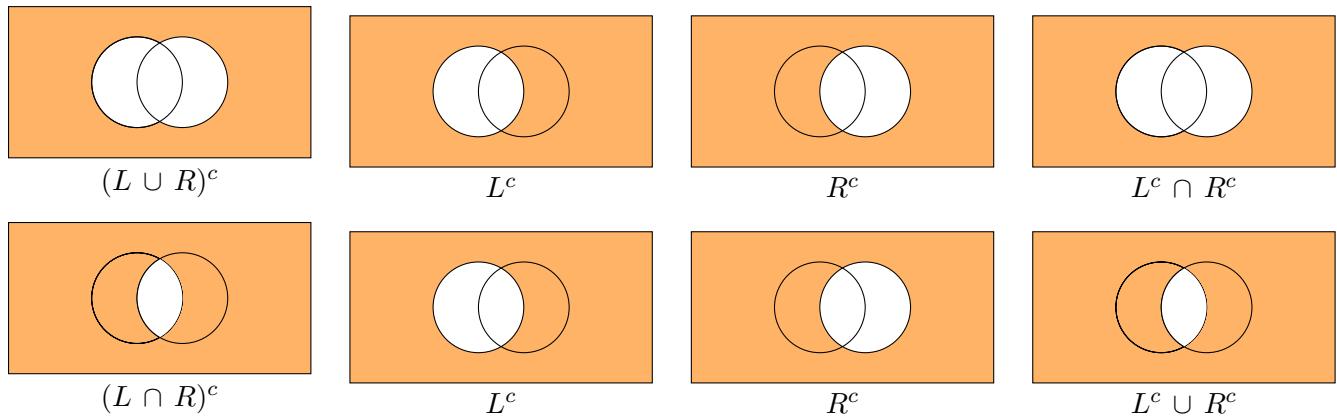
2.2.2 Venn Diagrams

[Venn diagrams](#) offer an easy way to visualize set operations.

In all the figures S is the region inside the large rectangle, L is the region inside the left circle and R is the region inside the right circle. The shaded region shows the set indicated underneath each figure.



Proof of DeMorgan's Laws



Example 3. Verify DeMorgan's laws for the subsets $A = \{1, 2, 3\}$ and $B = \{3, 4\}$ of the set $S = \{1, 2, 3, 4, 5\}$.

answer: For each law we just work through both sides of the equation and show they are the same.

1. $(A \cup B)^c = A^c \cap B^c$:

Right hand side: $A \cup B = \{1, 2, 3, 4\} \Rightarrow (A \cup B)^c = \{5\}$.

Left hand side: $A^c = \{4, 5\}$, $B^c = \{1, 2, 5\} \Rightarrow A^c \cap B^c = \{5\}$.

The two sides are equal. QED

2. $(A \cap B)^c = A^c \cup B^c$:

Right hand side: $A \cap B = \{3\} \Rightarrow (A \cap B)^c = \{1, 2, 4, 5\}$.

Left hand side: $A^c = \{4, 5\}$, $B^c = \{1, 2, 5\} \Rightarrow A^c \cup B^c = \{1, 2, 4, 5\}$.

The two sides are equal. QED

Think: Draw and label a Venn diagram with A the set of Brain and Cognitive Science majors and B the set of sophomores. Shade the region illustrating the first law. Can you express the first law in this case as a non-technical English sentence?

2.2.3 Products of sets

The [product of sets](#) S and T is the set of ordered pairs:

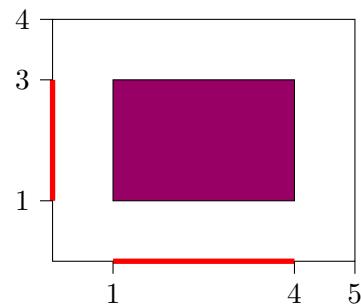
$$S \times T = \{(s, t) \mid s \in S, t \in T\}.$$

In words the right-hand side reads “the set of ordered pairs (s, t) such that s is in S and t is in T .

The following diagrams show two examples of the set product.

\times	1	2	3	4
1	(1,1)	(1,2)	(1,3)	(1,4)
2	(2,1)	(2,2)	(2,3)	(2,4)
3	(3,1)	(3,2)	(3,3)	(3,4)

$$\{1, 2, 3\} \times \{1, 2, 3, 4\}$$



$$[1, 4] \times [1, 3] \subset [0, 5] \times [0, 4]$$

The right-hand figure also illustrates that if $A \subset S$ and $B \subset T$ then $A \times B \subset S \times T$.

2.3 Counting

If S is finite, we use $|S|$ or $\#S$ to denote the number of elements of S .

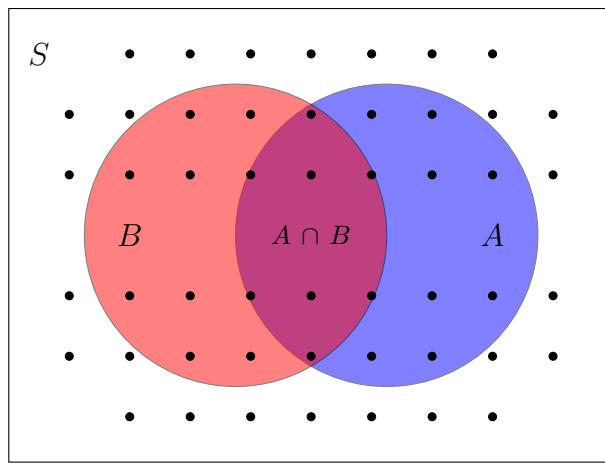
Two useful counting principles are the *inclusion-exclusion principle* and the *rule of product*.

2.3.1 Inclusion-exclusion principle

The *inclusion-exclusion principle* says

$$|A \cup B| = |A| + |B| - |A \cap B|.$$

We can illustrate this with a Venn diagram. S is all the dots, A is the dots in the blue circle, and B is the dots in the red circle.



$|A|$ is the number of dots in A and likewise for the other sets. The figure shows that $|A| + |B|$ double-counts $|A \cap B|$, which is why $|A \cap B|$ is subtracted off in the inclusion-exclusion formula.

Example 4. In a band of singers and guitarists, seven people sing, four play the guitar, and two do both. How big is the band?

answer: Let S be the set singers and G be the set guitar players. The inclusion-exclusion principle says

$$\text{size of band} = |S \cup G| = |S| + |G| - |S \cap G| = 7 + 4 - 2 = 9.$$

2.3.2 Rule of Product

The [Rule of Product](#) says:

If there are n ways to perform action 1 and then by m ways to perform action 2, then there are $n \cdot m$ ways to perform action 1 followed by action 2.

We will also call this the [multiplication](#) rule.

Example 5. If you have 3 shirts and 4 pants then you can make $3 \cdot 4 = 12$ outfits.

Think: An extremely important point is that the rule of product holds even if the ways to perform action 2 depend on action 1, as long as the *number* of ways to perform action 2 is independent of action 1. To illustrate this:

Example 6. There are 5 competitors in the 100m final at the Olympics. In how many ways can the gold, silver, and bronze medals be awarded?

answer: There are 5 ways to award the gold. Once that is awarded there are 4 ways to award the silver and then 3 ways to award the bronze: answer $5 \cdot 4 \cdot 3 = 60$ ways.

Note that the choice of gold medalist affects who can win the silver, but the number of possible silver medalists is always four.

2.4 Permutations and combinations

2.4.1 Permutations

A [permutation](#) of a set is a particular ordering of its elements. For example, the set $\{a, b, c\}$ has six permutations: $abc, acb, bac, bca, cab, cba$. We found the number of permutations by listing them all. We could also have found the number of permutations by using the rule of product. That is, there are 3 ways to pick the first element, then 2 ways for the second, and 1 for the first. This gives a total of $3 \cdot 2 \cdot 1 = 6$ permutations.

In general, the rule of product tells us that the number of permutations of a set of k elements is

$$k! = k \cdot (k - 1) \cdots 3 \cdot 2 \cdot 1.$$

We also talk about the permutations of k things out of a set of n things. We show what this means with an example.

Example 7. List all the permutations of 3 elements out of the set $\{a, b, c, d\}$. **answer:** This is a longer list,

$$\begin{array}{ccccccc} abc & acb & bac & bca & cab & cba \\ abd & adb & bad & bda & dab & dba \\ acd & adc & cad & cda & dac & dca \\ bcd & bdc & cbd & cdb & dbc & dcba \end{array}$$

Note that abc and acb count as distinct permutations. That is, **for permutations the order matters**.

There are 24 permutations. Note that the rule or product would have told us there are $4 \cdot 3 \cdot 2 = 24$ permutations without bothering to list them all.

2.4.2 Combinations

In contrast to permutations, in **combinations order does not matter**: **permutations are lists and combinations are sets**. We show what we mean with an example

Example 8. List all the combinations of 3 elements out of the set $\{a, b, c, d\}$.

answer: Such a combination is a collection of 3 elements without regard to order. So, abc and cab both represent the same combination. We can list all the combinations by listing all the subsets of exactly 3 elements.

$$\{a, b, c\} \quad \{a, b, d\} \quad \{a, c, d\} \quad \{b, c, d\}$$

There are only 4 combinations. Contrast this with the 24 permutations in the previous example. The factor of 6 comes because every combination of 3 things can be written in 6 different orders.

2.4.3 Formulas

We'll use the following notations.

nP_k = number of permutations (lists) of k distinct elements from a set of size n

$nC_k = \binom{n}{k}$ = number of combinations (subsets) of k elements from a set of size n

We emphasise that by the number of combinations of k elements we mean the number of subsets of size k .

These have the following notation and formulas:

$$\begin{aligned} \text{Permutations: } {}nP_k &= \frac{n!}{(n-k)!} = n(n-1)\cdots(n-k+1) \\ \text{Combinations: } {}nC_k &= \frac{n!}{k!(n-k)!} = \frac{{nP_k}}{k!} \end{aligned}$$

The notation nC_k is read “ n choose k ”. The formula for nP_k follows from the rule of product. It also implies the formula for nC_k because a subset of size k can be ordered in $k!$ ways.

We can illustrate the relation between permutations and combinations by lining up the results of the previous two examples.

abc	acb	bac	bca	cab	cba	$\{a, b, c\}$
abd	adb	bad	bda	dab	dba	$\{a, b, d\}$
acd	adc	cad	cda	dac	dca	$\{a, c, d\}$
bcd	bdc	cbd	cdb	dbc	dcb	$\{b, c, d\}$

Permutations: ${}_4P_3$

Combinations: ${}_4C_3$

Notice that each row in the permutations list consists of all $3!$ permutations of the corresponding set in the combinations list.

2.4.4 Examples

Example 9. Count the following:

- (i) The number of ways to choose 2 out of 4 things (order does not matter).
- (ii) The number of ways to list 2 out of 4 things.
- (iii) The number of ways to choose 3 out of 10 things.

answer: (i) This is asking for combinations: $\binom{4}{2} = \frac{4!}{2!2!} = 6$.

(ii) This is asking for permutations: ${}_4P_2 = \frac{4!}{2!} = 12$.

(iii) This is asking for combinations: $\binom{10}{3} = \frac{10!}{3!7!} = \frac{10 \cdot 9 \cdot 8}{3 \cdot 2 \cdot 1} = 120$.

Example 10. (i) Count the number of ways to get 3 heads in a sequence of 10 flips of a coin.

(ii) If the coin is fair, what is the probability of exactly 3 heads in 10 flips.

answer: (i) This asks for the number sequences of 10 flips (heads or tails) with exactly 3 heads. That is, we have to choose exactly 3 out of 10 flips to be heads. This is the same question as in the previous example.

$$\binom{10}{3} = \frac{10!}{3!7!} = \frac{10 \cdot 9 \cdot 8}{3 \cdot 2 \cdot 1} = 120.$$

(ii) Each flip has 2 possible outcomes (heads or tails). So the rule of product says there are $2^{10} = 1024$ sequences of 10 flips. Since the coin is fair each sequence is equally probable. So the probability of 3 heads is

$$\frac{120}{1024} = .117.$$

Probability: Terminology and Examples

Class 2, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Know the definitions of sample space, event and probability function.
2. Be able to organize a scenario with randomness into an experiment and sample space.
3. Be able to make basic computations using a probability function.

2 Terminology

2.1 Probability cast list

- **Experiment:** a repeatable procedure with well-defined possible outcomes.
- **Sample space:** the set of all possible outcomes. We usually denote the sample space by Ω , sometimes by S .
- **Event:** a subset of the sample space.
- **Probability function:** a function giving the probability for each outcome.

Later in the course we will learn about

- Probability density: a continuous distribution of probabilities.
- Random variable: a random numerical outcome.

2.2 Simple examples

Example 1. Toss a fair coin.

Experiment: toss the coin, report if it lands heads or tails.

Sample space: $\Omega = \{H, T\}$.

Probability function: $P(H) = .5, P(T) = .5$.

Example 2. Toss a fair coin 3 times.

Experiment: toss the coin 3 times, list the results.

Sample space: $\Omega = \{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\}$.

Probability function: Each outcome is equally likely with probability $1/8$.

For small sample spaces we can put the set of outcomes and probabilities into a **probability table**.

Outcomes	HHH	HHT	HTH	HTT	THH	THT	TTH	TTT
Probability	1/8	1/8	1/8	1/8	1/8	1/8	1/8	1/8

Example 3. Measure the mass of a proton

Experiment: follow some defined procedure to measure the mass and report the result.

Sample space: $\Omega = [0, \infty)$, i.e. in principle we can get any positive value.

Probability function: since there is a continuum of possible outcomes there is no probability function. Instead we need to use a *probability density*, which we will learn about later in the course.

Example 4. Taxis (An infinite discrete sample space)

Experiment: count the number of taxis that pass 77 Mass. Ave during an 18.05 class.

Sample space: $\Omega = \{0, 1, 2, 3, 4, \dots\}$.

This is often modeled with the following probability function known as the Poisson distribution. (Do not worry about mastering the Poisson distribution just yet):

$$P(k) = e^{-\lambda} \frac{\lambda^k}{k!},$$

where λ is the average number of taxis. We can put this in a table:

Outcome	0	1	2	3	...	k	...
Probability	$e^{-\lambda}$	$e^{-\lambda} \lambda$	$e^{-\lambda} \lambda^2/2$	$e^{-\lambda} \lambda^3/3!$...	$e^{-\lambda} \lambda^k/k!$...

Question: Accepting that this is a valid probability function, what is $\sum_{k=0}^{\infty} e^{-\lambda} \frac{\lambda^k}{k!}$?

answer: This is the total probability of all possible outcomes, so the sum equals 1. (Note, this also follows from the Taylor series $e^\lambda = \sum_{n=0}^{\infty} \frac{\lambda^n}{n!}$.)

In a given setup there can be more than one reasonable choice of sample space. Here is a simple example.

Example 5. Two dice (Choice of sample space)

Suppose you roll one die. Then the sample space and probability function are

Outcome	1	2	3	4	5	6
Probability:	1/6	1/6	1/6	1/6	1/6	1/6

Now suppose you roll two dice. What should be the sample space? Here are two options.

1. Record the pair of numbers showing on the dice (first die, second die).
2. Record the sum of the numbers on the dice. In this case there are 11 outcomes $\{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$. These outcomes are **not all equally likely**.

As above, we can put this information in tables. For the first case, the sample space is the product of the sample spaces for each die

$$\{(1, 1), (2, 1), (3, 1), \dots, (6, 6)\}$$

Each of the 36 outcomes is equally likely. (Why 36 outcomes?) For the probability function we will make a two dimensional table with the rows corresponding to the number on the first die, the columns the number on the second die and the entries the probability.

		Die 2					
		1	2	3	4	5	6
Die 1		1	1/36	1/36	1/36	1/36	1/36
		2	1/36	1/36	1/36	1/36	1/36
		3	1/36	1/36	1/36	1/36	1/36
		4	1/36	1/36	1/36	1/36	1/36
		5	1/36	1/36	1/36	1/36	1/36
		6	1/36	1/36	1/36	1/36	1/36

Two dice in a two dimensional table

In the second case we can present outcomes and probabilities in our usual table.

outcome	2	3	4	5	6	7	8	9	10	11	12
probability	1/36	2/36	3/36	4/36	5/36	6/36	5/36	4/36	3/36	2/36	1/36

The sum of two dice

Think: What is the relationship between the two probability tables above?

We will see that the best choice of sample space depends on the context. For now, simply note that given the outcome as a pair of numbers it is easy to find the sum.

Note. Listing the experiment, sample space and probability function is a good way to start working systematically with probability. It can help you avoid some of the common pitfalls in the subject.

Events.

An **event** is a collection of outcomes, i.e. an event is a subset of the sample space Ω . This sounds odd, but it actually corresponds to the common meaning of the word.

Example 6. Using the setup in Example 2 we would describe the event that you get exactly two heads in words by E = ‘exactly 2 heads’. Written as a subset this becomes

$$E = \{HHT, HTH, THH\}.$$

You should get comfortable moving between describing events in words and as subsets of the sample space.

The probability of an event E is computed by adding up the probabilities of all of the outcomes in E . In this example each outcome has probability $1/8$, so we have $P(E) = 3/8$.

2.3 Definition of a discrete sample space

Definition. A **discrete sample space** is one that is listable, it can be either finite or infinite.

Examples. $\{H, T\}$, $\{1, 2, 3\}$, $\{1, 2, 3, 4, \dots\}$, $\{2, 3, 5, 7, 11, 13, 17, \dots\}$ are all discrete sets. The first two are finite and the last two are infinite.

Example. The interval $0 \leq x \leq 1$ is *not* discrete, rather it is *continuous*. We will deal with continuous sample spaces in a few days.

2.4 The probability function

So far we've been using a casual definition of the probability function. Let's give a more precise one.

Careful definition of the probability function.

For a discrete sample space S a **probability function** P assigns to each outcome ω a number $P(\omega)$ called the probability of ω . P must satisfy two rules:

- Rule 1. $0 \leq P(\omega) \leq 1$ (probabilities are between 0 and 1).
- Rule 2. The sum of the probabilities of all possible outcomes is 1 (something must occur)

In symbols Rule 2 says: if $S = \{\omega_1, \omega_2, \dots, \omega_n\}$ then $P(\omega_1) + P(\omega_2) + \dots + P(\omega_n) = 1$. Or, using summation notation: $\sum_{j=1}^n P(\omega_j) = 1$.

The probability of an event E is the sum of the probabilities of all the outcomes in E . That is,

$$P(E) = \sum_{\omega \in E} P(\omega).$$

Think: Check Rules 1 and 2 on Examples 1 and 2 above.

Example 7. Flip until heads (A classic example)

Suppose we have a coin with probability p of heads and we have the following scenario.

Experiment: Toss the coin until the first heads. Report the number of tosses.

Sample space: $\Omega = \{1, 2, 3, \dots\}$.

Probability function: $P(n) = (1 - p)^{n-1}p$.

Challenge 1: show the sum of all the probabilities equals 1 (hint: geometric series).

Challenge 2: justify the formula for $P(n)$ (we will do this soon).

Stopping problems. The previous toy example is an uncluttered version of a general class of problems called **stopping rule problems**. A stopping rule is a rule that tells you when to end a certain process. In the toy example above the process was flipping a coin and we stopped after the first heads. A more practical example is a rule for ending a series of medical treatments. Such a rule could depend on how well the treatments are working, how the patient is tolerating them and the probability that the treatments would continue to be effective. One could ask about the probability of stopping within a certain number of treatments or the average number of treatments you should expect before stopping.

3 Some rules of probability

For events A , L and R contained in a sample space Ω .

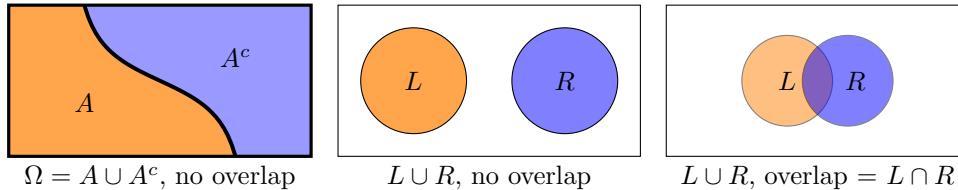
Rule 1. $P(A^c) = 1 - P(A)$.

Rule 2. If L and R are disjoint then $P(L \cup R) = P(L) + P(R)$.

Rule 3. If L and R are not disjoint, we have the **inclusion-exclusion principle**:

$$P(L \cup R) = P(L) + P(R) - P(L \cap R)$$

We visualize these rules using Venn diagrams.



We can also justify them logically.

Rule 1: A and A^c split Ω into two non-overlapping regions. Since the total probability $P(\Omega) = 1$ this rule says that the probability of A and the probability of 'not A ' are complementary, i.e. sum to 1.

Rule 2: L and R split $L \cup R$ into two non-overlapping regions. So the probability of $L \cup R$ is split between $P(L)$ and $P(R)$

Rule 3: In the sum $P(L) + P(R)$ the overlap $P(L \cap R)$ gets counted twice. So $P(L) + P(R) - P(L \cap R)$ counts everything in the union exactly once.

Think: Rule 2 is a special case of Rule 3.

For the following examples suppose we have an experiment that produces a random integer between 1 and 20. The probabilities are not necessarily uniform, i.e., not necessarily the same for each outcome.

Example 8. If the probability of an even number is .6 what is the probability of an odd number?

answer: Since being odd is complementary to being even, the probability of being odd is $1 - .6 = .4$.

Let's redo this example a bit more formally, so you see how it's done. First, so we can refer to it, let's name the random integer X . Let's also name the event 'X is even' as A . Then the event 'X is odd' is A^c . We are given that $P(A) = .6$. Therefore $P(A^c) = 1 - .6 = \boxed{.4}$.

Example 9. Consider the 2 events, A : 'X is a multiple of 2'; B : 'X is odd and less than 10'. Suppose $P(A) = .6$ and $P(B) = .25$.

(i) What is $A \cap B$?

(ii) What is the probability of $A \cup B$?

answer: (i) Since all numbers in A are even and all numbers in B are odd, these events are disjoint. That is, $\boxed{A \cap B = \emptyset}$.

(ii) Since A and B are disjoint $\boxed{P(A \cup B) = P(A) + P(B) = .85}$.

Example 10. Let A , B and C be the events X is a multiple of 2, 3 and 6 respectively. If $P(A) = .6$, $P(B) = .3$ and $P(C) = .2$ what is $P(A \text{ or } B)$?

answer: Note two things. First we used the word 'or' which means union: ' A or B ' = $A \cup B$. Second, an integer is divisible by 6 if and only if it is divisible by both 2 and 3.

This translates into $C = A \cap B$. So the inclusion-exclusion principle says

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) = .6 + .3 - .2 = \boxed{.7}.$$

Conditional Probability, Independence and Bayes' Theorem

Class 3, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Know the definitions of conditional probability and independence of events.
2. Be able to compute conditional probability directly from the definition.
3. Be able to use the multiplication rule to compute the total probability of an event.
4. Be able to check if two events are independent.
5. Be able to use Bayes' formula to 'invert' conditional probabilities.
6. Be able to organize the computation of conditional probabilities using trees and tables.
7. Understand the base rate fallacy thoroughly.

2 Conditional Probability

Conditional probability answers the question 'how does the probability of an event change if we have extra information'. We'll illustrate with an example.

Example 1. Toss a fair coin 3 times.

(a) What is the probability of 3 heads?

answer: Sample space $\Omega = \{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\}$.

All outcomes are equally likely, so $P(3 \text{ heads}) = 1/8$.

(b) Suppose we are told that the first toss was heads. Given this information how should we compute the probability of 3 heads?

answer: We have a new (reduced) sample space: $\Omega' = \{HHH, HHT, HTH, HTT\}$.

All outcomes are equally likely, so

$$P(3 \text{ heads given that the first toss is heads}) = 1/4.$$

This is called **conditional probability**, since it takes into account additional conditions. To develop the notation, we rephrase (b) in terms of *events*.

Rephrased (b) Let A be the event 'all three tosses are heads' = $\{HHH\}$.

Let B be the event 'the first toss is heads' = $\{HHH, HHT, HTH, HTT\}$.

The **conditional probability** of A knowing that B occurred is written

$$P(A|B)$$

This is read as

'the conditional probability of A given B '

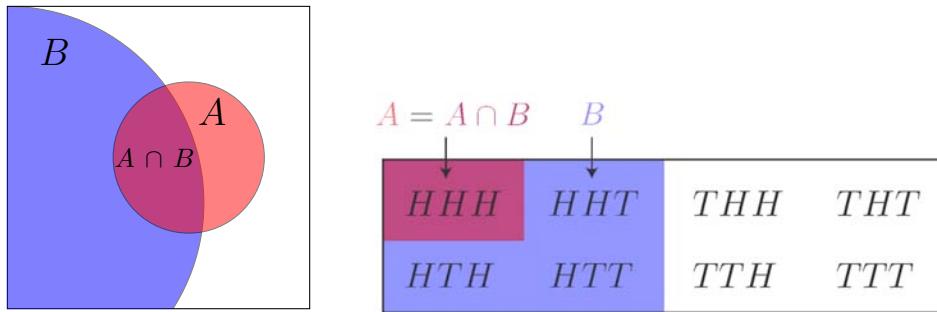
or

'the probability of A conditioned on B '

or simply

'the probability of A given B '.

We can visualize conditional probability as follows. Think of $P(A)$ as the proportion of the area of the *whole* sample space taken up by A . For $P(A|B)$ we restrict our attention to B . That is, $P(A|B)$ is the proportion of area of B taken up by A , i.e. $P(A \cap B)/P(B)$.



Conditional probability: Abstract visualization and coin example

Note, $A \subset B$ in the right-hand figure, so there are only two colors shown.

The formal definition of conditional probability catches the gist of the above example and visualization.

Formal definition of conditional probability

Let A and B be events. We define [the conditional probability](#) of A given B as

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \text{ provided } P(B) \neq 0. \quad (1)$$

Let's redo the coin tossing example using the definition in Equation (1). Recall $A = \text{'3 heads'}$ and $B = \text{'first toss is heads'}$. We have $P(A) = 1/8$ and $P(B) = 1/2$. Since $A \cap B = A$, we also have $P(A \cap B) = 1/8$. Now according to (1), $P(A|B) = \frac{1/8}{1/2} = 1/4$, which agrees with our answer in Example 1b.

3 Multiplication Rule

The following formula is called the [multiplication rule](#).

$$P(A \cap B) = P(A|B) \cdot P(B). \quad (2)$$

This is simply a rewriting of the definition in Equation (1) of conditional probability. We will see that our use of the multiplication rule is very similar to our use of the rule of product in counting. In fact, the multiplication rule is just a souped up version of the rule of product.

We start with a simple example where we can check all the probabilities directly by counting.

Example 2. Draw two cards from a deck. Define the events: $S_1 = \text{'first card is a spade'}$ and $S_2 = \text{'second card is a spade'}$. What is the $P(S_2|S_1)$?

answer: We can do this directly by counting: if the first card is a spade then of the 51 cards remaining, 12 are spades.

$$P(S_2|S_1) = 12/51.$$

Now, let's recompute this using formula (1). We have to compute $P(S_1)$, $P(S_2)$ and $P(S_1 \cap S_2)$: We know that $P(S_1) = 1/4$ because there are 52 equally likely ways to draw the first card and 13 of them are spades. The same logic says that there are 52 equally likely ways the second card can be drawn, so $P(S_2) = 1/4$.

Aside: The probability $P(S_2) = 1/4$ may seem surprising since the value of first card certainly affects the probabilities for the second card. However, if we look at *all* possible two card sequences we will see that every card in the deck has equal probability of being the second card. Since 13 of the 52 cards are spades we get $P(S_2) = 13/52 = 1/4$. Another way to say this is: if we are not given value of the first card then we have to consider all possibilities for the second card.

Continuing, we see that

$$P(S_1 \cap S_2) = \frac{13 \cdot 12}{52 \cdot 51} = 3/51.$$

This was found by counting the number of ways to draw a spade followed by a second spade and dividing by the number of ways to draw any card followed by any other card). Now, using (1) we get

$$P(S_2|S_1) = \frac{P(S_2 \cap S_1)}{P(S_1)} = \frac{3/51}{1/4} = 12/51.$$

Finally, we verify the multiplication rule by computing both sides of (2).

$$P(S_1 \cap S_2) = \frac{13 \cdot 12}{52 \cdot 51} = \frac{3}{51} \quad \text{and} \quad P(S_2|S_1) \cdot P(S_1) = \frac{12}{51} \cdot \frac{1}{4} = \frac{3}{51}. \quad \text{QED}$$

Think: For S_1 and S_2 in the previous example, what is $P(S_2|S_1^c)$?

4 Law of Total Probability

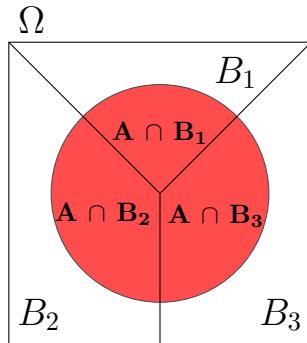
The law of total probability will allow us to use the multiplication rule to find probabilities in more interesting examples. It involves a lot of notation, but the idea is fairly simple. We state the law when the sample space is divided into 3 pieces. It is a simple matter to extend the rule when there are more than 3 pieces.

Law of Total Probability

Suppose the sample space Ω is divided into 3 disjoint events B_1 , B_2 , B_3 (see the figure below). Then for any event A :

$$\begin{aligned} P(A) &= P(A \cap B_1) + P(A \cap B_2) + P(A \cap B_3) \\ P(A) &= P(A|B_1) P(B_1) + P(A|B_2) P(B_2) + P(A|B_3) P(B_3) \end{aligned} \tag{3}$$

The top equation says 'if A is divided into 3 pieces then $P(A)$ is the sum of the probabilities of the pieces'. The bottom equation (3) is called **the law of total probability**. It is just a rewriting of the top equation using the multiplication rule.



The sample space Ω and the event A are each divided into 3 disjoint pieces.

The law holds if we divide Ω into any number of events, so long as they are *disjoint* and *cover* all of Ω . Such a division is often called a *partition* of Ω .

Our first example will be one where we already know the answer and can verify the law.

Example 3. An urn contains 5 red balls and 2 green balls. Two balls are drawn one after the other. What is the probability that the second ball is red?

answer: The sample space is $\Omega = \{\text{rr}, \text{rg}, \text{gr}, \text{gg}\}$.

Let R_1 be the event ‘the first ball is red’, G_1 = ‘first ball is green’, R_2 = ‘second ball is red’, G_2 = ‘second ball is green’. We are asked to find $P(R_2)$.

The fast way to compute this is just like $P(S_2)$ in the card example above. Every ball is equally likely to be the second ball. Since 5 out of 7 balls are red, $P(R_2) = 5/7$.

Let’s compute this same value using the law of total probability (3). First, we’ll find the conditional probabilities. This is a simple counting exercise.

$$P(R_2|R_1) = 4/6, \quad P(R_2|G_1) = 5/6.$$

Since R_1 and G_1 partition Ω the law of total probability says

$$\begin{aligned} P(R_2) &= P(R_2|R_1)P(R_1) + P(R_2|G_1)P(G_1) \\ &= \frac{4}{6} \cdot \frac{5}{7} + \frac{5}{6} \cdot \frac{2}{7} \\ &= \frac{30}{42} = \frac{5}{7}. \end{aligned} \tag{4}$$

Probability urns

The example above used probability urns. Their use goes back to the beginning of the subject and we would be remiss not to introduce them. This toy model is very useful. We quote from Wikipedia: http://en.wikipedia.org/wiki/Urn_problem

In probability and statistics, an urn problem is an idealized mental exercise in which some objects of real interest (such as atoms, people, cars, etc.) are represented as colored balls in an urn or other container. One pretends to draw (remove) one or more balls from the urn; the goal is to determine the probability of drawing one color or another, or some other properties. A key parameter is whether each ball is returned to the urn after each draw.

It doesn't take much to make an example where (3) is really the best way to compute the probability. Here is a game with slightly more complicated rules.

Example 4. An urn contains 5 red balls and 2 green balls. A ball is drawn. If it's green a red ball is added to the urn and if it's red a green ball is added to the urn. (The original ball is not returned to the urn.) Then a second ball is drawn. What is the probability the second ball is red?

answer: The law of total probability says that $P(R_2)$ can be computed using the expression in Equation (4). Only the values for the probabilities will change. We have

$$P(R_2|R_1) = 4/7, \quad P(R_2|G_1) = 6/7.$$

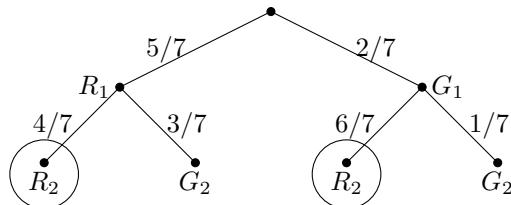
Therefore,

$$P(R_2) = P(R_2|R_1)P(R_1) + P(R_2|G_1)P(G_1) = \frac{4}{7} \cdot \frac{5}{7} + \frac{6}{7} \cdot \frac{2}{7} = \frac{32}{49}.$$

5 Using Trees to Organize the Computation

Trees are a great way to organize computations with conditional probability and the law of total probability. The figures and examples will make clear what we mean by a tree. As with the rule of product, the key is to organize the underlying process into a sequence of actions.

We start by redoing Example 4. The sequence of actions are: first draw ball 1 (and add the appropriate ball to the urn) and then draw ball 2.



You interpret this tree as follows. Each dot is called a **node**. The tree is organized by levels. The top node (**root node**) is at level 0. The next layer down is level 1 and so on. Each level shows the outcomes at one stage of the game. Level 1 shows the possible outcomes of the first draw. Level 2 shows the possible outcomes of the second draw starting from each node in level 1.

Probabilities are written along the branches. The probability of R_1 (red on the first draw) is $5/7$. It is written along the branch from the root node to the one labeled R_1 . At the next level we put in **conditional** probabilities. The probability along the branch from R_1 to R_2 is $P(R_2|R_1) = 4/7$. It represents the probability of going to node R_2 given that you are already at R_1 .

The multiplication rule says that the probability of getting to any node is just the product of the probabilities along the path to get there. For example, the node labeled R_2 at the far left really represents the event $R_1 \cap R_2$ because it comes from the R_1 node. The multiplication rule now says

$$P(R_1 \cap R_2) = P(R_1) \cdot P(R_2|R_1) = \frac{5}{7} \cdot \frac{4}{7},$$

which is exactly multiplying along the path to the node.

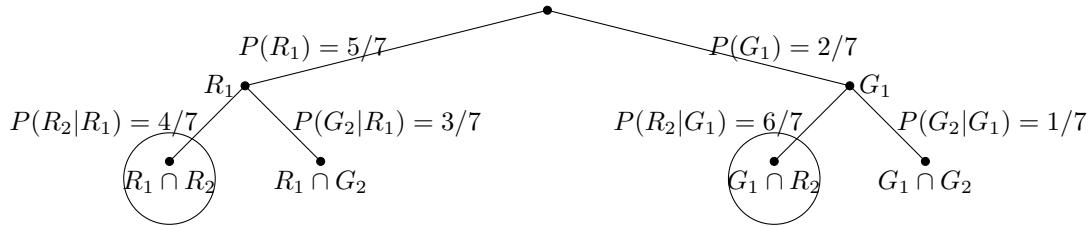
The law of total probability is just the statement that $P(R_2)$ is the sum of the probabilities of all paths leading to R_2 (the two circled nodes in the figure). In this case,

$$P(R_2) = \frac{5}{7} \cdot \frac{4}{7} + \frac{2}{7} \cdot \frac{6}{7} = \frac{32}{49},$$

exactly as in the previous example.

5.1 Shorthand vs. precise trees

The tree given above involves some shorthand. For example, the node marked R_2 at the far left really represents the event $R_1 \cap R_2$, since it ends the path from the root through R_1 to R_2 . Here is the same tree with everything labeled precisely. As you can see this tree is more cumbersome to make and use. We usually use the shorthand version of trees. You should make sure you know how to interpret them precisely.



6 Independence

Two events are independent if knowledge that one occurred does not change the probability that the other occurred. Informally, events are independent if they do not influence one another.

Example 5. Toss a coin twice. We expect the outcomes of the two tosses to be independent of one another. In real experiments this always has to be checked. If my coin lands in honey and I don't bother to clean it, then the second toss might be affected by the outcome of the first toss.

More seriously, the independence of experiments can be undermined by the failure to clean or recalibrate equipment between experiments or to isolate supposedly independent observers from each other or a common influence. We've all experienced hearing the same 'fact' from different people. Hearing it from different sources tends to lend it credence until we learn that they all heard it from a common source. That is, our sources were not independent.

Translating the verbal description of independence into symbols gives

$$A \text{ is independent of } B \quad \text{if} \quad P(A|B) = P(A). \tag{5}$$

That is, knowing that B occurred does not change the probability that A occurred. In terms of events as subsets, knowing that the realized outcome is in B does not change the probability that it is in A .

If A and B are independent in the above sense, then the multiplication rule gives $P(A \cap B) = P(A|B) \cdot P(B) = P(A) \cdot P(B)$. This justifies the following technical definition of independence.

Formal definition of independence: Two events A and B are **independent** if

$$P(A \cap B) = P(A) \cdot P(B) \quad (6)$$

This is a nice symmetric definition which makes clear that A is independent of B if and only if B is independent of A . Unlike the equation with conditional probabilities, this definition makes sense even when $P(B) = 0$. In terms of conditional probabilities, we have:

1. If $P(B) \neq 0$ then A and B are independent if and only if $P(A|B) = P(A)$.
2. If $P(A) \neq 0$ then A and B are independent if and only if $P(B|A) = P(B)$.

Independent events commonly arise as different trials in an experiment, as in the following example.

Example 6. Toss a fair coin twice. Let H_1 = ‘heads on first toss’ and let H_2 = ‘heads on second toss’. Are H_1 and H_2 independent?

answer: Since $H_1 \cap H_2$ is the event ‘both tosses are heads’ we have

$$P(H_1 \cap H_2) = 1/4 = P(H_1)P(H_2).$$

Therefore the events are independent.

We can ask about the independence of any two events, as in the following two examples.

Example 7. Toss a fair coin 3 times. Let H_1 = ‘heads on first toss’ and A = ‘two heads total’. Are H_1 and A independent?

answer: We know that $P(A) = 3/8$. Since this is not 0 we can check if the formula in Equation 5 holds. Now, $H_1 = \{\text{HHH, HHT, HTH, HTT}\}$ contains exactly two outcomes (HHT, HTH) from A , so we have $P(A|H_1) = 2/4$. Since $P(A|H_1) \neq P(A)$ these events are not independent.

Example 8. Draw one card from a standard deck of playing cards. Let’s examine the independence of 3 events ‘the card is an ace’, ‘the card is a heart’ and ‘the card is red’.

Define the events as A = ‘ace’, H = ‘hearts’, R = ‘red’.

(a) We know that $P(A) = 4/52 = 1/13$, $P(A|H) = 1/13$. Since $P(A) = P(A|H)$ we have that A is independent of H .

(b) $P(A|R) = 2/26 = 1/13$. So A is independent of R . That is, whether the card is an ace is independent of whether it’s red.

(c) Finally, what about H and R ? Since $P(H) = 1/4$ and $P(H|R) = 1/2$, H and R are not independent. We could also see this the other way around: $P(R) = 1/2$ and $P(R|H) = 1$, so H and R are not independent.

6.1 Paradoxes of Independence

An event A with probability 0 is independent of itself, since in this case both sides of equation (6) are 0. This appears paradoxical because knowledge that A occurred certainly

gives information about whether A occurred. We resolve the paradox by noting that since $P(A) = 0$ the statement ‘ A occurred’ is vacuous.

Think: For what other value(s) of $P(A)$ is A independent of itself?

7 Bayes' Theorem

Bayes' theorem is a pillar of both probability and statistics and it is central to the rest of this course. For two events A and B [Bayes' theorem](#) (also called [Bayes' rule](#) and [Bayes' formula](#)) says

$$P(B|A) = \frac{P(A|B) \cdot P(B)}{P(A)}. \quad (7)$$

Comments: 1. Bayes' rule tells us how to ‘invert’ conditional probabilities, i.e. to find $P(B|A)$ from $P(A|B)$.

2. In practice, $P(A)$ is often computed using the law of total probability.

Proof of Bayes' rule

The key point is that $A \cap B$ is symmetric in A and B . So the multiplication rule says

$$P(B|A) \cdot P(A) = P(A \cap B) = P(A|B) \cdot P(B).$$

Now divide through by $P(A)$ to get Bayes' rule.

A common mistake is to confuse $P(A|B)$ and $P(B|A)$. They can be very different. This is illustrated in the next example.

Example 9. Toss a coin 5 times. Let H_1 = ‘first toss is heads’ and let H_A = ‘all 5 tosses are heads’. Then $P(H_1|H_A) = 1$ but $P(H_A|H_1) = 1/16$.

For practice, let’s use Bayes' theorem to compute $P(H_1|H_A)$ using $P(H_A|A_1)$. The terms are $P(H_A|H_1) = 1/16$, $P(H_1) = 1/2$, $P(H_A) = 1/32$. So,

$$P(H_1|H_A) = \frac{P(H_A|H_1)P(H_1)}{P(H_A)} = \frac{(1/16) \cdot (1/2)}{1/32} = 1,$$

which agrees with our previous calculation.

7.1 The Base Rate Fallacy

The base rate fallacy is one of many examples showing that it’s easy to confuse the meaning of $P(B|A)$ and $P(A|B)$ when a situation is described in words. This is one of the key examples from probability and it will inform much of our practice and interpretation of statistics. You should strive to understand it thoroughly.

Example 10. The Base Rate Fallacy

Consider a routine screening test for a disease. Suppose the frequency of the disease in the population ([base rate](#)) is 0.5%. The test is highly accurate with a 5% false positive rate and a 10% false negative rate.

You take the test and it comes back positive. What is the probability that you have the disease?

answer: We will do the computation three times: using trees, tables and symbols. We'll use the following notation for the relevant events:

$D+$ = 'you have the disease'

$D-$ = 'you do not have the disease'

$T+$ = 'you tested positive'

$T-$ = 'you tested negative'.

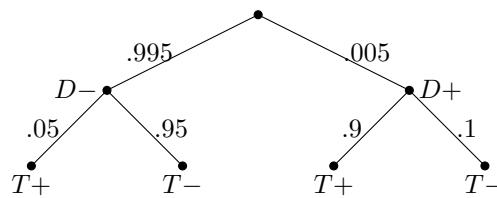
We are given $P(D+) = .005$ and therefore $P(D-) = .995$. The false positive and false negative rates are (by definition) conditional probabilities.

$$P(\text{false positive}) = P(T+|D-) = .05 \quad \text{and} \quad P(\text{false negative}) = P(T-|D+) = .1.$$

The complementary probabilities are known as the true negative and true positive rates:

$$P(T-|D-) = 1 - P(T+|D-) = .95 \quad \text{and} \quad P(T+|D+) = 1 - P(T-|D+) = .9.$$

Trees: All of these probabilities can be displayed quite nicely in a tree.



The question asks for the probability that you have the disease given that you tested positive, i.e. what is the value of $P(D+|T+)$. We aren't given this value, but we do know $P(T+|D+)$, so we can use Bayes' theorem.

$$P(D+|T+) = \frac{P(T+|D+) \cdot P(D+)}{P(T+)}$$

The two probabilities in the numerator are given. We compute the denominator $P(T+)$ using the law of total probability. Using the tree we just have to sum the probabilities for each of the nodes marked $T+$

$$P(T+) = .995 \times .05 + .005 \times .9 = .05425$$

Thus,

$$P(D+|T+) = \frac{.9 \times .005}{.05425} = 0.082949 \approx 8.3\%.$$

Remarks: This is called the base rate fallacy because the base rate of the disease in the population is so low that the vast majority of the people taking the test are healthy, and even with an accurate test most of the positives will be healthy people. Ask your doctor for his/her guess at the odds.

To summarize the base rate fallacy with specific numbers

95% of all tests are accurate does not imply 95% of positive tests are accurate

We will refer back to this example frequently. It and similar examples are at the heart of many statistical misunderstandings.

Other ways to work Example 10

Tables: Another trick that is useful for computing probabilities is to make a table. Let's redo the previous example using a table built with 10000 total people divided according to the probabilities in this example.

We construct the table as follows. Pick a number, say 10000 people, and place it as the grand total in the lower right. Using $P(D+) = .005$ we compute that 50 out of the 10000 people are sick ($D+$). Likewise 9950 people are healthy ($D-$). At this point the table looks like:

	$D+$	$D-$	total
$T+$			
$T-$			
total	50	9950	10000

Using $P(T+|D+) = .9$ we can compute that the number of sick people who tested positive as 90% of 50 or 45. The other entries are similar. At this point the table looks like the table below on the left. Finally we sum the $T+$ and $T-$ rows to get the completed table on the right.

	$D+$	$D-$	total
$T+$	45	498	
$T-$	5	9452	
total	50	9950	10000

	$D+$	$D-$	total
$T+$	45	498	543
$T-$	5	9452	9457
total	50	9950	10000

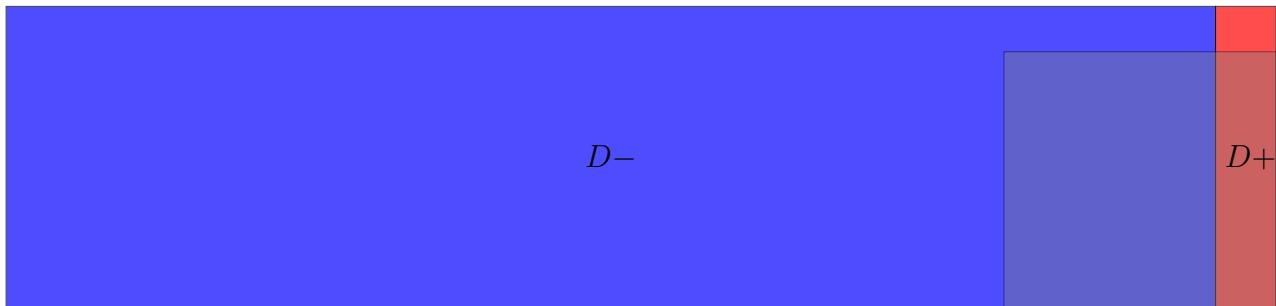
Using the complete table we can compute

$$P(D+|T+) = \frac{|D+ \cap T+|}{|T+|} = \frac{45}{543} = 8.3\%.$$

Symbols: For completeness, we show how the solution looks when written out directly in symbols.

$$\begin{aligned} P(D+|T+) &= \frac{P(T+|D+) \cdot P(D+)}{P(T+)} \\ &= \frac{P(T+|D+) \cdot P(D+)}{P(T+|D+) \cdot P(D+) + P(T+|D-) \cdot P(D-)} \\ &= \frac{.9 \times .005}{.9 \times .005 + .05 \times .995} \\ &= 8.3\% \end{aligned}$$

Visualization: The figure below illustrates the base rate fallacy. The large blue area represents all the healthy people. The much smaller red area represents the sick people. The shaded rectangle represents the the people who test positive. The shaded area covers most of the red area and only a small part of the blue area. Even so, the most of the shaded area is over the blue. That is, most of the positive tests are of healthy people.



7.2 Bayes' rule in 18.05

As we said at the start of this section, Bayes' rule is a pillar of probability and statistics. We have seen that Bayes' rule allows us to ‘invert’ conditional probabilities. When we learn statistics we will see that the art of statistical inference involves deciding how to proceed when one (or more) of the terms on the right side of Bayes' rule is unknown.

Discrete Random Variables

Class 4, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Know the definition of a discrete random variable.
2. Know the Bernoulli, binomial, and geometric distributions and examples of what they model.
3. Be able to describe the probability mass function and cumulative distribution function using tables and formulas.
4. Be able to construct new random variables from old ones.
5. Know how to compute expected value (mean).

2 Random Variables

This topic is largely about introducing some useful terminology, building on the notions of sample space and probability function. The key words are

1. Random variable
2. Probability mass function (pmf)
3. Cumulative distribution function (cdf)

2.1 Recap

A **discrete sample space** Ω is a finite or listable set of outcomes $\{\omega_1, \omega_2 \dots\}$. The probability of an outcome ω is denoted $P(\omega)$.

An **event** E is a subset of Ω . The **probability of an event** E is $P(E) = \sum_{\omega \in E} P(\omega)$.

2.2 Random variables as payoff functions

Example 1. A game with 2 dice.

Roll a die twice and record the outcomes as (i, j) , where i is the result of the first roll and j the result of the second. We can take the sample space to be

$$\Omega = \{(1, 1), (1, 2), (1, 3), \dots, (6, 6)\} = \{(i, j) \mid i, j = 1, \dots, 6\}.$$

The probability function is $P(i, j) = 1/36$.

In this game, you win \$500 if the sum is 7 and lose \$100 otherwise. We give this [payoff function](#) the name X and describe it formally by

$$X(i, j) = \begin{cases} 500 & \text{if } i + j = 7 \\ -100 & \text{if } i + j \neq 7. \end{cases}$$

Example 2. We can change the game by using a different payoff function. For example

$$Y(i, j) = ij - 10.$$

In this example if you roll (6, 2) then you win \$2. If you roll (2, 3) then you win -\$4 (i.e., lose \$4).

Question: Which game is the better bet?

answer: We will come back to this once we learn about expectation.

These payoff functions are examples of random variables. A [random variable assigns a number to each outcome in a sample space](#). More formally:

Definition: Let Ω be a sample space. A [discrete random variable](#) is a function

$$X : \Omega \rightarrow \mathbf{R}$$

that takes a discrete set of values. (Recall that \mathbf{R} stands for the real numbers.)

Why is X called a random variable? It's 'random' because its value depends on a random outcome of an experiment. And we treat X like we would a usual variable: we can add it to other random variables, square it, and so on.

2.3 Events and random variables

For any value a we write $X = a$ to mean the [event](#) consisting of all outcomes ω with $X(\omega) = a$.

Example 3. In Example 1 we rolled two dice and X was the random variable

$$X(i, j) = \begin{cases} 500 & \text{if } i + j = 7 \\ -100 & \text{if } i + j \neq 7. \end{cases}$$

The [event](#) $X = 500$ is the set $\{(1,6), (2,5), (3,4), (4,3), (5,2), (6,1)\}$, i.e. the set of all outcomes that sum to 7. So $P(X = 500) = 1/6$.

We allow a to be any value, even values that X never takes. In Example 1, we could look at the event $X = 1000$. Since X never equals 1000 this is just the [empty event](#) (or empty set)

$$'X = 1000' = \{\} = \emptyset \quad P(X = 1000) = 0.$$

2.4 Probability mass function and cumulative distribution function

It gets tiring and hard to read and write $P(X = a)$ for the probability that $X = a$. When we know we're talking about X we will simply write $p(a)$. If we want to make X explicit we will write $p_X(a)$. We spell this out in a definition.

Definition: The **probability mass function (pmf)** of a discrete random variable is the function $p(a) = P(X = a)$.

Note:

1. We always have $0 \leq p(a) \leq 1$.
2. We allow a to be any number. If a is a value that X never takes, then $p(a) = 0$.

Example 4. Let Ω be our earlier sample space for rolling 2 dice. Define the random variable M to be the **maximum value of the two dice**:

$$M(i, j) = \max(i, j).$$

For example, the roll (3,5) has maximum 5, i.e. $M(3, 5) = 5$.

We can describe a random variable by listing its possible values and the probabilities associated to these values. For the above example we have:

value	$a:$	1	2	3	4	5	6
pmf	$p(a):$	1/36	3/36	5/36	7/36	9/36	11/36

For example, $p(2) = 3/36$.

Question: What is $p(8)$? **answer:** $p(8) = 0$.

Think: What is the pmf for $Z(i, j) = i + j$? Does it look familiar?

2.5 Events and inequalities

Inequalities with random variables describe events. For example $X \leq a$ is the set of all outcomes ω such that $X(\omega) \leq a$.

Example 5. If our sample space is the set of all pairs of (i, j) coming from rolling two dice and $Z(i, j) = i + j$ is the sum of the dice then

$$Z \leq 4 = \{(1, 1), (1, 2), (1, 3), (2, 1), (2, 2), (3, 1)\}$$

2.6 The cumulative distribution function (cdf)

Definition: The **cumulative distribution function (cdf)** of a random variable X is the function F given by $F(a) = P(X \leq a)$. We will often shorten this to **distribution function**.

Note well that the definition of $F(a)$ uses the symbol less than or equal. This will be important for getting your calculations exactly right.

Example. Continuing with the example M , we have

value	$a:$	1	2	3	4	5	6
pmf	$p(a):$	1/36	3/36	5/36	7/36	9/36	11/36
cdf	$F(a):$	1/36	4/36	9/36	16/36	25/36	36/36

$F(a)$ is called the **cumulative** distribution function because $F(a)$ gives the total probability that accumulates by adding up the probabilities $p(b)$ as b runs from $-\infty$ to a . For example, in the table above, the entry $16/36$ in column 4 for the cdf is the sum of the values of the pmf from column 1 to column 4. In notation:

As events: ' $M \leq 4$ ' = {1, 2, 3, 4}; $F(4) = P(M \leq 4) = 1/36 + 3/36 + 5/36 + 7/36 = 16/36$.

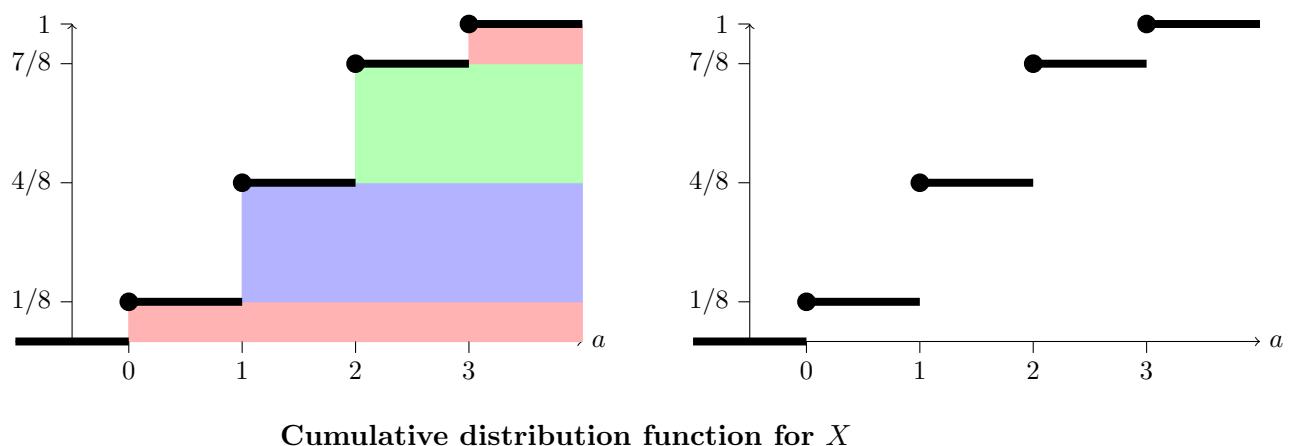
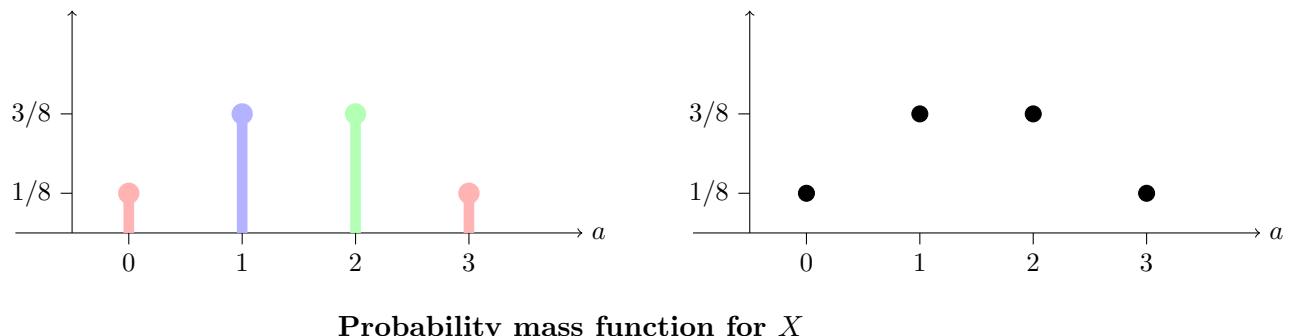
Just like the probability mass function, $F(a)$ is defined for all values a . In the above example, $F(8) = 1$, $F(-2) = 0$, $F(2.5) = 4/36$, and $F(\pi) = 9/36$.

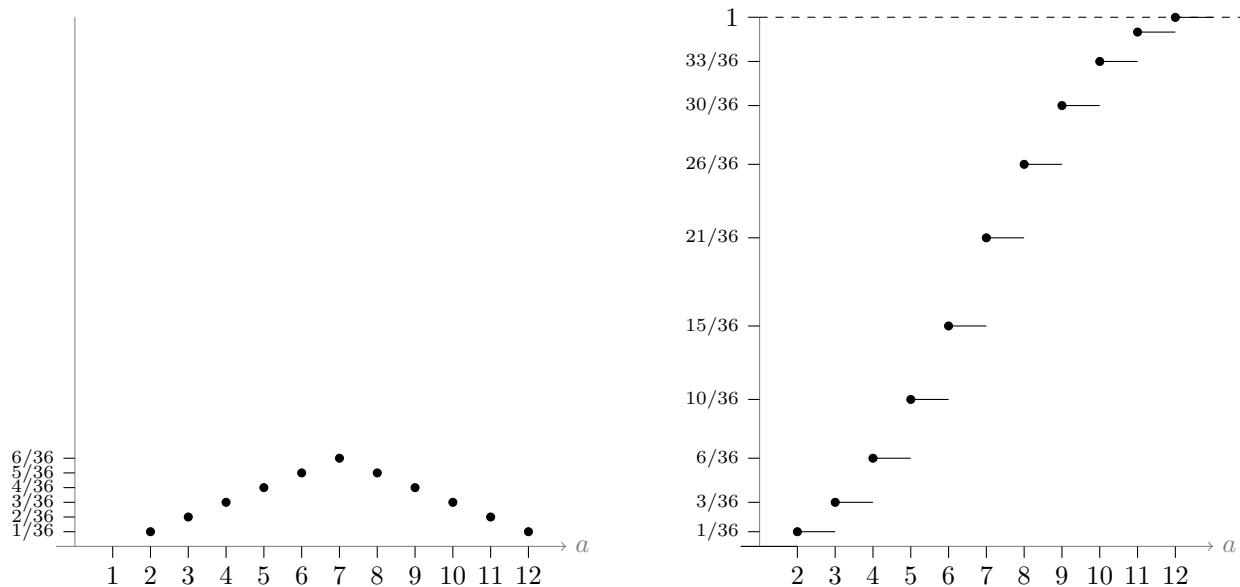
2.7 Graphs of $p(a)$ and $F(a)$

We can visualize the pmf and cdf with graphs. For example, let X be the number of heads in 3 tosses of a fair coin:

value a :	0	1	2	3
pmf $p(a)$:	$1/8$	$3/8$	$3/8$	$1/8$
cdf $F(a)$:	$1/8$	$4/8$	$7/8$	1

The colored graphs show how the cumulative distribution function is built by **accumulating** probability as a increases. The black and white graphs are the more standard presentations.





pmf and cdf for the maximum of two dice (Example 4)

Histograms: Later we will see another way to visualize the pmf using histograms. These require some care to do right, so we will wait until we need them.

2.8 Properties of the cdf F

The cdf F of a random variable satisfies several properties:

1. F is **non-decreasing**. That is, its graph never goes down, or symbolically if $a \leq b$ then $F(a) \leq F(b)$.
2. $0 \leq F(a) \leq 1$.
3. $\lim_{a \rightarrow \infty} F(a) = 1$, $\lim_{a \rightarrow -\infty} F(a) = 0$.

In words, (1) says the cumulative probability $F(a)$ increases or remains constant as a increases, but never decreases; (2) says the accumulated probability is always between 0 and 1; (3) says that as a gets very large, it becomes more and more certain that $X \leq a$ and as a gets very negative it becomes more and more certain that $X > a$.

Think: Why does a cdf satisfy each of these properties?

3 Specific Distributions

3.1 Bernoulli Distributions

Model: The Bernoulli distribution models one trial in an experiment that can result in either **success** or **failure**. This is the most important distribution and is also the simplest. A random variable X has a **Bernoulli distribution** with parameter p if:

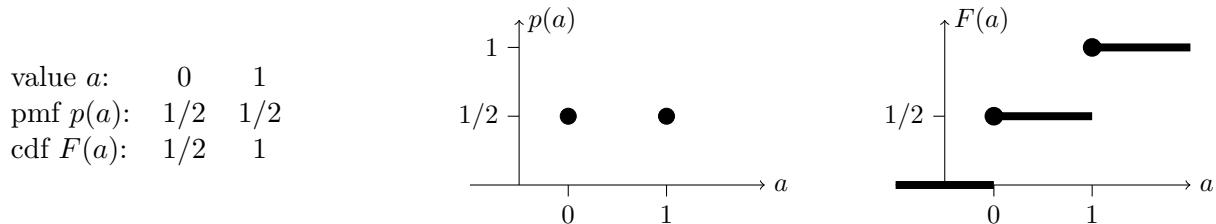
1. X takes the values 0 and 1.
2. $P(X = 1) = p$ and $P(X = 0) = 1 - p$.

We will write $X \sim \text{Bernoulli}(p)$ or $\text{Ber}(p)$, which is read “ X follows a Bernoulli distribution with parameter p ” or “ X is drawn from a Bernoulli distribution with parameter p ”.

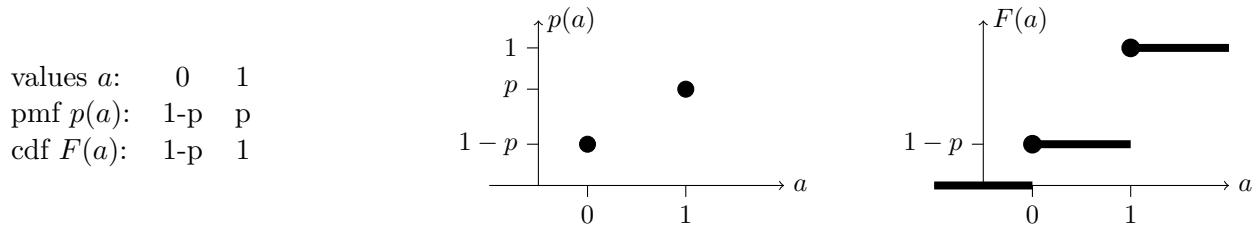
A simple model for the Bernoulli distribution is to flip a coin with probability p of heads, with $X = 1$ on heads and $X = 0$ on tails. The general terminology is to say X is 1 on **success** and 0 on **failure**, with success and failure defined by the context.

Many decisions can be modeled as a binary choice, such as votes for or against a proposal. If p is the proportion of the voting population that favors the proposal, than the vote of a random individual is modeled by a $\text{Bernoulli}(p)$.

Here are the table and graphs of the pmf and cdf for the $\text{Bernoulli}(1/2)$ distribution and below that for the general $\text{Bernoulli}(p)$ distribution.



Table, pmf and cmf for the $\text{Bernoulli}(1/2)$ distribution



Table, pmf and cmf for the $\text{Bernoulli}(p)$ distribution

3.2 Binomial Distributions

The **binomial distribution** $\text{Binomial}(n,p)$, or $\text{Bin}(n,p)$, models the number of successes in n independent $\text{Bernoulli}(p)$ trials.

There is a hierarchy here. A single Bernoulli trial is, say, one toss of a coin. A single binomial trial consists of n Bernoulli trials. For coin flips the sample space for a Bernoulli trial is $\{H, T\}$. The sample space for a binomial trial is all **sequences** of heads and tails of length n . Likewise a Bernoulli random variable takes values 0 and 1 and a binomial random variables takes values $0, 1, 2, \dots, n$.

Example 6. $\text{Binomial}(1,p)$ is the same as $\text{Bernoulli}(p)$.

Example 7. The number of heads in n flips of a coin with probability p of heads follows a $\text{Binomial}(n, p)$ distribution.

We describe $X \sim \text{Binomial}(n, p)$ by giving its values and probabilities. For notation we will use k to mean an arbitrary number between 0 and n .

We remind you that ‘ n choose k ’ = $\binom{n}{k} = {}_n C_k$ is the number of ways to choose k things out of a collection of n things and it has the formula

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}. \quad (1)$$

(It is also called a **binomial coefficient**.) Here is a table for the pmf of a $\text{Binomial}(n, k)$ random variable. We will explain how the binomial coefficients enter the pmf for the binomial distribution after a simple example.

values $a:$	0	1	2	\dots	k	\dots	n
pmf $p(a):$	$(1-p)^n$	$\binom{n}{1} p^1 (1-p)^{n-1}$	$\binom{n}{2} p^2 (1-p)^{n-2}$	\dots	$\binom{n}{k} p^k (1-p)^{n-k}$	\dots	p^n

Example 8. What is the probability of 3 or more heads in 5 tosses of a fair coin?

answer: The binomial coefficients associated with $n = 5$ are

$$\binom{5}{0} = 1, \quad \binom{5}{1} = \frac{5!}{1!4!} = \frac{5 \cdot 4 \cdot 3 \cdot 2 \cdot 1}{4 \cdot 3 \cdot 2 \cdot 1} = 5, \quad \binom{5}{2} = \frac{5!}{2!3!} = \frac{5 \cdot 4 \cdot 3 \cdot 2 \cdot 1}{2 \cdot 1 \cdot 3 \cdot 2 \cdot 1} = \frac{5 \cdot 4}{2} = 10,$$

and similarly

$$\binom{5}{3} = 10, \quad \binom{5}{4} = 5, \quad \binom{5}{5} = 1.$$

Using these values we get the following table for $X \sim \text{Binomial}(5, p)$.

values $a:$	0	1	2	3	4	5
pmf $p(a):$	$(1-p)^5$	$5p(1-p)^4$	$10p^2(1-p)^3$	$10p^3(1-p)^2$	$5p^4(1-p)$	p^5

We were told $p = 1/2$ so

$$P(X \geq 3) = 10 \left(\frac{1}{2}\right)^3 \left(\frac{1}{2}\right)^2 + 5 \left(\frac{1}{2}\right)^4 \left(\frac{1}{2}\right)^1 + \left(\frac{1}{2}\right)^5 = \frac{16}{32} = \frac{1}{2}.$$

Think: Why is the value of $1/2$ not surprising?

3.3 Explanation of the binomial probabilities

For concreteness, let $n = 5$ and $k = 2$ (the argument for arbitrary n and k is identical.) So $X \sim \text{binomial}(5, p)$ and we want to compute $p(2)$. The long way to compute $p(2)$ is to list all the ways to get exactly 2 heads in 5 coin flips and add up their probabilities. The list has 10 entries:

HHTTT, HTHTT, HTTHT, HTTTH, THHTT, THTHT, THTTH, TTHHT, TTHTH, TTTHH

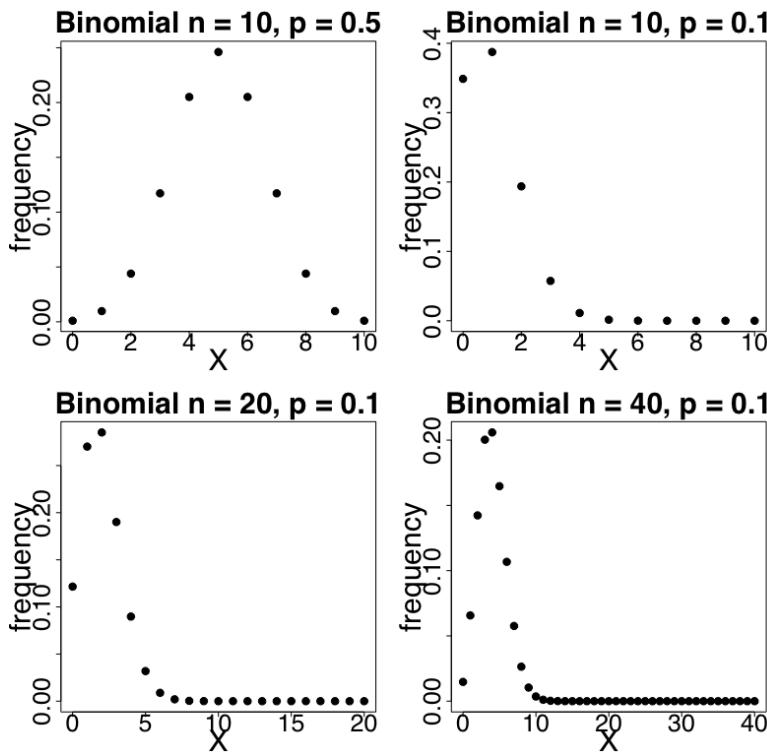
Each entry has the same probability of occurring, namely

$$p^2(1 - p)^3.$$

This is because each of the two heads has probability p and each of the 3 tails has probability $1 - p$. Because the individual tosses are independent we can multiply probabilities. Therefore, the total probability of exactly 2 heads is the sum of 10 identical probabilities, i.e. $p(2) = 10p^2(1 - p)^3$, as shown in the table.

This guides us to the shorter way to do the computation. We have to count the number of sequences with exactly 2 heads. To do this we need to choose 2 of the tosses to be heads and the remaining 3 to be tails. The number of such sequences is the number of ways to choose 2 out of 5 things, that is $\binom{5}{2}$. Since each such sequence has the same probability, $p^2(1 - p)^3$, we get the probability of exactly 2 heads $p(2) = \binom{5}{2}p^2(1 - p)^3$.

Here are some binomial probability mass function (here, frequency is the same as probability).



3.4 Geometric Distributions

A [geometric distribution](#) models the number of tails before the first head in a sequence of coin flips (Bernoulli trials).

Example 9. (a) Flip a coin repeatedly. Let X be the number of tails before the first heads. So, X can equal 0, i.e. the first flip is heads, 1, 2, In principle it take any nonnegative integer value.

(b) Give a flip of tails the value 0, and heads the value 1. In this case, X is the number of 0's before the first 1.

(c) Give a flip of tails the value 1, and heads the value 0. In this case, X is the number of 1's before the first 0.

(d) Call a flip of tails a success and heads a failure. So, X is the number of successes before the first failure.

(e) Call a flip of tails a failure and heads a success. So, X is the number of failures before the first success.

You can see this models many different scenarios of this type. The most neutral language is the number of tails before the first head.

Formal definition. The random variable X follows a [geometric distribution with parameter \$p\$](#) if

- X takes the values 0, 1, 2, 3, ...
- its pmf is given by $p(k) = P(X = k) = (1 - p)^k p$.

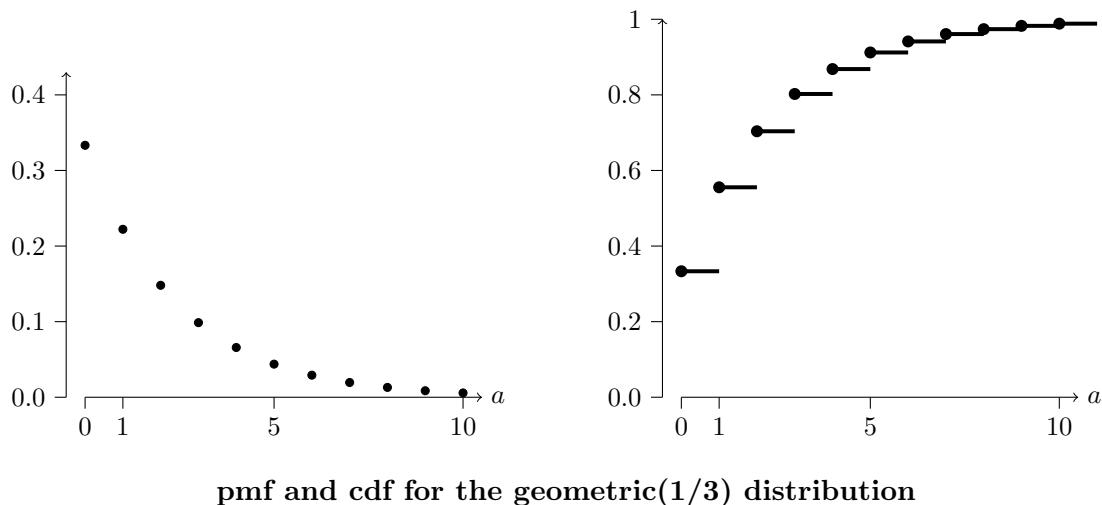
We denote this by $X \sim \text{geometric}(p)$ or $\text{geo}(p)$. In table form we have:

value	$a:$	0	1	2	3	...	k	...
pmf	$p(a):$	p	$(1 - p)p$	$(1 - p)^2 p$	$(1 - p)^3 p$...	$(1 - p)^k p$...

Table: $X \sim \text{geometric}(p)$: X = the number of 0s before the first 1.

We will show how this table was computed in an example below.

The geometric distribution is an example of a discrete distribution that takes an infinite number of possible values. Things can get confusing when we work with successes and failure since we might want to model the number of successes before the first failure or we might want the number of failures before the first success. To keep straight things straight you can translate to the neutral language of the number of tails before the first heads.



pmf and cdf for the $\text{geometric}(1/3)$ distribution

Example 10. Computing geometric probabilities. Suppose that the inhabitants of an island plan their families by having babies until the first girl is born. Assume the probability of having a girl with each pregnancy is 0.5 independent of other pregnancies, that all babies survive and there are no multiple births. What is the probability that a family has k boys?

answer: In neutral language we can think of boys as tails and girls as heads. Then the number of boys in a family is the number of tails before the first heads.

Let's practice using standard notation to present this. So, let X be the number of boys in a (randomly-chosen) family. So, X is a geometric random variable. We are asked to find $p(k) = P(X = k)$. A family has k boys if the sequence of children in the family from oldest to youngest is

$$BBB\dots BG$$

with the first k children being boys. The probability of this sequence is just the product of the probability for each child, i.e. $(1/2)^k \cdot (1/2) = (1/2)^{k+1}$. (Note: The assumptions of equal probability and independence are simplifications of reality.)

Think: What is the ratio of boys to girls on the island?

More geometric confusion. Another common definition for the geometric distribution is the number of tosses until the first heads. In this case X can take the values 1, i.e. the first flip is heads, 2, 3, This is just our geometric random variable plus 1. The methods of computing with it are just like the ones we used above.

3.5 Uniform Distribution

The uniform distribution models any situation where all the outcomes are equally likely.

$$X \sim \text{uniform}(N).$$

X takes values $1, 2, 3, \dots, N$, each with probability $1/N$. We have already seen this distribution many times when modeling to fair coins ($N = 2$), dice ($N = 6$), birthdays ($N = 365$), and poker hands ($N = \binom{52}{5}$).

3.6 Discrete Distributions Applet

The applet at <http://mathlets.org/mathlets/probability-distributions/> gives a dynamic view of some discrete distributions. The graphs will change smoothly as you move the various sliders. Try playing with the different distributions and parameters.

This applet is carefully color-coded. Two things with the same color represent the same or closely related notions. By understanding the color-coding and other details of the applet, you will acquire a stronger intuition for the distributions shown.

3.7 Other Distributions

There are a million other named distributions arising in various contexts. We don't expect you to memorize them (we certainly have not!), but you should be comfortable using a resource like Wikipedia to look up a pmf. For example, take a look at the info box at the top right of http://en.wikipedia.org/wiki/Hypergeometric_distribution. The info box lists many (surely unfamiliar) properties in addition to the pmf.

4 Arithmetic with Random Variables

We can do arithmetic with random variables. For example, we can add, subtract, multiply or square them.

There is a simple, but [extremely important](#) idea for counting. It says that if we have a sequence of numbers that are either 0 or 1 then the sum of the sequence is the number of 1s.

Example 11. Consider the sequence with five 1s

$$1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0.$$

It is easy to see that the sum of this sequence is 5 the number of 1s.

We illustrate this idea by counting the number of heads in n tosses of a coin.

Example 12. Toss a fair coin n times. Let X_j be 1 if the j th toss is heads and 0 if it's tails. So, X_j is a Bernoulli($1/2$) random variable. Let X be the total number of heads in the n tosses. Assuming the tosses are independent we know $X \sim \text{binomial}(n, 1/2)$. We can also write

$$X = X_1 + X_2 + X_3 + \dots + X_n.$$

Again, this is because the terms in the sum on the right are all either 0 or 1. So, the sum is exactly the number of X_j that are 1, i.e. the number of heads.

The important thing to see in the example above is that we've written the more complicated binomial random variable X as the sum of extremely simple random variables X_j . This will allow us to manipulate X algebraically.

Think: Suppose X and Y are independent and $X \sim \text{binomial}(n, 1/2)$ and $Y \sim \text{binomial}(m, 1/2)$. What kind of distribution does $X + Y$ follow? (**Answer:** $\text{binomial}(n + m, 1/2)$). Why?)

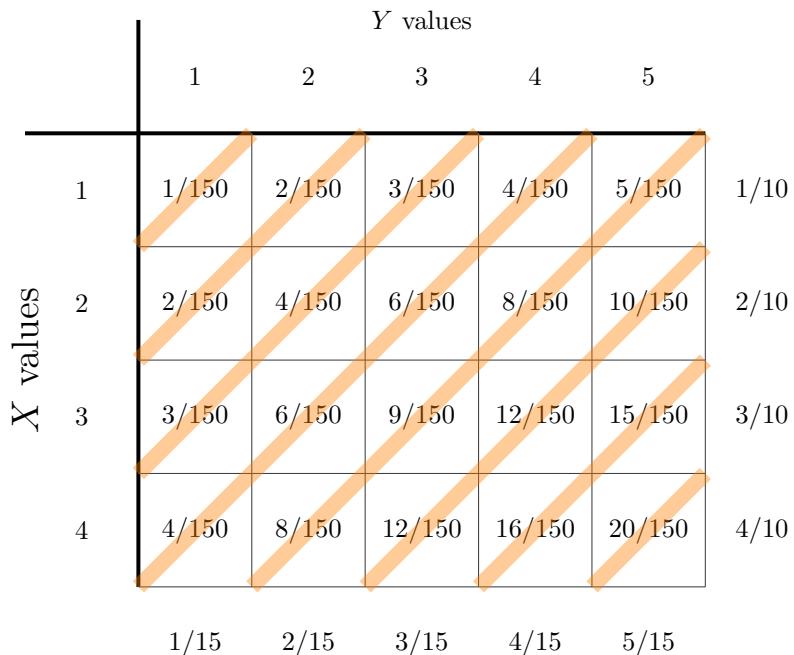
Example 13. Suppose X and Y are independent random variables with the following tables.

Values of X	$x:$	1	2	3	4
pmf	$p_X(x):$	1/10	2/10	3/10	4/10

Values of Y	$y:$	1	2	3	4	5
pmf	$p_Y(y):$	1/15	2/15	3/15	4/15	5/15

Check that the total probability for each random variable is 1. Make a table for the random variable $X + Y$.

answer: The first thing to do is make a two-dimensional table for the product sample space consisting of pairs (x, y) , where x is a possible value of X and y one of Y . To help do the computation, the probabilities for the X values are put in the far right column and those for Y are in the bottom row. Because X and Y are independent the probability for (x, y) pair is just the product of the individual probabilities.



The diagonal stripes show sets of squares where $X + Y$ is the same. All we have to do to compute the probability table for $X + Y$ is sum the probabilities for each stripe.

$$\begin{array}{ll}
 X + Y \text{ values:} & 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8 \quad 9 \\
 \text{pmf:} & 1/150 \quad 4/150 \quad 10/150 \quad 20/150 \quad 30/150 \quad 34/150 \quad 31/150 \quad 20/150
 \end{array}$$

When the tables are too big to write down we'll need to use purely algebraic techniques to compute the probabilities of a sum. We will learn how to do this in due course.

Discrete Random Variables: Expected Value

Class 4, 18.05

Jeremy Orloff and Jonathan Bloom

1 Expected Value

In the R reading questions for this lecture, you simulated the average value of rolling a die many times. You should have gotten a value close to the exact answer of 3.5. To motivate the formal definition of the average, or **expected value**, we first consider some examples.

Example 1. Suppose we have a six-sided die marked with five 3's and one 6. (This was the red one from our non-transitive dice.) What would you expect the average of 6000 rolls to be?

answer: If we knew the value of each roll, we could compute the average by summing the 6000 values and dividing by 6000. Without knowing the values, we can compute the **expected average** as follows.

Since there are five 3's and one six we expect roughly 5/6 of the rolls will give 3 and 1/6 will give 6. Assuming this to be exactly true, we have the following table of values and counts:

value:	3	6
expected counts:	5000	1000

The average of these 6000 values is then

$$\frac{5000 \cdot 3 + 1000 \cdot 6}{6000} = \frac{5}{6} \cdot 3 + \frac{1}{6} \cdot 6 = 3.5$$

We consider this the expected average in the sense that we ‘expect’ each of the possible values to occur with the given frequencies.

Example 2. We roll two standard 6-sided dice. You win \$1000 if the sum is 2 and lose \$100 otherwise. How much do you expect to win on average per trial?

answer: The probability of a 2 is 1/36. If you play N times, you can ‘expect’ $\frac{1}{36} \cdot N$ of the trials to give a 2 and $\frac{35}{36} \cdot N$ of the trials to give something else. Thus your total expected winnings are

$$1000 \cdot \frac{N}{36} - 100 \cdot \frac{35N}{36}.$$

To get the expected average per trial we divide the total by N :

$$\text{expected average} = 1000 \cdot \frac{1}{36} - 100 \cdot \frac{35}{36} = -69.44.$$

Think: Would you be willing to play this game one time? Multiple times?

Notice that in both examples the sum for the expected average consists of terms which are a value of the random variable times its probability. This leads to the following definition.

Definition: Suppose X is a discrete random variable that takes values x_1, x_2, \dots, x_n with probabilities $p(x_1), p(x_2), \dots, p(x_n)$. The **expected value** of X is denoted $E(X)$ and defined

From physics we know the center of mass is

$$\bar{x} = \frac{m_1 x_1 + m_2 x_2}{m_1 + m_2} = \frac{500 \cdot 3 + 100 \cdot 6}{600} = 3.5.$$

We call this formula a ‘weighted’ average of the x_1 and x_2 . Here x_1 is weighted more heavily because it has more mass.

Now look at the definition of expected value $E(X)$. It is a weighted average of the values of X with the weights being probabilities $p(x_i)$ rather than masses! We might say that “The expected value is the point at which the distribution would balance”. Note the similarity between the physics example and Example 1.

1.2 Algebraic properties of $E(X)$

When we add, scale or shift random variables the expected values do the same. The shorthand mathematical way of saying this is that $E(X)$ is [linear](#).

1. If X and Y are random variables on a sample space Ω then

$$E(X + Y) = E(X) + E(Y)$$

2. If a and b are constants then

$$E(aX + b) = aE(X) + b.$$

We will think of $aX + b$ as [scaling](#) X by a and [shifting](#) it by b .

Before proving these properties, let’s consider a few examples.

Example 6. Roll two dice and let X be the sum. Find $E(X)$.

answer: Let X_1 be the value on the first die and let X_2 be the value on the second die. Since $X = X_1 + X_2$ we have $E(X) = E(X_1) + E(X_2)$. Earlier we computed that $E(X_1) = E(X_2) = 3.5$, therefore $E(X) = 7$.

Example 7. Let $X \sim \text{binomial}(n, p)$. Find $E(X)$.

answer: Recall that X models the number of successes in n Bernoulli(p) random variables, which we’ll call X_1, \dots, X_n . The key fact, which we highlighted in the previous reading for this class, is that

$$X = \sum_{j=1}^n X_j.$$

Now we can use the Algebraic Property (1) to make the calculation simple.

$$X = \sum_{j=1}^n X_j \Rightarrow E(X) = \sum_j E(X_j) = \sum_j p = \boxed{np}.$$

We could have computed $E(X)$ directly as

$$E(X) = \sum_{k=0}^n kp(k) = \sum_{k=0}^n k \binom{n}{k} p^k (1-p)^{n-k}.$$

It is possible to show that the sum of this series is indeed np . We think you'll agree that the method using Property (1) is much easier.

Example 8. (For infinite random variables the mean does not always exist.) Suppose X has an infinite number of values according to the following table

$$\begin{array}{ccccccc} \text{values } x: & 2 & 2^2 & 2^3 & \dots & 2^k & \dots \\ \text{pmf } p(x): & 1/2 & 1/2^2 & 1/2^3 & \dots & 1/2^k & \dots \end{array} \quad \text{Try to compute the mean.}$$

answer: The mean is

$$E(X) = \sum_{k=1}^{\infty} 2^k \frac{1}{2^k} = \sum_{k=1}^{\infty} 1 = \infty.$$

The mean does not exist! This can happen with infinite series.

1.3 Proofs of the algebraic properties of $E(X)$

The proof of Property (1) is simple, but there is some subtlety in even understanding what it means to add two random variables. Recall that the value of random variable is a number determined by the outcome of an experiment. To add X and Y means to add the values of X and Y for the same outcome. In table form this looks like:

$$\begin{array}{cccccc} \text{outcome } \omega: & \omega_1 & \omega_2 & \omega_3 & \dots & \omega_n \\ \text{value of } X: & x_1 & x_2 & x_3 & \dots & x_n \\ \text{value of } Y: & y_1 & y_2 & y_3 & \dots & y_n \\ \text{value of } X + Y: & x_1 + y_1 & x_2 + y_2 & x_3 + y_3 & \dots & x_n + y_n \\ \text{prob. } P(\omega): & P(\omega_1) & P(\omega_2) & P(\omega_3) & \dots & P(\omega_n) \end{array}$$

The proof of (1) follows immediately:

$$E(X + Y) = \sum (x_i + y_i)P(\omega_i) = \sum x_i P(\omega_i) + \sum y_i P(\omega_i) = E(X) + E(Y).$$

The proof of Property (2) only takes one line.

$$E(aX + b) = \sum p(x_i)(ax_i + b) = a \sum p(x_i)x_i + b \sum p(x_i) = aE(X) + b.$$

The b term in the last expression follows because $\sum p(x_i) = 1$.

Example 9. Mean of a geometric distribution

Let $X \sim \text{geo}(p)$. Recall this means X takes values $k = 0, 1, 2, \dots$ with probabilities $p(k) = (1-p)^k p$. (X models the number of tails before the first heads in a sequence of Bernoulli trials.) The mean is given by

$$E(X) = \frac{1-p}{p}.$$

To see this requires a clever trick. Mathematicians love this sort of thing and we hope you are able to follow the logic. In this class we will not ask you to come up with something like this on an exam.

Here's the trick.: to compute $E(X)$ we have to sum the infinite series

$$E(X) = \sum_{k=0}^{\infty} k(1-p)^k p.$$

Here is the trick. We know the sum of the geometric series: $\sum_{k=0}^{\infty} x^k = \frac{1}{1-x}$.

Differentiate both sides: $\sum_{k=0}^{\infty} kx^{k-1} = \frac{1}{(1-x)^2}$.

Multiply by x : $\sum_{k=0}^{\infty} kx^k = \frac{x}{(1-x)^2}$.

Replace x by $1-p$: $\sum_{k=0}^{\infty} k(1-p)^k = \frac{1-p}{p^2}$.

Multiply by p : $\sum_{k=0}^{\infty} k(1-p)^k p = \frac{1-p}{p}$.

This last expression is the mean.

$$E(X) = \frac{1-p}{p}.$$

Example 10. Flip a fair coin until you get heads for the first time. What is the expected number of times you flipped tails?

answer: The number of tails before the first head is modeled by $X \sim \text{geo}(1/2)$. From the previous example $E(X) = \frac{1/2}{1/2} = 1$. This is a surprisingly small number.

Example 11. Michael Jordan, the greatest basketball player ever, made 80% of his free throws. In a game what is the expected number he would make before his first miss.

answer: Here is an example where we want the number of successes before the first failure. Using the neutral language of heads and tails: success is tails (probability $1-p$) and failure is heads (probability $= p$). Therefore $p = .2$ and the number of tails (made free throws) before the first heads (missed free throw) is modeled by a $X \sim \text{geo}(.2)$. We saw in Example 9 that this is

$$E(X) = \frac{1-p}{p} = \frac{.8}{.2} = 4.$$

1.4 Expected values of functions of a random variable

(The change of variables formula.)

If X is a discrete random variable taking values x_1, x_2, \dots and h is a function the $h(X)$ is a new random variable. Its expected value is

$$E(h(X)) = \sum_j h(x_j)p(x_j).$$

We illustrate this with several examples.

Example 12. Let X be the value of a roll of one die and let $Y = X^2$. Find $E(Y)$.

answer: Since there are a small number of values we can make a table.

X	1	2	3	4	5	6
Y	1	4	9	16	25	36
prob	1/6	1/6	1/6	1/6	1/6	1/6

Notice the probability for each Y value is the same as that of the corresponding X value.
So,

$$E(Y) = E(X^2) = 1^2 \cdot \frac{1}{6} + 2^2 \cdot \frac{1}{6} + \dots + 6^2 \cdot \frac{1}{6} = 15.167.$$

Example 13. Roll two dice and let X be the sum. Suppose the payoff function is given by $Y = X^2 - 6X + 1$. Is this a good bet?

answer: We have $E(Y) = \sum_{j=2}^{12} (j^2 - 6j + 1)p(j)$, where $p(j) = P(X = j)$.

We show the table, but really we'll use R to do the calculation.

X	2	3	4	5	6	7	8	9	10	11	12
Y	-7	-8	-7	-4	1	8	17	28	41	56	73
prob	1/36	2/36	3/36	4/36	5/36	6/36	5/36	4/36	3/36	2/36	1/36

Here's the R code I used to compute $E(Y) = 13.833$.

```
x = 2:12
y = x^2 - 6*x + 1
p = c(1 2 3 4 5 6 5 4 3 2 1)/36
ave = sum(p*y)
```

It gave $\text{ave} = 13.833$.

To answer the question above: since the expected payoff is positive it looks like a bet worth taking.

Quiz: If $Y = h(X)$ does $E(Y) = h(E(X))$? **answer: NO!!!** This is not true in general!

Think: Is it true in the previous example?

Quiz: If $Y = 3X + 77$ does $E(Y) = 3E(X) + 77$?

answer: Yes. By property (2), scaling and shifting does behave like this.

Variance of Discrete Random Variables

Class 5, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

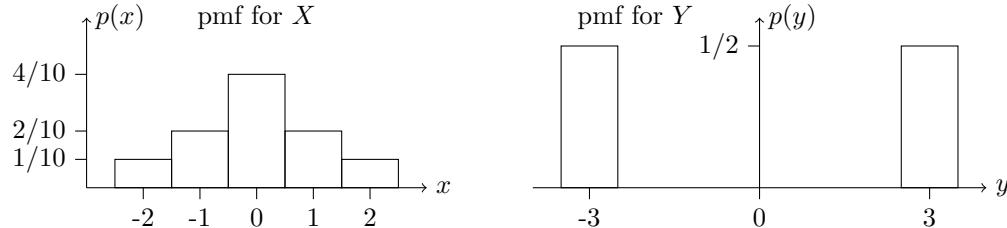
1. Be able to compute the variance and standard deviation of a random variable.
2. Understand that standard deviation is a measure of scale or spread.
3. Be able to compute variance using the properties of scaling and linearity.

2 Spread

The expected value (mean) of a random variable is a measure of [location or central tendency](#). If you had to summarize a random variable with a single number, the mean would be a good choice. Still, the mean leaves out a good deal of information. For example, the random variables X and Y below both have mean 0, but their probability mass is spread out about the mean quite differently.

values X	-2	-1	0	1	2	values Y	-3	3
pmf $p(x)$	1/10	2/10	4/10	2/10	1/10	pmf $p(y)$	1/2	1/2

It's probably a little easier to see the different spreads in plots of the probability mass functions. We use bars instead of dots to give a better sense of the mass.



pmf's for two different distributions both with mean 0

In the next section, we will learn how to quantify this spread.

3 Variance and standard deviation

Taking the mean as the center of a random variable's probability distribution, the [variance](#) is a measure of how much the probability mass is [spread](#) out around this center. We'll start with the formal definition of variance and then unpack its meaning.

Definition: If X is a random variable with mean $E(X) = \mu$, then the [variance](#) of X is defined by

$$\text{Var}(X) = E((X - \mu)^2).$$

The **standard deviation** σ of X is defined by

$$\sigma = \sqrt{\text{Var}(X)}.$$

If the relevant random variable is clear from context, then the variance and standard deviation are often denoted by σ^2 and σ ('sigma'), just as the mean is μ ('mu').

What does this mean? First, let's rewrite the definition explicitly as a sum. If X takes values x_1, x_2, \dots, x_n with probability mass function $p(x_i)$ then

$$\text{Var}(X) = E((X - \mu)^2) = \sum_{i=1}^n p(x_i)(x_i - \mu)^2.$$

In words, the formula for $\text{Var}(X)$ says to take a weighted average of the squared distance to the mean. By squaring, we make sure we are averaging only non-negative values, so that the spread to the right of the mean won't cancel that to the left. By using expectation, we are weighting high probability values more than low probability values. (See Example 2 below.)

Note on units:

1. σ has the same units as X .
2. $\text{Var}(X)$ has the same units as the square of X . So if X is in meters, then $\text{Var}(X)$ is in meters squared.

Because σ and X have the same units, the standard deviation is a natural measure of spread.

Let's work some examples to make the notion of variance clear.

Example 1. Compute the mean, variance and standard deviation of the random variable X with the following table of values and probabilities.

value x	1	3	5
pmf $p(x)$	1/4	1/4	1/2

answer: First we compute $E(X) = 7/2$. Then we extend the table to include $(X - 7/2)^2$.

value x	1	3	5
$p(x)$	1/4	1/4	1/2
$(x - 7/2)^2$	25/4	1/4	9/4

Now the computation of the variance is similar to that of expectation:

$$\text{Var}(X) = \frac{25}{4} \cdot \frac{1}{4} + \frac{1}{4} \cdot \frac{1}{4} + \frac{9}{4} \cdot \frac{1}{2} = \frac{11}{4}.$$

Taking the square root we have the standard deviation $\sigma = \sqrt{11/4}$.

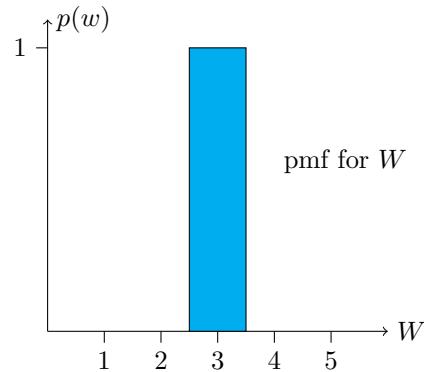
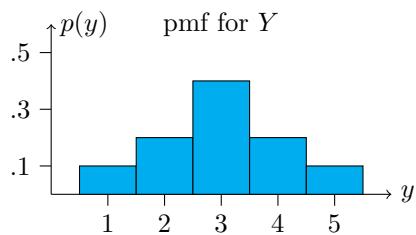
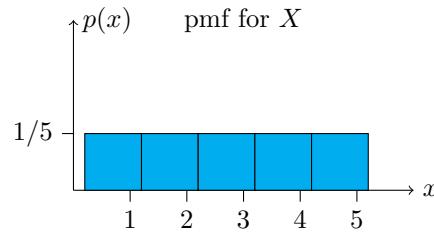
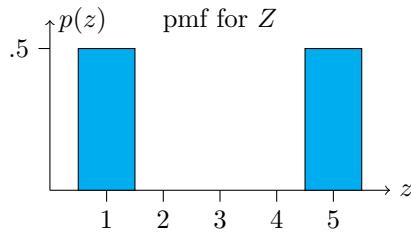
Example 2. For each random variable X , Y , Z , and W plot the pmf and compute the mean and variance.

(i)	value x	1	2	3	4	5
	pmf $p(x)$	1/5	1/5	1/5	1/5	1/5
(ii)	value y	1	2	3	4	5
	pmf $p(y)$	1/10	2/10	4/10	2/10	1/10

(iii)	value z	1	2	3	4	5
	pmf $p(z)$	5/10	0	0	0	5/10

(iv)	value w	1	2	3	4	5
	pmf $p(w)$	0	0	1	0	0

answer: Each random variable has the same mean 3, but the probability is spread out differently. In the plots below, we order the pmf's from largest to smallest variance: Z, X, Y, W .



Next we'll verify our visual intuition by computing the variance of each of the variables. All of them have mean $\mu = 3$. Since the variance is defined as an expected value, we can compute it using the tables.

(i)	value x	1	2	3	4	5
	pmf $p(x)$	1/5	1/5	1/5	1/5	1/5
	$(X - \mu)^2$	4	1	0	1	4

$$\text{Var}(X) = E((X - \mu)^2) = \frac{4}{5} + \frac{1}{5} + \frac{0}{5} + \frac{1}{5} + \frac{4}{5} = \boxed{2}.$$

(ii)	value y	1	2	3	4	5
	$p(y)$	1/10	2/10	4/10	2/10	1/10
	$(Y - \mu)^2$	4	1	0	1	4

$$\text{Var}(Y) = E((Y - \mu)^2) = \frac{4}{10} + \frac{2}{10} + \frac{0}{10} + \frac{2}{10} + \frac{4}{10} = \boxed{1.2}.$$

(iii)	value z	1	2	3	4	5
	pmf $p(z)$	5/10	0	0	0	5/10
	$(Z - \mu)^2$	4	1	0	1	4

$$\text{Var}(Z) = E((Z - \mu)^2) = \frac{20}{10} + \frac{20}{10} = \boxed{4}.$$

(iv)	value w	1	2	3	4	5
	pmf $p(w)$	0	0	1	0	0
	$(W - \mu)^2$	4	1	0	1	4

$\text{Var}(W) = \boxed{0}$. Note that W doesn't vary, so it has variance 0!

3.1 The variance of a $\text{Bernoulli}(p)$ random variable.

Bernoulli random variables are fundamental, so we should know their variance.

If $X \sim \text{Bernoulli}(p)$ then

$$\text{Var}(X) = p(1-p).$$

Proof: We know that $E(X) = p$. We compute $\text{Var}(X)$ using a table.

values X	0	1
pmf $p(x)$	$1-p$	p
$(X - \mu)^2$	$(0-p)^2$	$(1-p)^2$

$$\text{Var}(X) = (1-p)p^2 + p(1-p)^2 = (1-p)p(1-p+p) = \boxed{(1-p)p}.$$

As with all things Bernoulli, you should remember this formula.

Think: For what value of p does $\text{Bernoulli}(p)$ have the highest variance? Try to answer this by plotting the PMF for various p .

3.2 A word about independence

So far we have been using the notion of independent random variable without ever carefully defining it. For example, a binomial distribution is the sum of **independent** Bernoulli trials. This may (should?) have bothered you. Of course, we have an intuitive sense of what independence means for experimental trials. We also have the probabilistic sense that random variables X and Y are independent if knowing the value of X gives you no information about the value of Y .

In a few classes we will work with continuous random variables and joint probability functions. After that we will be ready for a full definition of independence. For now we can use the following definition, which is exactly what you expect and is valid for discrete random variables.

Definition: The discrete random variables X and Y are **independent** if

$$P(X = a, Y = b) = P(X = a)P(Y = b)$$

for any values a, b . That is, the probabilities multiply.

3.3 Properties of variance

The three most useful properties for computing variance are:

1. If X and Y are **independent** then $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$.

2. For constants a and b , $\text{Var}(aX + b) = a^2\text{Var}(X)$.
3. $\text{Var}(X) = E(X^2) - E(X)^2$.

For Property 1, note carefully the requirement that X and Y are independent. We will return to the proof of Property 1 in a later class.

Property 3 gives a formula for $\text{Var}(X)$ that is often easier to use in hand calculations. The computer is happy to use the definition! We'll prove Properties 2 and 3 after some examples.

Example 3. Suppose X and Y are independent and $\text{Var}(X) = 3$ and $\text{Var}(Y) = 5$. Find:

- (i) $\text{Var}(X + Y)$,
- (ii) $\text{Var}(3X + 4)$,
- (iii) $\text{Var}(X + X)$,
- (iv) $\text{Var}(X + 3Y)$.

answer: To compute these variances we make use of Properties 1 and 2.

(i) Since X and Y are **independent**, $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) = 8$.

(ii) Using Property 2, $\text{Var}(3X + 4) = 9 \cdot \text{Var}(X) = 27$.

(iii) Don't be fooled! Property 1 fails since X is certainly not independent of itself. We can use Property 2: $\text{Var}(X + X) = \text{Var}(2X) = 4 \cdot \text{Var}(X) = 12$. (Note: if we mistakenly used Property 1, we would get the wrong answer of 6.)

(iv) We use both Properties 1 and 2.

$$\text{Var}(X + 3Y) = \text{Var}(X) + \text{Var}(3Y) = 3 + 9 \cdot 5 = 48.$$

Example 4. Use Property 3 to compute the variance of $X \sim \text{Bernoulli}(p)$.

answer: From the table

X	0	1
$p(x)$	$1-p$	p
X^2	0	1

we have $E(X^2) = p$. So Property 3 gives

$$\text{Var}(X) = E(X^2) - E(X)^2 = p - p^2 = p(1-p).$$

This agrees with our earlier calculation.

Example 5. Redo Example 1 using Property 3.

answer: From the table

X	1	3	5
$p(x)$	$1/4$	$1/4$	$1/2$
X^2	1	9	25

we have $E(X) = 7/2$ and

$$E(X^2) = 1^2 \cdot \frac{1}{4} + 3^2 \cdot \frac{1}{4} + 5^2 \cdot \frac{1}{2} = \frac{60}{4} = 15.$$

So $\text{Var}(X) = 15 - (7/2)^2 = 11/4$ –as before in Example 1.

3.4 Variance of binomial(n,p)

Suppose $X \sim \text{binomial}(n, p)$. Since X is the sum of *independent* Bernoulli(p) variables and each Bernoulli variable has variance $p(1-p)$ we have

$$X \sim \text{binomial}(n, p) \Rightarrow \text{Var}(X) = np(1-p).$$

3.5 Proof of properties 2 and 3

Proof of Property 2: This follows from the properties of $E(X)$ and some algebra.

Let $\mu = E(X)$. Then $E(aX + b) = a\mu + b$ and

$$\text{Var}(aX+b) = E((aX+b-(a\mu+b))^2) = E((aX-a\mu)^2) = E(a^2(X-\mu)^2) = a^2E((X-\mu)^2) = a^2\text{Var}(X).$$

Proof of Property 3: We use the properties of $E(X)$ and a bit of algebra. Remember that μ is a constant and that $E(X) = \mu$.

$$\begin{aligned} E((X - \mu)^2) &= E(X^2 - 2\mu X + \mu^2) \\ &= E(X^2) - 2\mu E(X) + \mu^2 \\ &= E(X^2) - 2\mu^2 + \mu^2 \\ &= E(X^2) - \mu^2 \\ &= E(X^2) - E(X)^2. \quad \text{QED} \end{aligned}$$

4 Tables of Distributions and Properties

Distribution	range X	pmf $p(x)$	mean $E(X)$	variance $\text{Var}(X)$
Bernoulli(p)	0, 1	$p(0) = 1 - p, \quad p(1) = p$	p	$p(1 - p)$
Binomial(n, p)	0, 1, ..., n	$p(k) = \binom{n}{k} p^k (1-p)^{n-k}$	np	$np(1 - p)$
Uniform(n)	1, 2, ..., n	$p(k) = \frac{1}{n}$	$\frac{n+1}{2}$	$\frac{n^2 - 1}{12}$
Geometric(p)	0, 1, 2, ...	$p(k) = p(1-p)^k$	$\frac{1-p}{p}$	$\frac{1-p}{p^2}$

Let X be a discrete random variable with range x_1, x_2, \dots and pmf $p(x_j)$.

Expected Value:	Variance:
Synonyms: mean, average	
Notation: $E(X), \mu$	$\text{Var}(X), \sigma^2$
Definition: $E(X) = \sum_j p(x_j)x_j$	$E((X - \mu)^2) = \sum_j p(x_j)(x_j - \mu)^2$
Scale and shift: $E(aX + b) = aE(X) + b$	$\text{Var}(aX + b) = a^2\text{Var}(X)$
Linearity: (for any X, Y)	(for X, Y independent)
	$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$
Functions of X : $E(h(X)) = \sum p(x_j) h(x_j)$	
Alternative formula:	$\text{Var}(X) = E(X^2) - E(X)^2 = E(X^2) - \mu^2$

Continuous Random Variables

Class 5, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Know the definition of a continuous random variable.
2. Know the definition of the probability density function (pdf) and cumulative distribution function (cdf).
3. Be able to explain why we use probability density for continuous random variables.

2 Introduction

We now turn to **continuous random variables**. All random variables assign a number to each outcome in a sample space. Whereas discrete random variables take on a discrete set of possible values, continuous random variables have a continuous set of values.

Computationally, to go from discrete to continuous we simply replace sums by integrals. It will help you to keep in mind that (informally) an integral is just a continuous sum.

Example 1. Since time is continuous, the amount of time Jon is early (or late) for class is a continuous random variable. Let's go over this example in some detail.

Suppose you measure how early Jon arrives to class each day (in units of minutes). That is, the outcome of one trial in our experiment is a time in minutes. We'll assume there are random fluctuations in the exact time he shows up. Since in principle Jon could arrive, say, 3.43 minutes early, or 2.7 minutes late (corresponding to the outcome -2.7), or at any other time, the sample space consists of all real numbers. So the random variable which gives the outcome itself has a **continuous range** of possible values.

It is too cumbersome to keep writing ‘the random variable’, so in future examples we might write: Let T = “time in minutes that Jon is early for class on any given day.”

3 Calculus Warmup

While we will assume you can compute the most familiar forms of derivatives and integrals by hand, we do not expect you to be calculus whizzes. For tricky expressions, we'll let the computer do most of the calculating. Conceptually, you should be comfortable with two views of a definite integral.

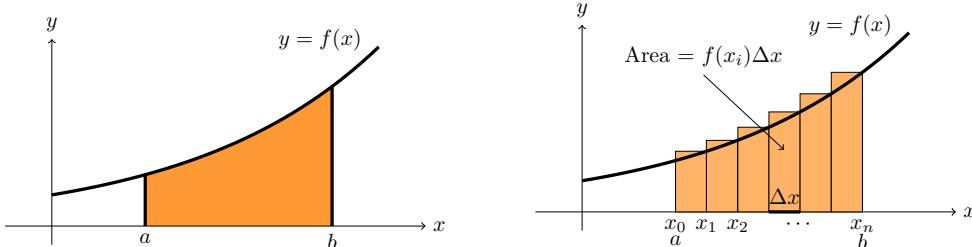
$$1. \int_a^b f(x) dx = \text{area under the curve } y = f(x).$$

$$2. \int_a^b f(x) dx = \text{'sum of } f(x) dx'.$$

The connection between the two is:

$$\text{area} \approx \text{sum of rectangle areas} = f(x_1)\Delta x + f(x_2)\Delta x + \dots + f(x_n)\Delta x = \sum_{i=1}^n f(x_i)\Delta x.$$

As the width Δx of the intervals gets smaller the approximation becomes better.



Area is approximately the sum of rectangles

Note: In calculus you learned to compute integrals by finding antiderivatives. This is important for calculations, but don't confuse this method for the reason we use integrals. Our interest in integrals comes primarily from its interpretation as a 'sum' and to a much lesser extent its interpretation as area.

4 Continuous Random Variables and Probability Density Functions

A continuous random variable takes a **range of values**, which may be finite or infinite in extent. Here are a few examples of ranges: $[0, 1]$, $[0, \infty)$, $(-\infty, \infty)$, $[a, b]$.

Definition: A random variable X is **continuous** if there is a function $f(x)$ such that for any $c \leq d$ we have

$$P(c \leq X \leq d) = \int_c^d f(x) dx. \quad (1)$$

The function $f(x)$ is called the **probability density function (pdf)**.

The pdf always satisfies the following properties:

1. $f(x) \geq 0$ (f is nonnegative).
2. $\int_{-\infty}^{\infty} f(x) dx = 1$ (This is equivalent to: $P(-\infty < X < \infty) = 1$).

The probability density function $f(x)$ of a continuous random variable is the analogue of the probability mass function $p(x)$ of a discrete random variable. Here are two important differences:

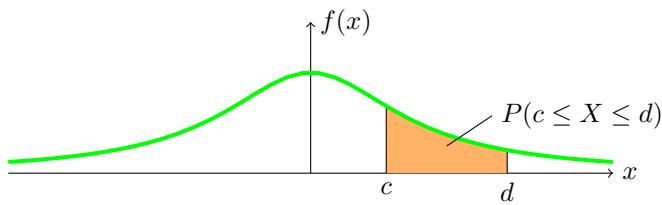
1. Unlike $p(x)$, the pdf $f(x)$ is *not* a probability. You have to integrate it to get probability. (See section 4.2 below.)
2. Since $f(x)$ is not a probability, there is no restriction that $f(x)$ be less than or equal to 1.

Note: In Property 2, we integrated over $(-\infty, \infty)$ since we did not know the range of values taken by X . Formally, this makes sense because we just define $f(x)$ to be 0 outside of the range of X . In practice, we would integrate between bounds given by the range of X .

4.1 Graphical View of Probability

If you graph the probability density function of a continuous random variable X then

$$P(c \leq X \leq d) = \text{area under the graph between } c \text{ and } d.$$



Think: What is the total area under the pdf $f(x)$?

4.2 The terms ‘probability mass’ and ‘probability density’

Why do we use the terms mass and density to describe the pmf and pdf? What is the difference between the two? The simple answer is that these terms are completely analogous to the mass and density you saw in physics and calculus. We'll review this first for the probability mass function and then discuss the probability density function.

Mass as a sum:

If masses m_1, m_2, m_3 , and m_4 are set in a row at positions x_1, x_2, x_3 , and x_4 , then the total mass is $m_1 + m_2 + m_3 + m_4$.

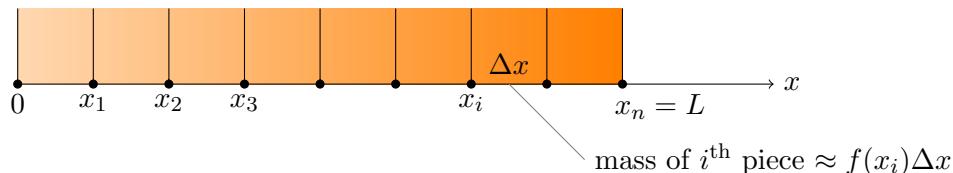


We can define a ‘mass function’ $p(x)$ with $p(x_j) = m_j$ for $j = 1, 2, 3, 4$, and $p(x) = 0$ otherwise. In this notation the total mass is $p(x_1) + p(x_2) + p(x_3) + p(x_4)$.

The **probability mass function** behaves in exactly the same way, except it has the dimension of probability instead of mass.

Mass as an integral of density:

Suppose you have a rod of length L meters with varying density $f(x)$ kg/m. (Note the units are mass/length.)



If the density varies continuously, we must find the total mass of the rod by integration:

$$\text{total mass} = \int_0^L f(x) dx.$$

This formula comes from dividing the rod into small pieces and 'summing' up the mass of each piece. That is:

$$\text{total mass} \approx \sum_{i=1}^n f(x_i) \Delta x$$

In the limit as Δx goes to zero the sum becomes the integral.

The **probability density function** behaves exactly the same way, except it has units of probability/(unit x) instead of kg/m. Indeed, equation (1) is exactly analogous to the above integral for total mass.

While we're on a physics kick, note that for both discrete and continuous random variables, the expected value is simply the **center of mass** or balance point.

Example 2. Suppose X has pdf $f(x) = 3$ on $[0, 1/3]$ (this means $f(x) = 0$ outside of $[0, 1/3]$). Graph the pdf and compute $P(.1 \leq X \leq .2)$ and $P(.1 \leq X \leq 1)$.

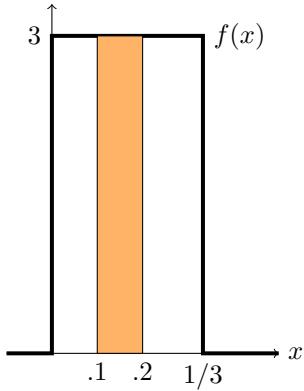
answer: $P(.1 \leq X \leq .2)$ is shown below at left. We can compute the integral:

$$P(.1 \leq X \leq .2) = \int_{.1}^{.2} f(x) dx = \int_{.1}^{.2} 3 dx = .3.$$

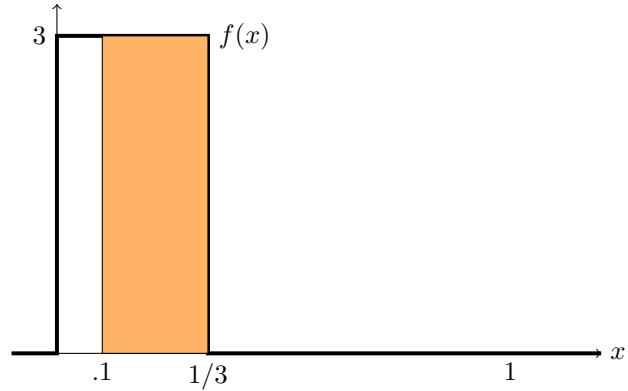
Or we can find the area geometrically:

$$\text{area of rectangle} = 3 \cdot .1 = .3.$$

$P(.1 \leq X \leq 1)$ is shown below at right. Since there is only area under $f(x)$ up to $1/3$, we have $P(.1 \leq X \leq 1) = 3 \cdot (1/3 - .1) = .7$.



$$P(.1 \leq X \leq .2)$$



$$P(.1 \leq X \leq 1)$$

Think: In the previous example $f(x)$ takes values greater than 1. Why does this not violate the rule that probabilities are always between 0 and 1?

Note on notation. We can define a random variable by giving its range and probability density function. For example we might say, let X be a random variable with range $[0,1]$

and pdf $f(x) = x/2$. Implicitly, this means that X has no probability density outside of the given range. If we wanted to be absolutely rigorous, we would say explicitly that $f(x) = 0$ outside of $[0,1]$, but in practice this won't be necessary.

Example 3. Let X be a random variable with range $[0,1]$ and pdf $f(x) = Cx^2$. What is the value of C ?

answer: Since the total probability must be 1, we have

$$\int_0^1 f(x) dx = 1 \quad \Leftrightarrow \quad \int_0^1 Cx^2 dx = 1.$$

By evaluating the integral, the equation at right becomes

$$C/3 = 1 \Rightarrow C = 3.$$

Note: We say the constant C above is needed to **normalize** the density so that the total probability is 1.

Example 4. Let X be the random variable in the Example 3. Find $P(X \leq 1/2)$.

answer: $P(X \leq 1/2) = \int_0^{1/2} 3x^2 dx = x^3 \Big|_0^{1/2} = \boxed{\frac{1}{8}}.$

Think: For this X (or any continuous random variable):

- What is $P(a \leq X \leq a)$?
- What is $P(X = 0)$?
- Does $P(X = a) = 0$ mean that X can never equal a ?

In words the above questions get at the fact that the probability that a random person's height is exactly 5'9" (to infinite precision, i.e. no rounding!) is 0. Yet it is still possible that someone's height is exactly 5'9". So the answers to the thinking questions are 0, 0, and No.

4.3 Cumulative Distribution Function

The **cumulative distribution function (cdf)** of a continuous random variable X is defined in exactly the same way as the cdf of a discrete random variable.

$$F(b) = P(X \leq b).$$

Note well that the definition is about probability. When using the cdf you should first think of it as a probability. Then when you go **to calculate** it you can use

$$F(b) = P(X \leq b) = \int_{-\infty}^b f(x) dx, \quad \text{where } f(x) \text{ is the pdf of } X.$$

Notes:

1. For discrete random variables, we defined the cumulative distribution function but did

not have much occasion to use it. The cdf plays a far more prominent role for continuous random variables.

2. As before, we started the integral at $-\infty$ because we did not know the precise range of X . Formally, this still makes sense since $f(x) = 0$ outside the range of X . In practice, we'll know the range and start the integral at the start of the range.
3. In practice we often say ‘ X has distribution $F(x)$ ’ rather than ‘ X has cumulative distribution function $F(x)$ ’.

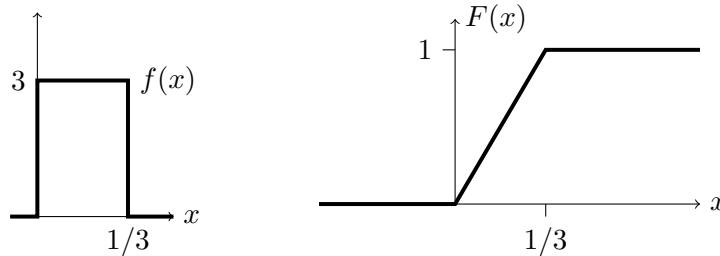
Example 5. Find the cumulative distribution function for the density in Example 2.

answer: For a in $[0,1/3]$ we have $F(a) = \int_0^a f(x) dx = \int_0^a 3 dx = 3a$.

Since $f(x)$ is 0 outside of $[0,1/3]$ we know $F(a) = P(X \leq a) = 0$ for $a < 0$ and $F(a) = 1$ for $a > 1/3$. Putting this all together we have

$$F(a) = \begin{cases} 0 & \text{if } a < 0 \\ 3a & \text{if } 0 \leq a \leq 1/3 \\ 1 & \text{if } 1/3 < a. \end{cases}$$

Here are the graphs of $f(x)$ and $F(x)$.



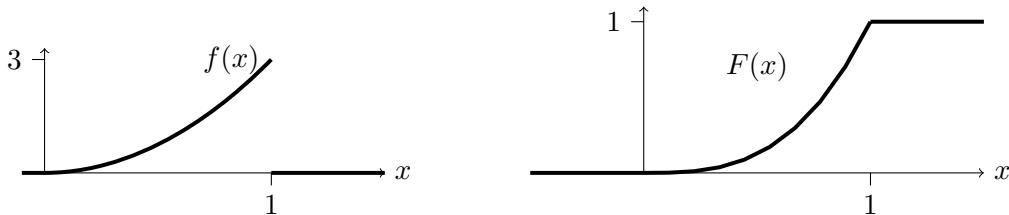
Note the different scales on the vertical axes. Remember that the vertical axis for the pdf represents probability density and that of the cdf represents probability.

Example 6. Find the cdf for the pdf in Example 3, $f(x) = 3x^2$ on $[0, 1]$. Suppose X is a random variable with this distribution. Find $P(X < 1/2)$.

answer: $f(x) = 3x^2$ on $[0,1] \Rightarrow F(a) = \int_0^a 3x^2 dx = a^3$ on $[0,1]$. Therefore,

$$F(a) = \begin{cases} 0 & \text{if } a < 0 \\ a^3 & \text{if } 0 \leq a \leq 1 \\ 1 & \text{if } 1 < a \end{cases}$$

Thus, $P(X < 1/2) = F(1/2) = 1/8$. Here are the graphs of $f(x)$ and $F(x)$:



4.4 Properties of cumulative distribution functions

Here is a summary of the most important properties of cumulative distribution functions (cdf)

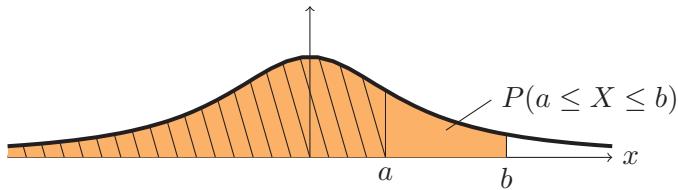
1. (Definition) $F(x) = P(X \leq x)$
2. $0 \leq F(x) \leq 1$
3. $F(x)$ is non-decreasing, i.e. if $a \leq b$ then $F(a) \leq F(b)$.
4. $\lim_{x \rightarrow \infty} F(x) = 1$ and $\lim_{x \rightarrow -\infty} F(x) = 0$
5. $P(a \leq X \leq b) = F(b) - F(a)$
6. $F'(x) = f(x)$.

Properties 2, 3, 4 are identical to those for discrete distributions. The graphs in the previous examples illustrate them.

Property 5 can be seen algebraically:

$$\begin{aligned} \int_{-\infty}^b f(x) dx &= \int_{-\infty}^a f(x) dx + \int_a^b f(x) dx \\ \Leftrightarrow \int_a^b f(x) dx &= \int_{-\infty}^b f(x) dx - \int_{-\infty}^a f(x) dx \\ \Leftrightarrow P(a \leq X \leq b) &= F(b) - F(a). \end{aligned}$$

Property 5 can also be seen geometrically. The orange region below represents $F(b)$ and the striped region represents $F(a)$. Their difference is $P(a \leq X \leq b)$.



Property 6 is the fundamental theorem of calculus.

4.5 Probability density as a dartboard

We find it helpful to think of sampling values from a continuous random variable as throwing darts at a funny dartboard. Consider the region underneath the graph of a pdf as a dartboard. Divide the board into small equal size squares and suppose that when you throw a dart you are equally likely to land in any of the squares. The probability the dart lands in a given region is the fraction of the total area under the curve taken up by the region. Since the total area equals 1, this fraction is just the area of the region. If X represents the x -coordinate of the dart, then the probability that the dart lands with x -coordinate between a and b is just

$$P(a \leq X \leq b) = \text{area under } f(x) \text{ between } a \text{ and } b = \int_a^b f(x) dx.$$

Gallery of Continuous Random Variables

Class 5, 18.05, Spring 2014

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to give examples of what uniform, exponential and normal distributions are used to model.
2. Be able to give the range and pdf's of uniform, exponential and normal distributions.

2 Introduction

Here we introduce a few fundamental continuous distributions. These will play important roles in the statistics part of the class. For each distribution, we give the range, the pdf, the cdf, and a short description of situations that it models. These distributions all depend on parameters, which we specify.

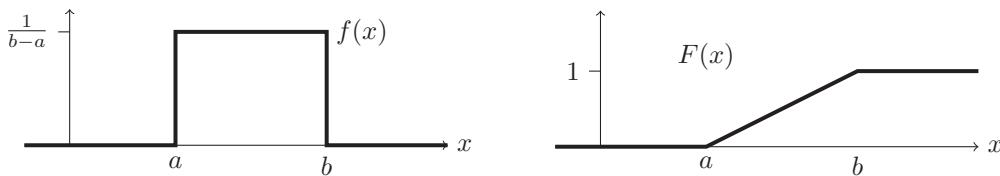
As you look through each distribution do not try to memorize all the details; you can always look those up. Rather, focus on the shape of each distribution and what it models.

Although it comes towards the end, we call your attention to the normal distribution. It is easily the most important distribution defined here.

3 Uniform distribution

1. Parameters: a, b .
2. Range: $[a, b]$.
3. Notation: $\text{uniform}(a, b)$ or $U(a, b)$.
4. Density: $f(x) = \frac{1}{b-a}$ for $a \leq x \leq b$.
5. Distribution: $F(x) = (x-a)/(b-a)$ for $a \leq x \leq b$.
6. Models: All outcomes in the range have equal probability (more precisely all outcomes have the same probability density).

Graphs:



pdf and cdf for $\text{uniform}(a, b)$ distribution.

Examples. 1. Suppose we have a tape measure with markings at each millimeter. If we measure (to the nearest marking) the length of items that are roughly a meter long, the rounding error will uniformly distributed between -0.5 and 0.5 millimeters.

2. Many boardgames use spinning arrows (spinners) to introduce randomness. When spun, the arrow stops at an angle that is uniformly distributed between 0 and 2π radians.
3. In most pseudo-random number generators, the basic generator simulates a uniform distribution and all other distributions are constructed by transforming the basic generator.

4 Exponential distribution

1. Parameter: λ .
2. Range: $[0, \infty)$.
3. Notation: $\text{exponential}(\lambda)$ or $\exp(\lambda)$.
4. Density: $f(x) = \lambda e^{-\lambda x}$ for $0 \leq x$.
5. Distribution: (easy integral)

$$F(x) = 1 - e^{-\lambda x} \text{ for } x \geq 0$$

6. *Right tail distribution:* $P(X > x) = 1 - F(x) = e^{-\lambda x}$.
7. Models: The waiting time for a continuous process to change state.

Examples. 1. If I step out to 77 Mass Ave after class and wait for the next taxi, my waiting time in minutes is exponentially distributed. We will see that in this case λ is given by one over the average number of taxis that pass per minute (on weekday afternoons).

2. The exponential distribution models the waiting time until an unstable isotope undergoes nuclear decay. In this case, the value of λ is related to the half-life of the isotope.

Memorylessness: There are other distributions that also model waiting times, but the exponential distribution has the additional property that it is memoryless. Here's what this means in the context of Example 1. Suppose that the probability that a taxi arrives within the first five minutes is p . If I wait five minutes and in fact no taxi arrives, then the probability that a taxi arrives within the next five minutes is still p .

By contrast, suppose I were to instead go to Kendall Square subway station and wait for the next inbound train. Since the trains are coordinated to follow a schedule (e.g., roughly 12 minutes between trains), if I wait five minutes without seeing a train then there is a far greater probability that a train will arrive in the next five minutes. In particular, waiting time for the subway is not memoryless, and a better model would be the uniform distribution on the range $[0, 12]$.

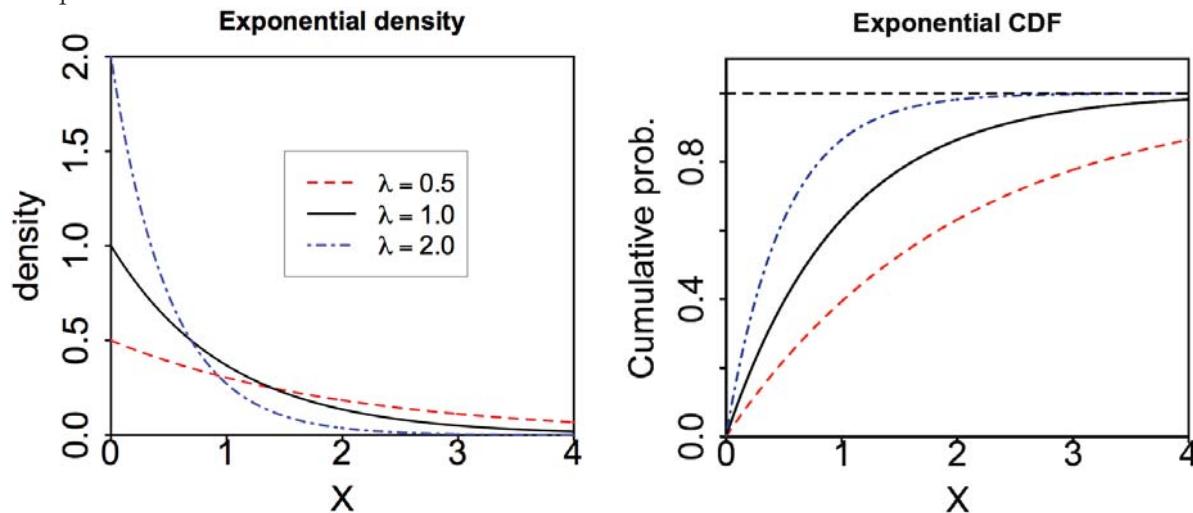
The memorylessness of the exponential distribution is analogous to the memorylessness of the (discrete) geometric distribution, where having flipped 5 tails in a row gives no information about the next 5 flips. Indeed, the exponential distribution is the precisely the

continuous counterpart of the geometric distribution, which models the waiting time for a discrete process to change state. More formally, memoryless means that the probability of waiting t more minutes is unaffected by having already waited s minutes without incident. In symbols, $P(X > s + t | X > s) = P(X > t)$.

Proof of memorylessness: Since $(X > s + t) \cap (X > s) = (X > s + t)$ we have

$$P(X > s + t | X > s) = \frac{P(X > s + t)}{P(X > s)} = \frac{e^{-\lambda(s+t)}}{e^{-\lambda s}} = e^{-\lambda t} = P(X > t). \text{ QED}$$

Graphs:



5 Normal distribution

In 1809, Carl Friedrich Gauss published a monograph introducing several notions that have become fundamental to statistics: the normal distribution, maximum likelihood estimation, and the method of least squares (we will cover all three in this course). For this reason, the normal distribution is also called the *Gaussian* distribution, and it is the most important continuous distribution.

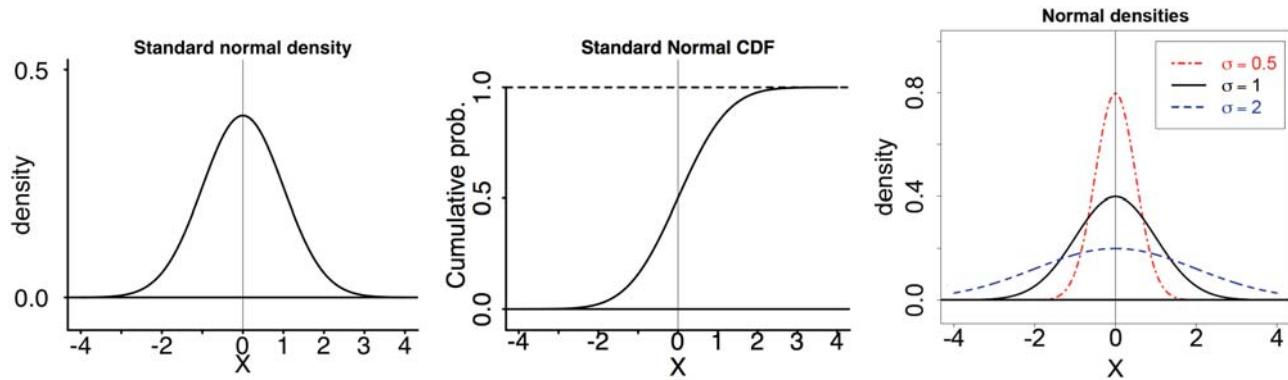
1. Parameters: μ, σ .
2. Range: $(-\infty, \infty)$.
3. Notation: $\text{normal}(\mu, \sigma^2)$ or $N(\mu, \sigma^2)$.
4. Density: $f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$.
5. Distribution: $F(x)$ has no formula, so use tables or software such as `pnorm` in R to compute $F(x)$.
6. Models: Measurement error, intelligence/ability, height, averages of lots of data.

The **standard normal distribution** $N(0, 1)$ has mean 0 and variance 1. We reserve Z for a standard normal random variable, $\phi(z) = \frac{1}{\sqrt{2\pi}}e^{-z^2/2}$ for the standard normal density, and $\Phi(z)$ for the standard normal distribution.

Note: we will define mean and variance for continuous random variables next time. They have the same interpretations as in the discrete case. As you might guess, the normal distribution $N(\mu, \sigma^2)$ has mean μ , variance σ^2 , and standard deviation σ .

Here are some graphs of normal distributions. Note they are shaped like a *bell curve*. Note also that as σ increases they become more spread out.

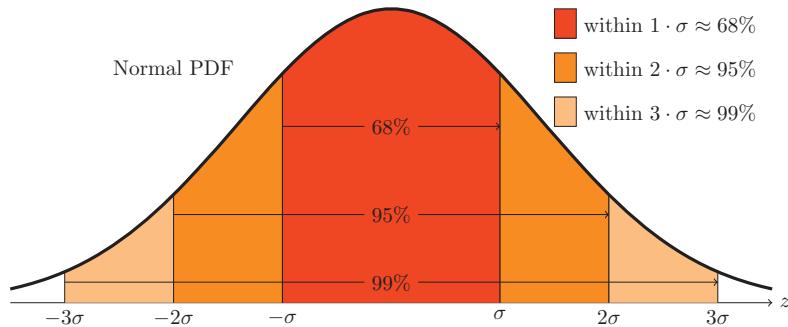
Graphs: (the *bell curve*):



5.1 Normal probabilities

To make approximations it is useful to remember the following rule of thumb for three approximate probabilities

$$P(-1 \leq Z \leq 1) \approx .68, \quad P(-2 \leq Z \leq 2) \approx .95, \quad P(-3 \leq Z \leq 3) \approx .99$$

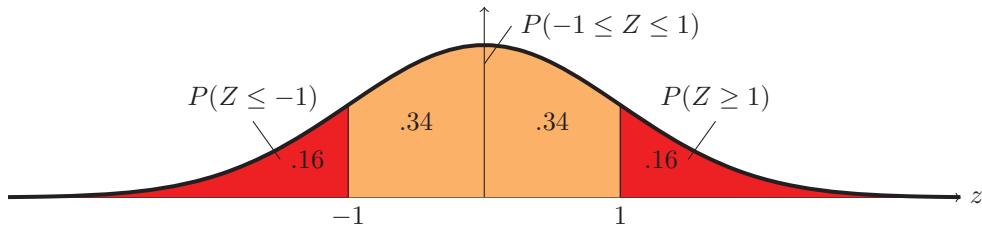


Symmetry calculations

We can use the symmetry of the standard normal distribution about $x = 0$ to make some calculations.

Example 1. The rule of thumb says $P(-1 \leq Z \leq 1) \approx .68$. Use this to estimate $\Phi(1)$.

answer: $\Phi(1) = P(Z \leq 1)$. In the figure, the two tails (in red) have combined area $1 - .68 = .32$. By symmetry the left tail has area $.16$ (half of $.32$), so $P(Z \leq 1) \approx .68 + .16 = .84$.



5.2 Using R to compute $\Phi(z)$.

```
# Use the R function pnorm(x,μ,σ) to compute F(x) for N(μ,σ²)
pnorm(1,0,1)
[1] 0.8413447

pnorm(0,0,1)
[1] 0.5
pnorm(1,0,2)
[1] 0.6914625

pnorm(1,0,1) - pnorm(-1,0,1)
[1] 0.6826895
pnorm(5,0,5) - pnorm(-5,0,5)
[1] 0.6826895

# Of course z can be a vector of values
pnorm(c(-3,-2,-1,0,1,2,3),0,1)
[1] 0.001349898 0.022750132 0.158655254 0.500000000 0.841344746 0.977249868 0.998650102
```

Note: The R function $\text{pnorm}(x, \mu, \sigma)$ uses σ whereas our notation for the normal distribution $N(\mu, \sigma^2)$ uses σ^2 .

Here's a table of values with fewer decimal points of accuracy

$z:$	-2	-1	0	.3	.5	1	2	3
$\Phi(z):$	0.0228	0.1587	0.5000	0.6179	0.6915	0.8413	0.9772	0.9987

Example 2. Use R to compute $P(-1.5 \leq Z \leq 2)$.

answer: This is $\Phi(2) - \Phi(-1.5) = \text{pnorm}(2,0,1) - \text{pnorm}(-1.5,0,1) = 0.91044$

6 Pareto and other distributions

In 18.05, we only have time to work with a few of the many wonderful distributions that are used in probability and statistics. We hope that after this course you will feel comfortable learning about new distributions and their properties when you need them. Wikipedia is often a great starting point.

The Pareto distribution is one common, beautiful distribution that we will not have time to cover in depth.

1. Parameters: $m > 0$ and $\alpha > 0$.

2. Range: $[m, \infty)$.

3. Notation: Pareto(m, α).

4. Density: $f(x) = \frac{\alpha m^\alpha}{x^{\alpha+1}}$.

5. Distribution: (easy integral)

$$F(x) = 1 - \frac{m^\alpha}{x^\alpha}, \text{ for } x \geq m$$

6. Tail distribution: $P(X > x) = m^\alpha/x^\alpha$, for $x \geq m$.

7. Models: The Pareto distribution models a **power law**, where the probability that an event occurs varies as a power of some attribute of the event. Many phenomena follow a power law, such as the size of meteors, income levels across a population, and population levels across cities. See Wikipedia for loads of examples:

http://en.wikipedia.org/wiki/Pareto_distribution#Applications

Manipulating Continuous Random Variables

Class 5, 18.05, Spring 2014

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to find the pdf and cdf of a random variable defined in terms of a random variable with known pdf and cdf.

2 Transformations of Random Variables

If $Y = aX + b$ then the properties of expectation and variance tell us that $E(Y) = aE(X) + b$ and $\text{Var}(Y) = a^2\text{Var}(X)$. But what is the distribution function of Y ? If Y is continuous, what is its pdf?

Often, when looking at transforms of discrete random variables we work with tables. For continuous random variables transforming the pdf is just change of variables (' u -substitution') from calculus. Transforming the cdf makes direct use of the definition of the cdf.

Let's remind ourselves of the basics:

1. The cdf of X is $F_X(x) = P(X \leq x)$.
2. The pdf of X is related to F_X by $f_X(x) = F'_X(x)$.

Example 1. Let $X \sim U(0, 2)$, so $f_X(x) = 1/2$ and $F_X(x) = x/2$ on $[0, 2]$. What is the range, pdf and cdf of $Y = X^2$?

answer: The range is easy: $[0, 4]$.

To find the cdf we work systematically from the definition.

$$F_Y(y) = P(Y \leq y) = P(X^2 \leq y) = P(X \leq \sqrt{y}) = F_X(\sqrt{y}) = \sqrt{y}/2.$$

To find the pdf we can just differentiate the cdf

$$f_Y(y) = \frac{d}{dy} F_Y(y) = \boxed{\frac{1}{4\sqrt{y}}}.$$

An alternative way to find the pdf directly is by change of variables. The trick here is to remember that it is $f_X(x)dx$ which gives probability ($f_X(x)$ by itself is probability density). Here is how the calculation goes in this example.

$$\begin{aligned} y = x^2 &\Rightarrow dy = 2x dx \Rightarrow dx = \frac{dy}{2\sqrt{y}} \\ f_X(x) dx &= \frac{dx}{2} = \frac{dy}{4\sqrt{y}} = f_Y(y) dy \end{aligned}$$

Therefore
$$\boxed{f_Y(y) = \frac{dy}{4\sqrt{y}}}$$

Example 2. Let $X \sim \exp(\lambda)$, so $f_X(x) = \lambda e^{-\lambda x}$ on $[0, \infty]$. What is the density of $Y = X^2$?

answer: Let's do this using the change of variables.

$$\begin{aligned} y = x^2 &\Rightarrow dy = 2x dx \Rightarrow dx = \frac{dy}{2\sqrt{y}} \\ f_X(x) dx &= \lambda e^{-\lambda x} dx = \lambda e^{-\lambda\sqrt{y}} \frac{dy}{2\sqrt{y}} = f_Y(y) dy \end{aligned}$$

Therefore
$$f_Y(y) = \frac{\lambda}{2\sqrt{y}} e^{-\lambda\sqrt{y}}.$$

Example 3. Assume $X \sim N(5, 3^2)$. Show that $Z = \frac{X - 5}{3}$ is standard normal, i.e., $Z \sim N(0, 1)$.

answer: Again using the change of variables and the formula for $f_X(x)$ we have

$$\begin{aligned} z &= \frac{x - 5}{3} \Rightarrow dz = \frac{dx}{3} \Rightarrow dx = 3 dz \\ f_X(x) dx &= \frac{1}{3\sqrt{2\pi}} e^{-(x-5)^2/(2\cdot3^2)} dx = \frac{1}{3\sqrt{2\pi}} e^{-z^2/2} 3 dz = \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz = f_Z(z) dz \end{aligned}$$

Therefore $f_Z(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$. Since this is exactly the density for $N(0, 1)$ we have shown that Z is standard normal.

This example shows an important general property of normal random variables which we give in the next example.

Example 4. Assume $X \sim N(\mu, \sigma^2)$. Show that $Z = \frac{X - \mu}{\sigma}$ is standard normal, i.e., $Z \sim N(0, 1)$.

answer: This is exactly the same computation as the previous example with μ replacing 5 and σ replacing 3. We show the computation without comment.

$$\begin{aligned} z &= \frac{x - \mu}{\sigma} \Rightarrow dz = \frac{dx}{\sigma} \Rightarrow dx = \sigma dz \\ f_X(x) dx &= \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/(2\cdot\sigma^2)} dx = \frac{1}{\sigma\sqrt{2\pi}} e^{-z^2/2} \sigma dz = \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz = f_Z(z) dz \end{aligned}$$

Therefore $f_Z(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$. This shows Z is standard normal.

Expectation, Variance and Standard Deviation for Continuous Random Variables

Class 6, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to compute and interpret expectation, variance, and standard deviation for continuous random variables.
2. Be able to compute and interpret quantiles for discrete and continuous random variables.

2 Introduction

So far we have looked at expected value, standard deviation, and variance for discrete random variables. These summary statistics have the same meaning for continuous random variables:

- The expected value $\mu = E(X)$ is a measure of location or central tendency.
- The standard deviation σ is a measure of the spread or scale.
- The variance $\sigma^2 = \text{Var}(X)$ is the square of the standard deviation.

To move from discrete to continuous, we will simply replace the sums in the formulas by integrals. We will do this carefully and go through many examples in the following sections. In the last section, we will introduce another type of [summary statistic](#), [quantiles](#). You may already be familiar with the .5 quantile of a distribution, otherwise known as the [median](#) or 50th percentile.

3 Expected value of a continuous random variable

Definition: Let X be a continuous random variable with range $[a, b]$ and probability density function $f(x)$. The [expected value](#) of X is defined by

$$E(X) = \int_a^b x f(x) dx.$$

Let's see how this compares with the formula for a discrete random variable:

$$E(X) = \sum_{i=1}^n x_i p(x_i).$$

The discrete formula says to take a weighted sum of the values x_i of X , where the weights are the probabilities $p(x_i)$. Recall that $f(x)$ is a probability [density](#). Its units are prob/(unit of X).

So $f(x) dx$ represents the probability that X is in an infinitesimal range of width dx around x . Thus we can interpret the formula for $E(X)$ as a weighted integral of the values x of X , where the weights are the probabilities $f(x) dx$.

As before, the expected value is also called the **mean** or **average**.

3.1 Examples

Let's go through several example computations. Where the solution requires an integration technique, we push the computation of the integral to the appendix.

Example 1. Let $X \sim \text{uniform}(0, 1)$. Find $E(X)$.

answer: X has range $[0, 1]$ and density $f(x) = 1$. Therefore,

$$E(X) = \int_0^1 x dx = \frac{x^2}{2} \Big|_0^1 = \boxed{\frac{1}{2}}.$$

Not surprisingly the mean is at the midpoint of the range.

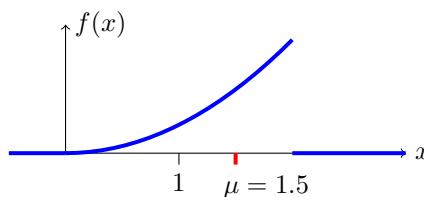
Example 2. Let X have range $[0, 2]$ and density $\frac{3}{8}x^2$. Find $E(X)$.

answer:

$$E(X) = \int_0^2 xf(x) dx = \int_0^2 \frac{3}{8}x^3 dx = \frac{3x^4}{32} \Big|_0^2 = \boxed{\frac{3}{2}}.$$

Does it make sense that this X has mean is in the right half of its range?

answer: Yes. Since the probability density increases as x increases over the range, the average value of x should be in the right half of the range.

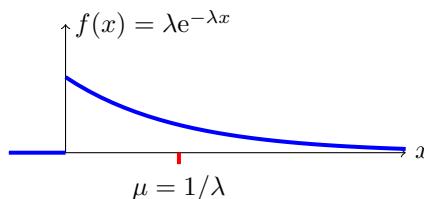


μ is “pulled” to the right of the midpoint 1 because there is more mass to the right.

Example 3. Let $X \sim \exp(\lambda)$. Find $E(X)$.

answer: The range of X is $[0, \infty)$ and its pdf is $f(x) = \lambda e^{-\lambda x}$. So (details in appendix)

$$E(X) = \int_0^\infty \lambda e^{-\lambda x} dx = -\lambda e^{-\lambda x} - \frac{e^{-\lambda x}}{\lambda} \Big|_0^\infty = \boxed{\frac{1}{\lambda}}.$$

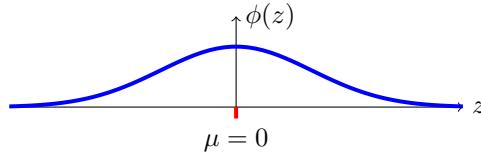


Mean of an exponential random variable

Example 4. Let $Z \sim N(0, 1)$. Find $E(Z)$.

answer: The range of Z is $(-\infty, \infty)$ and its pdf is $\phi(z) = \frac{1}{\sqrt{2\pi}}e^{-z^2/2}$. So (details in appendix)

$$E(Z) = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}}ze^{-z^2/2} dz = -\frac{1}{\sqrt{2\pi}}e^{-z^2/2} \Big|_{-\infty}^{\infty} = [0].$$



The standard normal distribution is symmetric and has mean 0.

3.2 Properties of $E(X)$

The properties of $E(X)$ for continuous random variables are the same as for discrete ones:

1. If X and Y are random variables on a sample space Ω then

$$E(X + Y) = E(X) + E(Y). \quad (\text{linearity I})$$

2. If a and b are constants then

$$E(aX + b) = aE(X) + b. \quad (\text{linearity II})$$

Example 5. In this example we verify that for $X \sim N(\mu, \sigma)$ we have $E(X) = \mu$.

answer: Example (4) showed that for standard normal Z , $E(Z) = 0$. We could mimic the calculation there to show that $E(X) = \mu$. Instead we will use the linearity properties of $E(X)$. In the class 5 notes on manipulating random variables we showed that if $X \sim N(\mu, \sigma^2)$ is a normal random variable we can **standardize** it:

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1).$$

Inverting this formula we have $X = \sigma Z + \mu$. The linearity of expected value now gives

$$E(X) = E(\sigma Z + \mu) = \sigma E(Z) + \mu = \mu$$

3.3 Expectation of Functions of X

This works exactly the same as the discrete case. if $h(x)$ is a function then $Y = h(X)$ is a random variable and

$$E(Y) = E(h(X)) = \int_{-\infty}^{\infty} h(x)f_X(x) dx.$$

Example 6. Let $X \sim \exp(\lambda)$. Find $E(X^2)$.

answer: Using integration by parts we have

$$E(X^2) = \int_0^{\infty} x^2 \lambda e^{-\lambda x} dx = \left[-x^2 e^{-\lambda x} - \frac{2x}{\lambda} e^{-\lambda x} - \frac{2}{\lambda^2} e^{-\lambda x} \right]_0^{\infty} = \boxed{\frac{2}{\lambda^2}}.$$

4 Variance

Now that we've defined expectation for continuous random variables, the definition of variance is identical to that of discrete random variables.

Definition: Let X be a continuous random variable with mean μ . The **variance** of X is

$$\text{Var}(X) = E((X - \mu)^2).$$

4.1 Properties of Variance

These are exactly the same as in the discrete case.

1. If X and Y are **independent** then $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$.
2. For constants a and b , $\text{Var}(aX + b) = a^2\text{Var}(X)$.
3. **Theorem:** $\text{Var}(X) = E(X^2) - E(X)^2 = E(X^2) - \mu^2$.

For Property 1, note carefully the requirement that X and Y are **independent**.

Property 3 gives a formula for $\text{Var}(X)$ that is often easier to use in hand calculations. The proofs of properties 2 and 3 are essentially identical to those in the discrete case. We will not give them here.

Example 7. Let $X \sim \text{uniform}(0, 1)$. Find $\text{Var}(X)$ and σ_X .

answer: In Example 1 we found $\mu = 1/2$. Next we compute

$$\text{Var}(X) = E((X - \mu)^2) = \int_0^1 (x - 1/2)^2 dx = \boxed{\frac{1}{12}}.$$

Example 8. Let $X \sim \exp(\lambda)$. Find $\text{Var}(X)$ and σ_X .

answer: In Examples 3 and 6 we computed

$$E(X) = \int_0^\infty x \lambda e^{-\lambda x} dx = \frac{1}{\lambda} \quad \text{and} \quad E(X^2) = \int_0^\infty x^2 \lambda e^{-\lambda x} dx = \frac{2}{\lambda^2}.$$

So by Property 3,

$$\text{Var}(X) = E(X^2) - E(X)^2 = \frac{2}{\lambda^2} - \frac{1}{\lambda^2} = \frac{1}{\lambda^2} \quad \text{and} \quad \sigma_X = \frac{1}{\lambda}.$$

We could have skipped Property 3 and computed this directly from $\text{Var}(X) = \int_0^\infty (x - 1/\lambda)^2 \lambda e^{-\lambda x} dx$.

Example 9. Let $Z \sim N(0, 1)$. Show $\text{Var}(Z) = 1$.

Note: The notation for normal variables is $X \sim N(\mu, \sigma^2)$. This is certainly suggestive, but as mathematicians we need to prove that $E(X) = \mu$ and $\text{Var}(X) = \sigma^2$. Above we showed $E(X) = \mu$. This example shows that $\text{Var}(Z) = 1$, just as the notation suggests. In the next example we'll show $\text{Var}(X) = \sigma^2$.

answer: Since $E(Z) = 0$, we have

$$\text{Var}(Z) = E(Z^2) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} z^2 e^{-z^2/2} dz.$$

(using integration by parts with $u = z$, $v' = ze^{-z^2/2} \Rightarrow u' = 1$, $v = -e^{-z^2/2}$)

$$= \frac{1}{\sqrt{2\pi}} \left(-ze^{-z^2/2} \Big|_{-\infty}^{\infty} \right) + \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-z^2/2} dz.$$

The first term equals 0 because the exponential goes to zero much faster than z grows at both $\pm\infty$. The second term equals 1 because it is exactly the total probability integral of the pdf $\varphi(z)$ for $N(0, 1)$. So $\text{Var}(X) = 1$.

Example 10. Let $X \sim N(\mu, \sigma^2)$. Show $\text{Var}(X) = \sigma^2$.

answer: This is an exercise in change of variables. Letting $z = (x - \mu)/\sigma$, we have

$$\begin{aligned} \text{Var}(X) &= E((X - \mu)^2) = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} (x - \mu)^2 e^{-(x-\mu)^2/2\sigma^2} dx \\ &= \frac{\sigma^2}{\sqrt{2\pi}} \int_{-\infty}^{\infty} z^2 e^{-z^2/2} dz = \sigma^2. \end{aligned}$$

The integral in the last line is the same one we computed for $\text{Var}(Z)$.

5 Quantiles

Definition: The median of X is the value x for which $P(X \leq x) = 0.5$, i.e. the value of x such that $P(X \leq x) = P(X \geq x)$. In other words, X has equal probability of being above or below the median, and each probability is therefore $1/2$. In terms of the cdf $F(x) = P(X \leq x)$, we can equivalently define the median as the value x satisfying $F(x) = 0.5$.

Think: What is the median of Z ?

answer: By symmetry, the median is 0.

Example 11. Find the median of $X \sim \exp(\lambda)$.

answer: The cdf of X is $F(x) = 1 - e^{-\lambda x}$. So the median is the value of x for which $F(x) = 1 - e^{-\lambda x} = 0.5$. Solving for x we find: $x = (\ln 2)/\lambda$.

Think: In this case the median does not equal the mean of $\mu = 1/\lambda$. Based on the graph of the pdf of X can you argue why the median is to the left of the mean.

Definition: The p^{th} quantile of X is the value q_p such that $P(X \leq q_p) = p$.

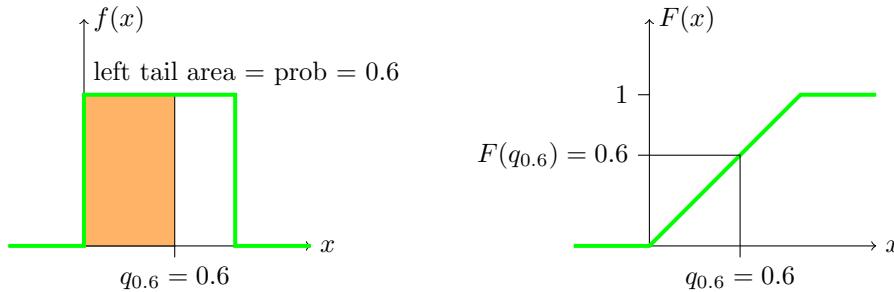
Notes. 1. In this notation the median is $q_{0.5}$.

2. We will usually write this in terms of the cdf: $F(q_p) = p$.

With respect to the pdf $f(x)$, the quantile q_p is the value such that there is an area of p to the left of q_p and an area of $1 - p$ to the right of q_p . In the examples below, note how we can represent the quantile graphically using either area of the pdf or height of the cdf.

Example 12. Find the 0.6 quantile for $X \sim U(0, 1)$.

answer: The cdf for X is $F(x) = x$ on the range $[0,1]$. So $q_{0.6} = 0.6$.

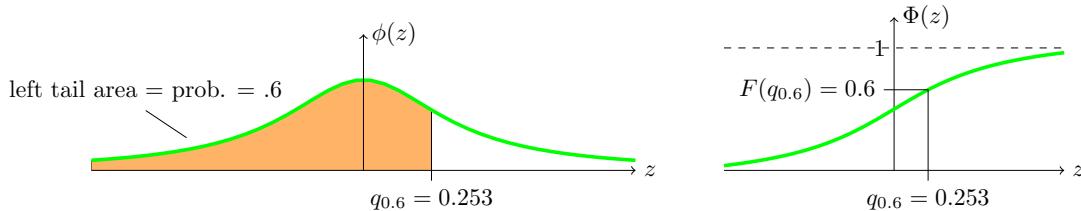


$$q_{0.6}: \text{left tail area} = 0.6 \Leftrightarrow F(q_{0.6}) = 0.6$$

Example 13. Find the 0.6 quantile of the standard normal distribution.

answer: We don't have a formula for the cdf, so we use the R 'quantile function' `qnorm`.

$$q_{0.6} = \text{qnorm}(0.6, 0, 1) = 0.25335$$



$$q_{0.6}: \text{left tail area} = 0.6 \Leftrightarrow F(q_{0.6}) = 0.6$$

Quantiles give a useful measure of **location** for a random variable. We will use them more in coming lectures.

5.1 Percentiles, deciles, quartiles

For convenience, quantiles are often described in terms of percentiles, deciles or quartiles. The 60th **percentile** is the same as the 0.6 quantile. For example you are in the 60th percentile for height if you are taller than 60 percent of the population, i.e. the **probability** that you are taller than a randomly chosen person is 60 percent.

Likewise, **deciles** represent steps of 1/10. The third decile is the 0.3 quantile. **Quartiles** are in steps of 1/4. The third quartile is the 0.75 quantile and the 75th percentile.

6 Appendix: Integral Computation Details

Example 3: Let $X \sim \exp(\lambda)$. Find $E(X)$.

The range of X is $[0, \infty)$ and its pdf is $f(x) = \lambda e^{-\lambda x}$. Therefore

$$E(X) = \int_0^\infty x f(x) dx = \int_0^\infty \lambda x e^{-\lambda x} dx$$

(using integration by parts with $u = x, v' = \lambda e^{-\lambda x} \Rightarrow u' = 1, v = -e^{-\lambda x}$)

$$\begin{aligned} &= -xe^{-\lambda x} \Big|_0^\infty + \int_0^\infty e^{-\lambda x} dx \\ &= 0 - \frac{e^{-\lambda x}}{\lambda} \Big|_0^\infty = \frac{1}{\lambda}. \end{aligned}$$

We used the fact that $xe^{-\lambda x}$ and $e^{-\lambda x}$ go to 0 as $x \rightarrow \infty$.

Example 4: Let $Z \sim N(0, 1)$. Find $E(Z)$.

The range of Z is $(-\infty, \infty)$ and its pdf is $\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$. By symmetry the mean must be 0. The only mathematically tricky part is to show that the integral converges, i.e. that the mean exists at all (some random variable do not have means, but we will not encounter this very often.) For completeness we include the argument, though this is not something we will ask you to do. We first compute the integral from 0 to ∞ :

$$\int_0^\infty z\phi(z) dz = \frac{1}{\sqrt{2\pi}} \int_0^\infty ze^{-z^2/2} dz.$$

The u -substitution $u = z^2/2$ gives $du = z dz$. So the integral becomes

$$\frac{1}{\sqrt{2\pi}} \int_0^\infty ze^{-z^2/2} dz = \frac{1}{\sqrt{2\pi}} \int_0^\infty e^{-u} du = -e^{-u} \Big|_0^\infty = 1$$

Similarly, $\int_{-\infty}^0 z\phi(z) dz = -1$. Adding the two pieces together gives $E(Z) = 0$.

Example 6: Let $X \sim \exp(\lambda)$. Find $E(X^2)$.

$$E(X^2) = \int_0^\infty x^2 f(x) dx = \int_0^\infty \lambda x^2 e^{-\lambda x} dx$$

(using integration by parts with $u = x^2, v' = \lambda e^{-\lambda x} \Rightarrow u' = 2x, v = -e^{-\lambda x}$)

$$= -x^2 e^{-\lambda x} \Big|_0^\infty + \int_0^\infty 2x e^{-\lambda x} dx$$

(the first term is 0, for the second term use integration by parts: $u = 2x, v' = e^{-\lambda x} \Rightarrow u' = 2, v = -\frac{e^{-\lambda x}}{\lambda}$)

$$\begin{aligned} &= -2x \frac{e^{-\lambda x}}{\lambda} \Big|_0^\infty + \int_0^\infty \frac{e^{-\lambda x}}{\lambda} dx \\ &= 0 - 2 \frac{e^{-\lambda x}}{\lambda^2} \Big|_0^\infty = \frac{2}{\lambda^2}. \end{aligned}$$

Central Limit Theorem and the Law of Large Numbers

Class 6, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Understand the statement of the law of large numbers.
2. Understand the statement of the central limit theorem.
3. Be able to use the central limit theorem to approximate probabilities of averages and sums of independent identically-distributed random variables.

2 Introduction

We all understand intuitively that the average of many measurements of the same unknown quantity tends to give a better estimate than a single measurement. Intuitively, this is because the random error of each measurement cancels out in the average. In these notes we will make this intuition precise in two ways: the law of large numbers (LoLN) and the central limit theorem (CLT).

Briefly, both the law of large numbers and central limit theorem are about many independent samples from same distribution. The LoLN tells us two things:

1. The average of many independent samples is (with high probability) close to the mean of the underlying distribution.
2. This density histogram of many independent samples is (with high probability) close to the graph of the density of the underlying distribution.

To be absolutely correct mathematically we need to make these statements more precise, but as stated they are a good way to think about the law of large numbers.

The central limit theorem says that the sum or average of many independent copies of a random variable is approximately a normal random variable. The CLT goes on to give precise values for the mean and standard deviation of the normal variable.

These are both remarkable facts. Perhaps just as remarkable is the fact that often in practice n does not have to be all that large. Values of $n > 30$ often suffice.

2.1 There is more to experimentation than mathematics

The mathematics of the LoLN says that the average of a lot of independent samples from a random variable will almost certainly approach the mean of the variable. The mathematics **cannot** tell us if the tool or experiment is producing data worth averaging. For example, if the measuring device is defective or poorly calibrated then the average of many measurements will be a highly accurate estimate of the wrong thing! This is an example of

systematic error or sampling bias, as opposed to the random error controlled by the law of large numbers.

3 The law of large numbers

Suppose X_1, X_2, \dots, X_n are independent random variables with the same underlying distribution. In this case, we say that the X_i are **independent and identically-distributed**, or **i.i.d.**. In particular, the X_i all have the same mean μ and standard deviation σ .

Let \bar{X}_n be the average of X_1, \dots, X_n :

$$\bar{X}_n = \frac{X_1 + X_2 + \dots + X_n}{n} = \frac{1}{n} \sum_{i=1}^n X_i.$$

Note that \bar{X}_n is itself a random variable. The law of large numbers and central limit theorem tell us about the value and distribution of \bar{X}_n , respectively.

LoLN: As n grows, the probability that \bar{X}_n is close to μ goes to 1.

CLT: As n grows, the distribution of \bar{X}_n converges to the normal distribution $N(\mu, \sigma^2/n)$.

Before giving a more formal statement of the LoLN, let's unpack its meaning through a concrete example (we'll return to the CLT later on).

Example 1. Averages of Bernoulli random variables

Suppose each X_i is an independent flip of a fair coin, so $X_i \sim \text{Bernoulli}(0.5)$ and $\mu = 0.5$. Then \bar{X}_n is the proportion of heads in n flips, and we expect that this proportion is close to 0.5 for large n . Randomness being what it is, this is not guaranteed; for example we could get 1000 heads in 1000 flips, though the probability of this occurring is very small.

So our intuition translates to: **with high probability** the sample average \bar{X}_n is close to the mean 0.5 for large n . We'll demonstrate by doing some calculations in R. You can find the code used for 'class 6 prep' in the usual place on our site.

To start we'll look at the probability of being within 0.1 of the mean. We can express this probability as

$$P(|\bar{X}_n - 0.5| < 0.1) \quad \text{or equivalently} \quad P(0.4 \leq \bar{X}_n \leq 0.6)$$

The law of large numbers says that this probability goes to 1 as the number of flips n gets large. Our R code produces the following values for $P(0.4 \leq \bar{X}_n \leq 0.6)$.

$n = 10:$	$\text{pbinom}(6, 10, 0.5) - \text{pbinom}(3, 10, 0.5)$	$= 0.65625$
$n = 50:$	$\text{pbinom}(30, 50, 0.5) - \text{pbinom}(19, 50, 0.5)$	$= 0.8810795$
$n = 100:$	$\text{pbinom}(60, 100, 0.5) - \text{pbinom}(39, 100, 0.5)$	$= 0.9647998$
$n = 500:$	$\text{pbinom}(300, 500, 0.5) - \text{pbinom}(199, 500, 0.5)$	$= 0.9999941$
$n = 1000:$	$\text{pbinom}(600, 1000, 0.5) - \text{pbinom}(399, 1000, 0.5)$	$= 1$

As predicted by the LoLN the probability goes to 1 as n grows.

We redo the computations to see the probability of being within 0.01 of the mean. Our R code produces the following values for $P(0.49 \leq \bar{X}_n \leq 0.51)$.

n = 10:	<code>pbinom(5, 10, 0.5) - pbinom(4, 10, 0.5)</code>	= 0.2460937
n = 100:	<code>pbinom(51, 100, 0.5) - pbinom(48, 100, 0.5)</code>	= 0.2356466
n = 1000:	<code>pbinom(510, 1000, 0.5) - pbinom(489, 1000, 0.5)</code>	= 0.49334
n = 10000:	<code>pbinom(5100, 10000, 0.5) - pbinom(4899, 10000, 0.5)</code>	= 0.9555742

Again we see the probability of being close to the mean going to 1 as n grows. Since 0.01 is smaller than 0.1 it takes larger values of n to raise the probability to near 1.

This convergence of the probability to 1 is the LoLN in action! Whenever you're confused, it will help you to keep this example in mind. So we see that the LoLN says that [with high probability](#) the average of a large number of independent trials from the same distribution will be very close to the underlying mean of the distribution. Now we're ready for the formal statement.

3.1 Formal statement of the law of large numbers

Theorem (Law of Large Numbers): Suppose $X_1, X_2, \dots, X_n, \dots$ are i.i.d. random variables with mean μ and variance σ^2 . For each n , let \bar{X}_n be the average of the first n variables. Then for any $a > 0$, we have

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| < a) = 1.$$

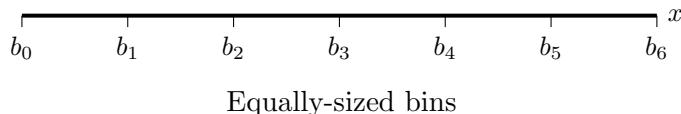
This says precisely that as n increases the probability of being within a of the mean goes to 1. Think of a as a small tolerance of error from the true mean μ . In our example, if we want the probability to be at least $p = 0.99999$ that the proportion of heads \bar{X}_n is within $a = 0.1$ of $\mu = 0.5$, then $n > N = 500$ is large enough. If we decrease the tolerance a and/or increase the probability p , then N will need to be larger.

4 Histograms

We can summarize multiple samples x_1, \dots, x_n of a random variable in a [histogram](#). Here we want to carefully construct histograms so that they resemble the area under the pdf. We will then see how the LoLN applies to histograms.

The step-by-step instructions for constructing a density histogram are as follows.

1. Pick an interval of the real line and divide it into m intervals, with endpoints b_0, b_1, \dots, b_m . Usually these are equally sized, so let's assume this to start.



Each of the intervals is called a [bin](#). For example, in the figure above the first bin is $[b_0, b_1]$ and the last bin is $[b_5, b_6]$. Each bin has a [bin width](#), e.g. $b_1 - b_0$ is the first bin width. Usually the bins all have the same width, called the bin width of the histogram.

2. Place each x_i into the bin that contains its value. If x_i lies on the boundary of two bins, we'll put it in the left bin (this is the R default, though it can be changed).

3. To draw a **frequency histogram**: put a vertical bar above each bin. The **height** of the bar should equal the number of x_i in the bin.
4. To draw a **density histogram**: put a vertical bar above each bin. The **area** of the bar should equal the fraction of all data points that lie in the bin.

Notes:

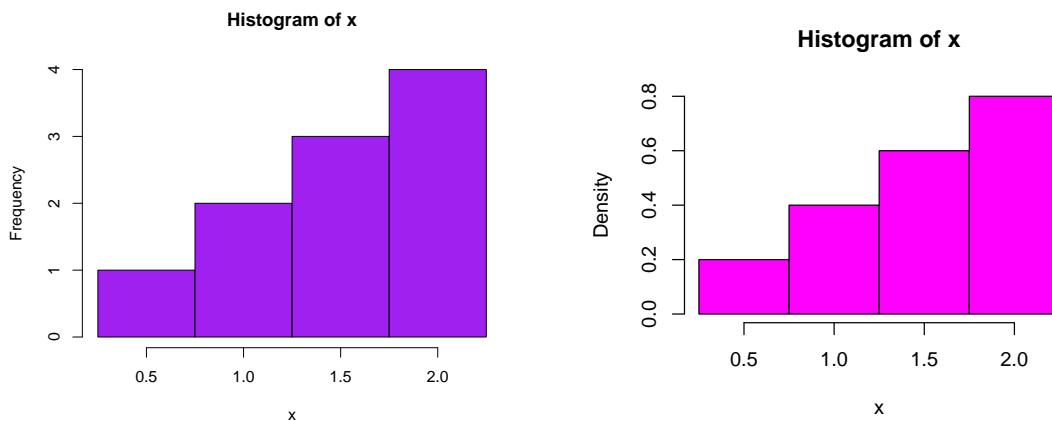
1. When all the bins have the same width, the frequency histogram bars have area proportional to the count. So the density histogram results from simply by dividing the height of each bar by the total area of the frequency histogram. **Ignoring the vertical scale, the two histograms look identical.**

2. Caution: if the bin widths differ, the frequency and density histograms may look very different. There is an example below. Don't let anyone fool you by manipulating bin widths to produce a histogram that suits their mischievous purposes!

In 18.05, we'll stick with equally-sized bins. In general, we prefer the density histogram since its vertical scale is the same as that of the pdf.

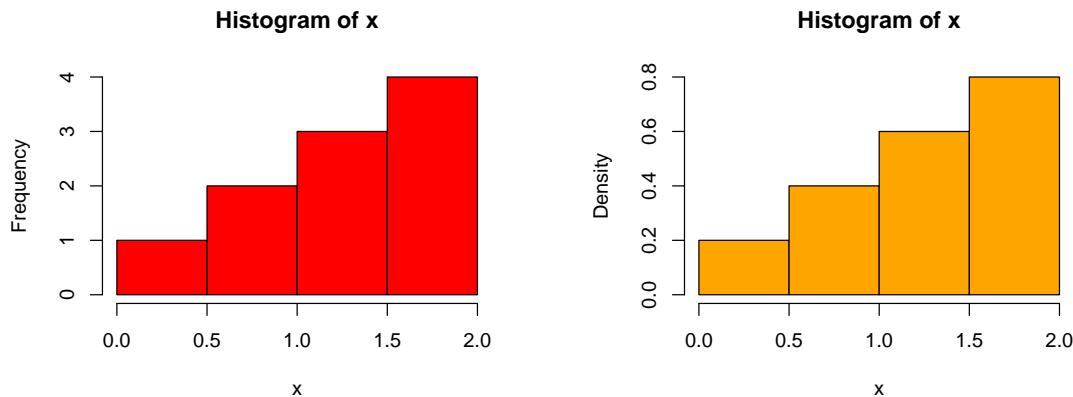
Examples. Here are some examples of histograms, all with the data $[0.5, 1, 1, 1.5, 1.5, 1.5, 2, 2, 2, 2]$. The R code that drew them is in the R file 'class6-prep.r'. You can find the file in the usual place on our site.

1. Here the frequency and density plots look the same but have different vertical scales.



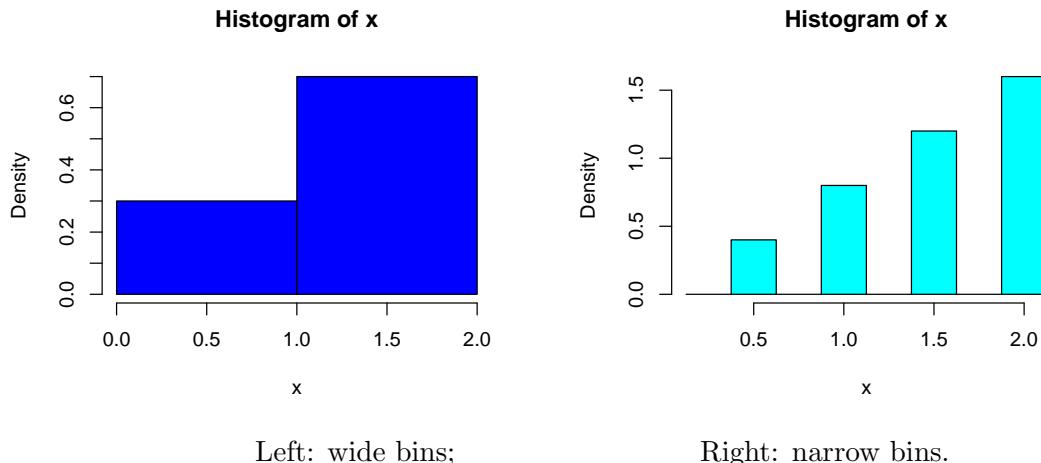
Bins centered at 0.5, 1, 1.5, 2, i.e. width 0.5, bounds at 0.25, 0.75, 1.25, 1.75, 2.25.

2. Here each value is on a bin boundary. Note the values are all on the bin boundaries and are put into the left-hand bin. That is, the bins are **right-closed**, e.g the first bin is for values in the right-closed interval $(0, 0.5]$.



Bin bounds at 0, 0.5, 1, 1.5, 2.

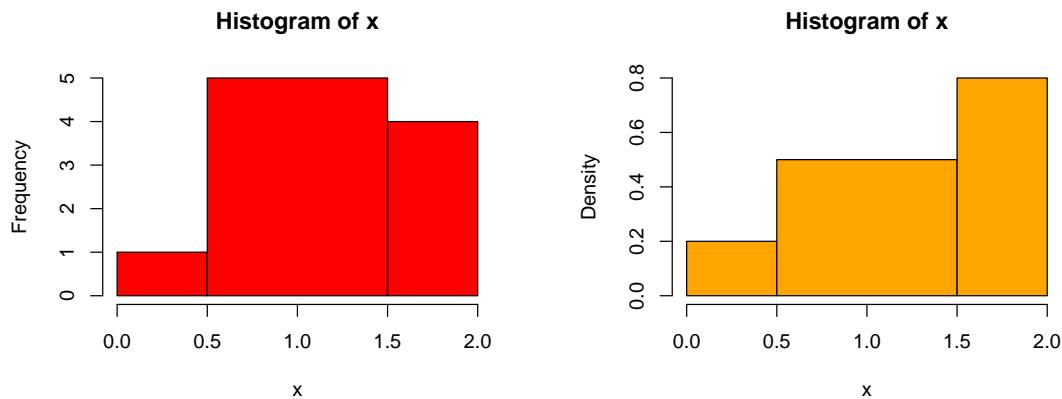
3. Here we show density histograms based on different bin widths. Note that the scale keeps the total area equal to 1. The gaps are bins with zero counts.



Left: wide bins;

Right: narrow bins.

4. Here we use unequal bin widths, so the frequency and density histograms look different



Don't be fooled! These are based on the same data.

The density histogram is the better choice with unequal bin widths. In fact, R will complain

if you try to make a frequency histogram with unequal bin widths. Compare the frequency histogram with unequal bin widths with all the other histograms we drew for this data. It clearly looks different. What happened is that by combining the data in bins $(0.5, 1]$ and $(1, 1.5]$ into one bin $(0.5, 1.5]$ we effectively made the height of both smaller bins greater.

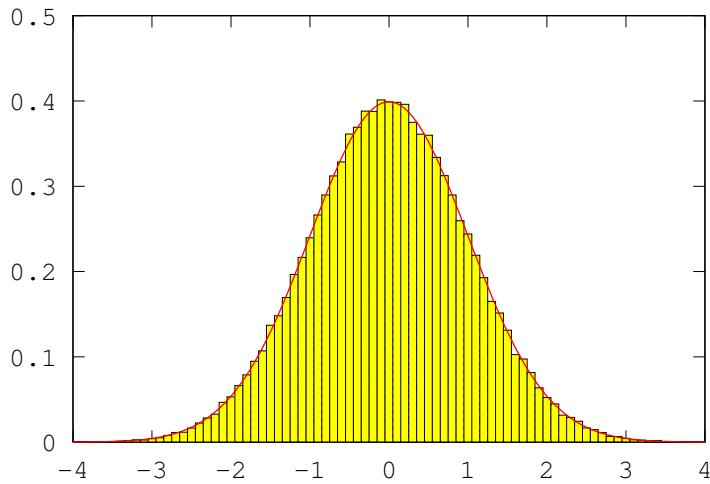
The reason the density histogram is nice is discussed in the next section.

4.1 The law of large numbers and histograms

The law of large number has an important consequence for density histograms.

LoLN for histograms: With high probability the density histogram of a large number of samples from a distribution is a good approximation of the graph of the underlying pdf $f(x)$.

Let's illustrate this by generating a density histogram with bin width 0.1 from 100000 draws from a standard normal distribution. As you can see, the density histogram very closely tracks the graph of the standard normal pdf $\phi(z)$.



Density histogram of 10000 draws from a standard normal distribution, with $\phi(z)$ in red.

5 The Central Limit Theorem

We now prepare for the statement of the CLT.

5.1 Standardization

Given a random variable X with mean μ and standard deviation σ , we define its standardization of X as the new random variable

$$Z = \frac{X - \mu}{\sigma}.$$

Note that Z has mean 0 and standard deviation 1. Note also that if X has a normal distribution, then the standardization of X is the standard normal distribution Z with

mean 0 and variance 1. This explains the term ‘standardization’ and the notation of Z above.

5.2 Statement of the Central Limit Theorem

Suppose $X_1, X_2, \dots, X_n, \dots$ are i.i.d. random variables each having mean μ and standard deviation σ . For each n let S_n denote the sum and let \bar{X}_n be the average of X_1, \dots, X_n .

$$\begin{aligned} S_n &= X_1 + X_2 + \dots + X_n = \sum_{i=1}^n X_i \\ \bar{X}_n &= \frac{X_1 + X_2 + \dots + X_n}{n} = \frac{S_n}{n}. \end{aligned}$$

The properties of mean and variance show

$$\begin{aligned} E(S_n) &= n\mu, & \text{Var}(S_n) &= n\sigma^2, & \sigma_{S_n} &= \sqrt{n}\sigma \\ E(\bar{X}_n) &= \mu, & \text{Var}(\bar{X}_n) &= \frac{\sigma^2}{n}, & \sigma_{\bar{X}_n} &= \frac{\sigma}{\sqrt{n}}. \end{aligned}$$

Since they are multiples of each other, S_n and \bar{X}_n have the same standardization

$$Z_n = \frac{S_n - n\mu}{\sigma\sqrt{n}} = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$$

Central Limit Theorem: For large n ,

$$\bar{X}_n \approx N(\mu, \sigma^2/n), \quad S_n \approx N(n\mu, n\sigma^2), \quad Z_n \approx N(0, 1).$$

Notes: 1. In words: \bar{X}_n is approximately a normal distribution with the same mean as X but a smaller variance.

2. S_n is approximately normal.

3. Standardized \bar{X}_n and S_n are approximately standard normal.

The central limit theorem allows us to approximate a sum or average of i.i.d random variables by a normal random variable. This is extremely useful because it is usually easy to do computations with the normal distribution.

A precise statement of the CLT is that the cdf's of Z_n converge to $\Phi(z)$:

$$\lim_{n \rightarrow \infty} F_{Z_n}(z) = \Phi(z).$$

The proof of the Central Limit Theorem is more technical than we want to get in 18.05. It is accessible to anyone with a decent calculus background.

5.3 Standard Normal Probabilities

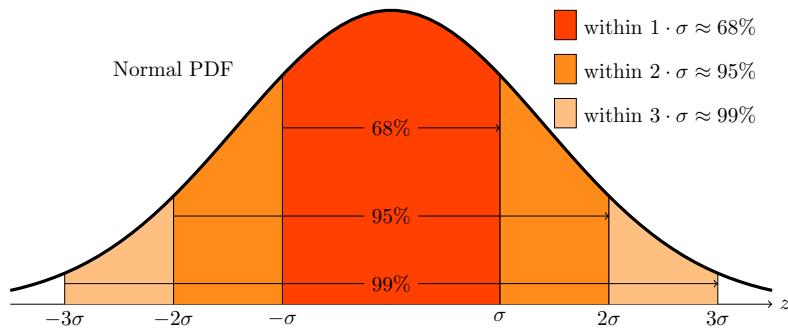
To apply the CLT, we will want to have some normal probabilities at our fingertips. The following probabilities appeared in Class 5. Let $Z \sim N(0, 1)$, a standard normal random variable. Then with rounding we have:

1. $P(|Z| < 1) = 0.68$
2. $P(|Z| < 2) = 0.95$; more precisely $P(|Z| < 1.96) \approx .95$.
3. $P(|Z| < 3) = 0.997$

These numbers are easily compute in R using `pnorm`. However, they are well worth remembering as rules of thumb. You should think of them as:

1. The probability that a normal random variable is within 1 standard deviation of its mean is 0.68.
2. The probability that a normal random variable is within 2 standard deviations of its mean is 0.95.
3. The probability that a normal random variable is within 3 standard deviations of its mean is 0.997.

This is shown graphically in the following figure.



Claim: From these numbers we can derive:

1. $P(Z < 1) \approx 0.84$
2. $P(Z < 2) \approx 0.977$
3. $P(Z < 3) \approx 0.999$

Proof: We know $P(|Z| < 1) = 0.68$. The remaining probability of 0.32 is in the two regions $Z > 1$ and $Z < -1$. These regions are referred to as the [right-hand tail](#) and the [left-hand tail](#) respectively. By symmetry each tail has area 0.16. Thus,

$$P(Z < 1) = P(|Z| < 1) + P(\text{left-hand tail}) = 0.84$$

The other two cases are handled similarly.

5.4 Applications of the CLT

Example 2. Flip a fair coin 100 times. Estimate the probability of more than 55 heads.

answer: Let X_j be the result of the j^{th} flip, so $X_j = 1$ for heads and $X_j = 0$ for tails. The total number of heads is

$$S = X_1 + X_2 + \dots + X_{100}.$$

We know $E(X_j) = 0.5$ and $\text{Var}(X_j) = 1/4$. Since $n = 100$, we have

$$E(S) = 50, \quad \text{Var}(S) = 25 \quad \text{and} \quad \sigma_S = 5.$$

The central limit theorem says that the standardization of S is approximately $N(0, 1)$. The question asks for $P(S > 55)$. Standardizing and using the CLT we get

$$P(S > 55) = P\left(\frac{S - 50}{5} > \frac{55 - 50}{5}\right) \approx P(Z > 1) = 0.16.$$

Here Z is a standard normal random variable and $P(Z > 1) = 1 - P(Z < 1) \approx 0.16$.

Example 3. Estimate the probability of more than 220 heads in 400 flips.

answer: This is nearly identical to the previous example. Now $\mu_S = 200$ and $\sigma_S = 10$ and we want $P(S > 220)$. Standardizing and using the CLT we get:

$$P(S > 220) = P\left(\frac{S - \mu_S}{\sigma_S} > \frac{220 - 200}{10}\right) \approx P(Z > 2) = .025.$$

Again, $Z \sim N(0, 1)$ and the rules of thumb show $P(Z > 2) = .025$.

Note: Even though $55/100 = 220/400$, the probability of more than 55 heads in 100 flips is larger than the probability of more than 220 heads in 400 flips. This is due to the LoLN and the larger value of n in the latter case.

Example 4. Estimate the probability of between 40 and 60 heads in 100 flips.

answer: As in the first example, $E(S) = 50$, $\text{Var}(S) = 25$ and $\sigma_S = 5$. So

$$P(40 \leq S \leq 60) = P\left(\frac{40 - 100}{5} \leq \frac{S - 50}{5} \leq \frac{60 - 50}{5}\right) \approx P(-2 \leq Z \leq 2)$$

We can compute the right-hand side using our rule of thumb. For a more accurate answer we use R:

$$\text{pnorm}(2) - \text{pnorm}(-2) = 0.954\dots$$

Recall that in Section 3 we used the binomial distribution to compute an answer of 0.965.... So our approximate answer using CLT is off by about 1%.

Think: Would you expect the CLT method to give a better or worse approximation of $P(200 < S < 300)$ with $n = 500$?

We encourage you to check your answer using R.

Example 5. Polling. When taking a political poll the results are often reported as a number with a margin of error. For example $52\% \pm 3\%$ favor candidate A. The rule of thumb is that if you poll n people then the margin of error is $\pm 1/\sqrt{n}$. We will now see exactly what this means and that it is an application of the central limit theorem.

Suppose there are 2 candidates A and B. Suppose further that the fraction of the population who prefer A is p_0 . That is, if you ask a random person who they prefer then the probability they'll answer A is p_0 .

To run the poll a pollster selects n people at random and asks 'Do you support candidate A or candidate B. Thus we can view the poll as a sequence of n independent Bernoulli(p_0) trials, X_1, X_2, \dots, X_n , where X_i is 1 if the i^{th} person prefers A and 0 if they prefer B. The fraction of people polled that prefer A is just the average \bar{X} .

We know that each $X_i \sim \text{Bernoulli}(p_0)$ so,

$$E(X_i) = p_0 \quad \text{and} \quad \sigma_{X_i} = \sqrt{p_0(1 - p_0)}.$$

Therefore, the central limit theorem tells us that

$$\bar{X} \approx N(p_0, \sigma/\sqrt{n}), \quad \text{where } \sigma = \sqrt{p_0(1 - p_0)}.$$

In a normal distribution 95% of the probability is within 2 standard deviations of the mean. This means that in 95% of polls of n people the sample mean \bar{X} will be within $2\sigma/\sqrt{n}$ of the true mean p_0 . The final step is to note that for any value of p_0 we have $\sigma \leq 1/2$. (It is an easy calculus exercise to see that $1/4$ is the maximum value of $\sigma^2 = p_0(1 - p_0)$.) This means that we can conservatively say that in 95% of polls of n people the sample mean \bar{X} is within $1/\sqrt{n}$ of the true mean. The frequentist statistician then takes the interval $\bar{X} \pm 1/\sqrt{n}$ and calls it the 95% confidence interval for p_0 .

A word of caution: it is tempting and common, but wrong, to think that there is a 95% probability the true fraction p_0 is in the confidence interval. This is subtle, but the error is the same one as thinking you have a disease if a 95% accurate test comes back positive. It's true that 95% of people taking the test get the correct result. It's not necessarily true that 95% of positive tests are correct.

5.5 Why use the CLT

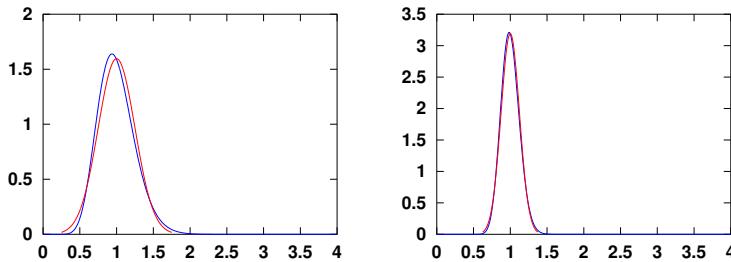
Since the probabilities in the above examples can be computed exactly using the binomial distribution, you may be wondering what is the point of finding an approximate answer using the CLT. In fact, we were only able to compute these probabilities exactly because the X_i were Bernoulli and so the sum S was binomial. In general, the distribution of the S will not be familiar, so you will not be able to compute the probabilities for S exactly; it can also happen that the exact computation is possible in theory but too computationally intensive in practice, even for a computer. The power of the CLT is that it applies when X_i has almost any distribution. Though we will see in the next section that some distributions may require larger n for the approximation to be a good one).

5.6 How big does n have to be to apply the CLT?

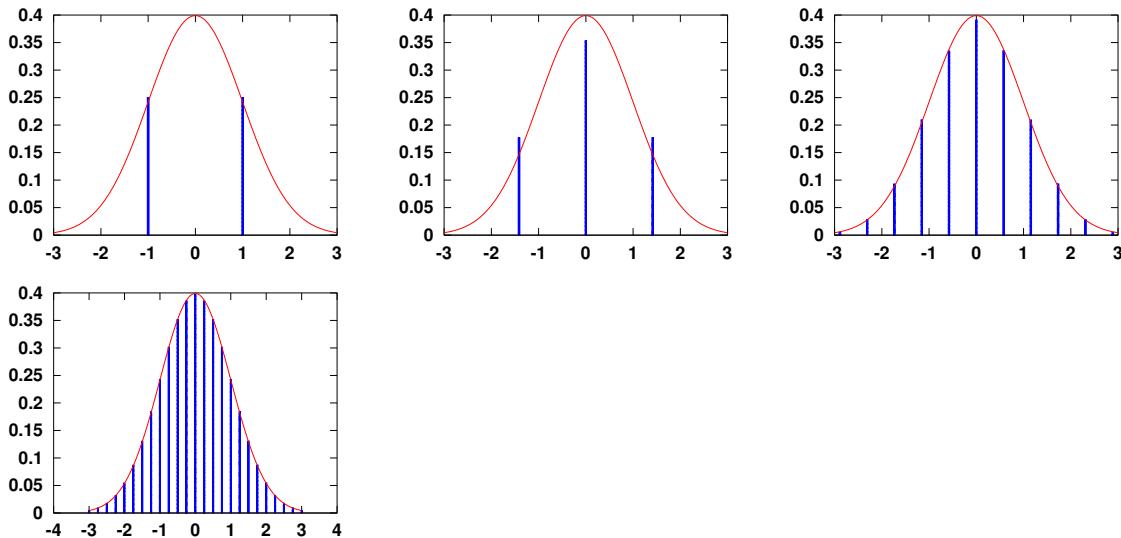
Short answer: often, not that big.

The following sequences of pictures show the convergence of averages to a normal distribution.

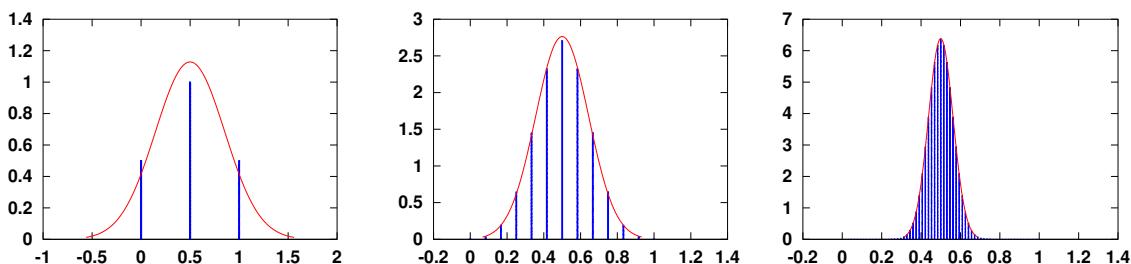
First we show the standardized average of n i.i.d. **uniform** random variables with $n = 1, 2, 4, 8, 12$. The pdf of the average is in blue and the standard normal pdf is in red. By the time $n = 12$ the fit between the standardized average and the true normal looks very good.



The central limit theorem works for discrete variables also. Here is the standardized average of n i.i.d. Bernoulli(.5) random variables with $n = 1, 2, 12, 64$. Notice that as n grows, the average can take more values, which allows the discrete distribution to 'fill in' the normal density.



Finally we show the (non-standardized) average of n Bernoulli(.5) random variables, with $n = 4, 12, 64$. Notice how the standard deviation gets smaller resulting in a spikier (more peaked) density.



Appendix

Class 6, 18.05

Jeremy Orloff and Jonathan Bloom

1 Introduction

In this appendix we give more formal mathematical material that is not strictly a part of 18.05. This will not be on homework or tests. We give this material to emphasize that in doing mathematics we should be careful to specify our hypotheses completely and give clear deductive arguments to prove our claims. We hope you find it interesting and illuminating.

2 With high probability the density histogram resembles the graph of the probability density function:

We stated that one consequence of the law of large numbers is that as the number of samples increases the density histogram of the samples has an increasing probability of matching the graph of the underlying pdf or pmf. This is a good rule of thumb, but it is rather imprecise. It is possible to make more precise statements. It will take some care to make a sensible and precise statement, which will not be quite so sweeping.

Suppose we have an experiment that produces data according to the random variable X and suppose we generate n independent samples from X . Call them

$$x_1, x_2, \dots, x_n.$$

By a bin we mean a range of values, i.e. $[x_k, x_{k+1})$. To make a density histogram of the data we divide the range of X into m bins and calculate the fraction of the data in each bin.

Now, let p_k be the probability a random data point is in the k th bin. This is the probability for an [indicator](#) (Bernoulli) random variable $B_{k,j}$ which is 1 if the j th data point is in the bin and 0 otherwise.

Statement 1. Let \bar{p}_k be the fraction of the data in bin k . As the number n of data points gets large the probability that \bar{x}_k is close to p_k approaches 1. Said differently, given any small number, call it a the probability $P(|\bar{p}_k - p_k| < a)$ depends on n , and as n goes to infinity this probability goes to 1.

Proof. Let \bar{B}_k be the average of $B_{k,j}$. Since $E(B_{k,j}) = p_k$, the law of large number says exactly that

$$P(|\bar{B}_k - p_k| < a) \quad \text{approaches 1 as } n \text{ goes to infinity.}$$

But, since the $B_{k,j}$ are indicator variables, their average is exactly \bar{p}_k , the fraction of the data in bin k . Replacing \bar{B}_k by \bar{p}_k in the above equation gives

$$P(|\bar{p}_k - p_k| < a) \quad \text{approaches 1 as } n \text{ goes to infinity.}$$

This is exactly what statement 1 claimed.

Statement 2. The same statement holds for a finite number of bins simultaneously. That is, for bins 1 to m we have

$$P((|\bar{B}_1 - p_1| < a), (|\bar{B}_2 - p_2| < a), \dots, (|\bar{B}_m - p_m| < a)) \text{ approaches 1 as } n \text{ goes to infinity.}$$

Proof. First we note the following probability rule, which is a consequence of the inclusion exclusion principle: If two events A and B have $P(A) = 1 - \alpha_1$ and $P(B) = 1 - \alpha_2$ then $P(A \cap B) \geq 1 - (\alpha_1 + \alpha_2)$.

Now, Statement 1 says that for any α we can find n large enough that $P(|\bar{B}_k - p_k| < a) > 1 - \alpha/m$ for each bin separately. By the probability rule, the probability of the intersection of all these events is at least $1 - \alpha$. Since we can let α be as small as we want by letting n go to infinity, in the limit we get probability 1 as claimed.

Statement 3. If $f(x)$ is a continuous probability density with range $[a, b]$ then by taking enough data and having a small enough bin width we can insure that with high probability the density histogram is as close as we want to the graph of $f(x)$.

Proof. We will only sketch the argument. Assume the bin around x has width is Δx . If Δx is small enough then the probability a data point is in the bin is approximately $f(x)\Delta x$. Statement 2 guarantees that if n is large enough then with high probability the fraction of data in the bin is also approximately $f(x)\Delta x$. Since this is the area of the bin we see that its height will be approximately $f(x)$. That is, with high probability the height of the histogram over any point x is close to $f(x)$. This is what Statement 3 claimed.

Note. If the range is infinite or the density goes to infinity at some point we need to be more careful. There are statements we could make for these cases.

3 The Chebyshev inequality

One proof of the LoLN follows from the following key inequality.

The Chebyshev inequality. Suppose Y is a random variable with mean μ and variance σ^2 . Then for any positive value a , we have

$$P(|Y - \mu| \geq a) \leq \frac{\text{Var}(Y)}{a^2}.$$

In words, the Chebyshev inequality says that the probability that Y differs from the mean by more than a is bounded by $\text{Var}(Y)/a^2$. Morally, the smaller the variance of Y , the smaller the probability that Y is far from its mean.

Proof of the LoLN: Since $\text{Var}(\bar{X}_n) = \text{Var}(X)/n$, the variance of the average \bar{X}_n goes to zero as n goes to infinity. So the Chebyshev inequality for $Y = \bar{X}_n$ and fixed a implies that as n grows, the probability that \bar{X}_n is farther than a from μ goes to 0. Hence the probability that \bar{X}_n is within a of μ goes to 1, which is the LoLN.

Proof of the Chebyshev inequality: The proof is essentially the same for discrete and continuous Y . We'll assume Y is continuous and also that $\mu = 0$, since replacing Y by

$Y - \mu$ does not change the variance. So

$$\begin{aligned} P(|Y| \geq a) &= \int_{-\infty}^{-a} f(y) dy + \int_a^{\infty} f(y) dy \leq \int_{-\infty}^{-a} \frac{y^2}{a^2} f(y) dy + \int_a^{\infty} \frac{y^2}{a^2} f(y) dy \\ &\leq \int_{-\infty}^{\infty} \frac{y^2}{a^2} f(y) dy = \frac{\text{Var}(Y)}{a^2}. \end{aligned}$$

The first inequality uses that $y^2/a^2 \geq 1$ on the intervals of integration. The second inequality follows because including the range $[-a, a]$ only makes the integral larger, since the integrand is positive.

4 The need for variance

We didn't lie to you, but we did gloss over one technical fact. Throughout we assumed that the underlying distributions had a variance. For example, the proof of the law of large numbers made use of the variance by way of the Chebyshev inequality. But there are distributions which do not have a variance because the sum or integral for the variance does not converge to a finite number. For such distributions the law of large numbers may not be true. In 18.05 we won't have to worry about this, but if you go deeper into statistics this may become important.

Joint Distributions, Independence

Class 7, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Understand what is meant by a **joint** pmf, pdf and cdf of two random variables.
2. Be able to compute probabilities and marginals from a joint pmf or pdf.
3. Be able to test whether two random variables are independent.

2 Introduction

In science and in real life, we are often interested in two (or more) random variables at the same time. For example, we might measure the height and weight of giraffes, or the IQ and birthweight of children, or the frequency of exercise and the rate of heart disease in adults, or the level of air pollution and rate of respiratory illness in cities, or the number of Facebook friends and the age of Facebook members.

Think: What relationship would you expect in each of the five examples above? Why?

In such situations the random variables have a **joint distribution** that allows us to compute probabilities of events involving both variables and understand the relationship between the variables. This is simplest when the variables are **independent**. When they are not, we use **covariance** and **correlation** as measures of the nature of the dependence between them.

3 Joint Distribution

3.1 Discrete case

Suppose X and Y are two discrete random variables and that X takes values $\{x_1, x_2, \dots, x_n\}$ and Y takes values $\{y_1, y_2, \dots, y_m\}$. The ordered pair (X, Y) take values in the product $\{(x_1, y_1), (x_1, y_2), \dots, (x_n, y_m)\}$. The **joint probability mass function** (joint pmf) of X and Y is the function $p(x_i, y_j)$ giving the probability of the joint outcome $X = x_i, Y = y_j$.

We organize this in a **joint probability table** as shown:

$X \setminus Y$	y_1	y_2	\dots	y_j	\dots	y_m
x_1	$p(x_1, y_1)$	$p(x_1, y_2)$	\dots	$p(x_1, y_j)$	\dots	$p(x_1, y_m)$
x_2	$p(x_2, y_1)$	$p(x_2, y_2)$	\dots	$p(x_2, y_j)$	\dots	$p(x_2, y_m)$
\dots	\dots	\dots	\dots	\dots	\dots	\dots
\dots	\dots	\dots	\dots	\dots	\dots	\dots
x_i	$p(x_i, y_1)$	$p(x_i, y_2)$	\dots	$p(x_i, y_j)$	\dots	$p(x_i, y_m)$
\dots	\dots	\dots	\dots	\dots	\dots	\dots
x_n	$p(x_n, y_1)$	$p(x_n, y_2)$	\dots	$p(x_n, y_j)$	\dots	$p(x_n, y_m)$

Example 1. Roll two dice. Let X be the value on the first die and let Y be the value on the second die. Then both X and Y take values 1 to 6 and the joint pmf is $p(i, j) = 1/36$ for all i and j between 1 and 6. Here is the joint probability table:

$X \setminus Y$	1	2	3	4	5	6
1	1/36	1/36	1/36	1/36	1/36	1/36
2	1/36	1/36	1/36	1/36	1/36	1/36
3	1/36	1/36	1/36	1/36	1/36	1/36
4	1/36	1/36	1/36	1/36	1/36	1/36
5	1/36	1/36	1/36	1/36	1/36	1/36
6	1/36	1/36	1/36	1/36	1/36	1/36

Example 2. Roll two dice. Let X be the value on the first die and let T be the total on both dice. Here is the joint probability table:

$X \setminus T$	2	3	4	5	6	7	8	9	10	11	12
1	1/36	1/36	1/36	1/36	1/36	1/36	0	0	0	0	0
2	0	1/36	1/36	1/36	1/36	1/36	1/36	0	0	0	0
3	0	0	1/36	1/36	1/36	1/36	1/36	1/36	0	0	0
4	0	0	0	1/36	1/36	1/36	1/36	1/36	1/36	0	0
5	0	0	0	0	1/36	1/36	1/36	1/36	1/36	1/36	0
6	0	0	0	0	0	1/36	1/36	1/36	1/36	1/36	1/36

A joint probability mass function must satisfy two properties:

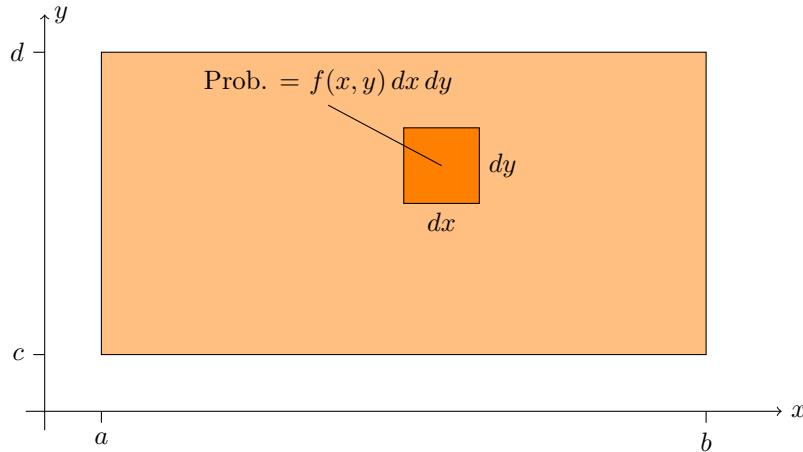
1. $0 \leq p(x_i, y_j) \leq 1$
2. The total probability is 1. We can express this as a [double sum](#):

$$\sum_{i=1}^n \sum_{j=1}^m p(x_i, y_j) = 1$$

3.2 Continuous case

The continuous case is essentially the same as the discrete case: we just replace discrete sets of values by continuous intervals, the joint probability mass function by a [joint probability density function](#), and the sums by integrals.

If X takes values in $[a, b]$ and Y takes values in $[c, d]$ then the pair (X, Y) takes values in the product $[a, b] \times [c, d]$. The [joint probability density function](#) (joint pdf) of X and Y is a function $f(x, y)$ giving the probability density at (x, y) . That is, the probability that (X, Y) is in a small rectangle of width dx and height dy around (x, y) is $f(x, y) dx dy$.



A joint probability density function must satisfy two properties:

1. $0 \leq f(x, y)$
2. The total probability is 1. We now express this as a [double integral](#):

$$\int_c^d \int_a^b f(x, y) dx dy = 1$$

Note: as with the pdf of a single random variable, the joint pdf $f(x, y)$ can take values greater than 1; it is a probability density, **not** a probability.

In 18.05 we won't expect you to be experts at double integration. Here's what we will expect.

- You should understand double integrals conceptually as double sums.
- You should be able to compute double integrals over rectangles.
- For a non-rectangular region, when $f(x, y) = c$ is constant, you should know that the double integral is the same as the $c \times$ (the area of the region).

3.3 Events

Random variables are useful for describing events. Recall that an event is a set of outcomes and that random variables assign numbers to outcomes. For example, the event ' $X > 1$ ' is the set of all outcomes for which X is greater than 1. These concepts readily extend to pairs of random variables and joint outcomes.

Example 3. In Example 1, describe the event $B = 'Y - X \geq 2'$ and find its probability.

answer: We can describe B as a set of (X, Y) pairs:

$$B = \{(1, 3), (1, 4), (1, 5), (1, 6), (2, 4), (2, 5), (2, 6), (3, 5), (3, 6), (4, 6)\}.$$

We can also describe it visually

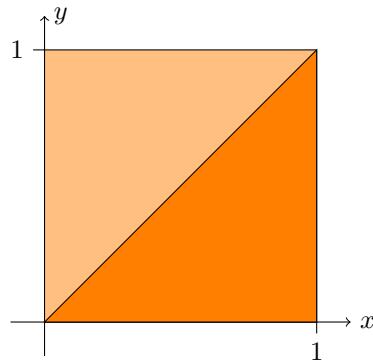
$X \setminus Y$	1	2	3	4	5	6
1	1/36	1/36	1/36	1/36	1/36	1/36
2	1/36	1/36	1/36	1/36	1/36	1/36
3	1/36	1/36	1/36	1/36	1/36	1/36
4	1/36	1/36	1/36	1/36	1/36	1/36
5	1/36	1/36	1/36	1/36	1/36	1/36
6	1/36	1/36	1/36	1/36	1/36	1/36

The event B consists of the outcomes in the shaded squares.

The probability of B is the sum of the probabilities in the orange shaded squares, so $P(B) = 10/36$.

Example 4. Suppose X and Y both take values in $[0,1]$ with uniform density $f(x, y) = 1$. Visualize the event ' $X > Y$ ' and find its probability.

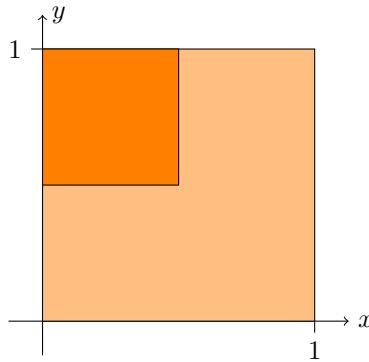
answer: Jointly X and Y take values in the unit square. The event ' $X > Y$ ' corresponds to the shaded lower-right triangle below. Since the density is constant, the probability is just the fraction of the total area taken up by the event. In this case, it is clearly 0.5.



The event ' $X > Y$ ' in the unit square.

Example 5. Suppose X and Y both take values in $[0,1]$ with density $f(x, y) = 4xy$. Show $f(x, y)$ is a valid joint pdf, visualize the event $A = 'X < 0.5 \text{ and } Y > 0.5'$ and find its probability.

answer: Jointly X and Y take values in the unit square.



The event A in the unit square.

To show $f(x, y)$ is a valid joint pdf we must check that it is positive (which it clearly is) and that the total probability is 1.

$$\text{Total probability} = \int_0^1 \int_0^1 4xy \, dx \, dy = \int_0^1 [2x^2y]_0^1 \, dy = \int_0^1 2y \, dy = 1. \quad \text{QED}$$

The event A is just the upper-left-hand quadrant. Because the density is not constant we must compute an integral to find the probability.

$$P(A) = \int_0^{.5} \int_{.5}^1 4xy \, dy \, dx = \int_0^{.5} [2xy^2]_{.5}^1 \, dx = \int_0^{.5} \frac{3x}{2} \, dx = \boxed{\frac{3}{16}}.$$

3.4 Joint cumulative distribution function

Suppose X and Y are jointly-distributed random variables. We will use the notation ‘ $X \leq x, Y \leq y$ ’ to mean the event ‘ $X \leq x$ and $Y \leq y$ ’. The **joint cumulative distribution function** (joint cdf) is defined as

$$F(x, y) = P(X \leq x, Y \leq y)$$

Continuous case: If X and Y are continuous random variables with joint density $f(x, y)$ over the range $[a, b] \times [c, d]$ then the joint cdf is given by the double integral

$$F(x, y) = \int_c^y \int_a^x f(u, v) \, du \, dv.$$

To recover the joint pdf, we differentiate the joint cdf. Because there are two variables we need to use partial derivatives:

$$f(x, y) = \frac{\partial^2 F}{\partial x \partial y}(x, y).$$

Discrete case: If X and Y are discrete random variables with joint pmf $p(x_i, y_j)$ then the joint cdf is give by the double sum

$$F(x, y) = \sum_{x_i \leq x} \sum_{y_j \leq y} p(x_i, y_j).$$

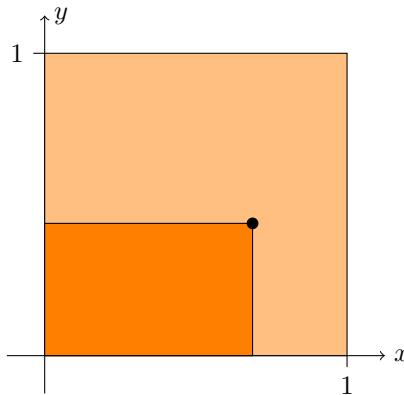
3.5 Properties of the joint cdf

The joint cdf $F(x, y)$ of X and Y must satisfy several properties:

1. $F(x, y)$ is non-decreasing: i.e. if x or y increase then $F(x, y)$ must stay constant or increase.
2. $F(x, y) = 0$ at the lower-left of the joint range.
If the lower left is $(-\infty, -\infty)$ then this means $\lim_{(x,y) \rightarrow (-\infty, -\infty)} F(x, y) = 0$.
3. $F(x, y) = 1$ at the upper-right of the joint range.
If the upper-right is (∞, ∞) then this means $\lim_{(x,y) \rightarrow (\infty, \infty)} F(x, y) = 1$.

Example 6. Find the joint cdf for the random variables in Example 5.

answer: The event ' $X \leq x$ and $Y \leq y$ ' is a rectangle in the unit square.



To find the cdf $F(x, y)$ we compute a double integral:

$$F(x, y) = \int_0^y \int_0^x 4uv \, du \, dv = \boxed{x^2 y^2}.$$

Example 7. In Example 1, compute $F(3.5, 4)$.

answer: We redraw the joint probability table. Notice how similar the picture is to the one in the previous example.

$F(3.5, 4)$ is the probability of the event ' $X \leq 3.5$ and $Y \leq 4$ '. We can visualize this event as the shaded rectangles in the table:

$X \setminus Y$	1	2	3	4	5	6
1	1/36	1/36	1/36	1/36	1/36	1/36
2	1/36	1/36	1/36	1/36	1/36	1/36
3	1/36	1/36	1/36	1/36	1/36	1/36
4	1/36	1/36	1/36	1/36	1/36	1/36
5	1/36	1/36	1/36	1/36	1/36	1/36
6	1/36	1/36	1/36	1/36	1/36	1/36

The event ' $X \leq 3.5$ and $Y \leq 4$ '.

Adding up the probability in the shaded squares we get $F(3.5, 4) = 12/36 = 1/3$.

Note. One unfortunate difference between the continuous and discrete visualizations is that for continuous variables the value increases as we go up in the vertical direction while the opposite is true for the discrete case. We have experimented with changing the discrete tables to match the continuous graphs, but it causes too much confusion. We will just have to live with the difference!

3.6 Marginal distributions

When X and Y are jointly-distributed random variables, we may want to consider only one of them, say X . In that case we need to find the pmf (or pdf or cdf) of X without Y . This is called a **marginal pmf** (or pdf or cdf). The next example illustrates the way to compute this and the reason for the term ‘marginal’.

3.7 Marginal pmf

Example 8. In Example 2 we rolled two dice and let X be the value on the first die and T be the total on both dice. Compute the marginal pmf of X and of T .

answer: In the table each row represents a single value of X . So the event ' $X = 3$ ' is the third row of the table. To find $P(X = 3)$ we simply have to sum up the probabilities in this row. We put the sum in the **right-hand margin** of the table. Likewise $P(T = 5)$ is just the sum of the column with $T = 5$. We put the sum in the **bottom margin** of the table.

$X \setminus T$	2	3	4	5	6	7	8	9	10	11	12	$p(x_i)$
1	1/36	1/36	1/36	1/36	1/36	1/36	0	0	0	0	0	1/6
2	0	1/36	1/36	1/36	1/36	1/36	1/36	0	0	0	0	1/6
3	0	0	1/36	1/36	1/36	1/36	1/36	1/36	0	0	0	1/6
4	0	0	0	1/36	1/36	1/36	1/36	1/36	1/36	0	0	1/6
5	0	0	0	0	1/36	1/36	1/36	1/36	1/36	1/36	0	1/6
6	0	0	0	0	0	1/36	1/36	1/36	1/36	1/36	1/36	1/6
$p(t_j)$	1/36	2/36	3/36	4/36	5/36	6/36	5/36	4/36	3/36	2/36	1/36	1

Computing the marginal probabilities $P(X = 3) = 1/6$ and $P(T = 5) = 4/36$.

Note: Of course in this case we already knew the pmf of X and of T . It is good to see that our computation here is in agreement!

As motivated by this example, marginal pmf's are obtained from the joint pmf by summing:

$$p_X(x_i) = \sum_j p(x_i, y_j), \quad p_Y(y_j) = \sum_i p(x_i, y_j)$$

The term **marginal** refers to the fact that the values are written in the margins of the table.

3.8 Marginal pdf

For a continuous joint density $f(x, y)$ with range $[a, b] \times [c, d]$, the marginal pdf's are:

$$f_X(x) = \int_c^d f(x, y) dy, \quad f_Y(y) = \int_a^b f(x, y) dx.$$

Compare these with the marginal pmf's above; as usual the sums are replaced by integrals. We say that to obtain the marginal for X , we **integrate out** Y from the joint pdf and vice versa.

Example 9. Suppose (X, Y) takes values on the square $[0, 1] \times [1, 2]$ with joint pdf $f(x, y) = \frac{8}{3}x^3y$. Find the marginal pdf's $f_X(x)$ and $f_Y(y)$.

answer: To find $f_X(x)$ we integrate out y and to find $f_Y(y)$ we integrate out x .

$$\begin{aligned} f_X(x) &= \int_1^2 \frac{8}{3}x^3y dy = \left[\frac{4}{3}x^3y^2 \right]_1^2 = \boxed{4x^3} \\ f_Y(y) &= \int_0^1 \frac{8}{3}x^3y dx = \left[\frac{2}{3}x^4y \right]_0^1 = \boxed{\frac{2}{3}y}. \end{aligned}$$

Example 10. Suppose (X, Y) takes values on the unit square $[0, 1] \times [0, 1]$ with joint pdf $f(x, y) = \frac{3}{2}(x^2 + y^2)$. Find the marginal pdf $f_X(x)$ and use it to find $P(X < 0.5)$.

answer:

$$\begin{aligned} f_X(x) &= \int_0^1 \frac{3}{2}(x^2 + y^2) dy = \left[\frac{3}{2}x^2y + \frac{y^3}{2} \right]_0^1 = \boxed{\frac{3}{2}x^2 + \frac{1}{2}}. \\ P(X < 0.5) &= \int_0^{0.5} f_X(x) dx = \int_0^{0.5} \frac{3}{2}x^2 + \frac{1}{2} dx = \left[\frac{1}{2}x^3 + \frac{1}{2}x \right]_0^{0.5} = \boxed{\frac{5}{16}}. \end{aligned}$$

3.9 Marginal cdf

Finding the marginal cdf from the joint cdf is easy. If X and Y jointly take values on $[a, b] \times [c, d]$ then

$$F_X(x) = F(x, d), \quad F_Y(y) = F(b, y).$$

If d is ∞ then this becomes a limit $F_X(x) = \lim_{y \rightarrow \infty} F(x, y)$. Likewise for $F_Y(y)$.

Example 11. The joint cdf in the last example was $F(x, y) = \frac{1}{2}(x^3 + xy^3)$ on $[0, 1] \times [0, 1]$. Find the marginal cdf's and use $F_X(x)$ to compute $P(X < 0.5)$.

answer: We have $F_X(x) = F(x, 1) = \frac{1}{2}(x^3 + x)$ and $F_Y(y) = F(1, y) = \frac{1}{2}(y + y^3)$. So $P(X < 0.5) = F_X(0.5) = \frac{1}{2}(0.5^3 + 0.5) = \frac{5}{16}$: exactly the same as before.

3.10 3D visualization

We visualized $P(a < X < b)$ as the area under the pdf $f(x)$ over the interval $[a, b]$. Since the range of values of (X, Y) is already a two dimensional region in the plane, the graph of

$f(x, y)$ is a surface over that region. We can then visualize probability as **volume** under the surface.

Think: Summoning your inner artist, sketch the graph of the joint pdf $f(x, y) = 4xy$ and visualize the probability $P(A)$ as a volume for Example 5.

4 Independence

We are now ready to give a careful mathematical definition of independence. Of course, it will simply capture the notion of independence we have been using up to now. But, it is nice to finally have a solid definition that can support complicated probabilistic and statistical investigations.

Recall that events A and B are independent if

$$P(A \cap B) = P(A)P(B).$$

Random variables X and Y define events like ' $X \leq 2$ ' and ' $Y > 5$ '. So, X and Y are independent if **any** event defined by X is independent of **any** event defined by Y . The formal definition that guarantees this is the following.

Definition: Jointly-distributed random variables X and Y are **independent** if their joint cdf is the product of the marginal cdf's:

$$F(X, Y) = F_X(x)F_Y(y).$$

For discrete variables this is equivalent to the joint pmf being the product of the marginal pmf's.:

$$p(x_i, y_j) = p_X(x_i)p_Y(y_j).$$

For continuous variables this is equivalent to the joint pdf being the product of the marginal pdf's.:

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

Once you have the joint distribution, checking for independence is usually straightforward although it can be tedious.

Example 12. For **discrete variables** independence means the probability in a cell must be the product of the marginal probabilities of its row and column. In the first table below this is true: every marginal probability is $1/6$ and every cell contains $1/36$, i.e. the product of the marginals. Therefore X and Y are independent.

In the second table below most of the cell probabilities are not the product of the marginal probabilities. For example, none of marginal probabilities are 0, so none of the cells with 0 probability can be the product of the marginals.

$X \setminus Y$	1	2	3	4	5	6	$p(x_i)$
$p(y_j)$	1/6	1/6	1/6	1/6	1/6	1/6	1
1	1/36	1/36	1/36	1/36	1/36	1/36	1/6
2	1/36	1/36	1/36	1/36	1/36	1/36	1/6
3	1/36	1/36	1/36	1/36	1/36	1/36	1/6
4	1/36	1/36	1/36	1/36	1/36	1/36	1/6
5	1/36	1/36	1/36	1/36	1/36	1/36	1/6
6	1/36	1/36	1/36	1/36	1/36	1/36	1/6

$X \setminus T$	2	3	4	5	6	7	8	9	10	11	12	$p(x_i)$
$p(y_j)$	1/36	2/36	3/36	4/36	5/36	6/36	5/36	4/36	3/36	2/36	1/36	1
1	1/36	1/36	1/36	1/36	1/36	1/36	0	0	0	0	0	1/6
2	0	1/36	1/36	1/36	1/36	1/36	1/36	0	0	0	0	1/6
3	0	0	1/36	1/36	1/36	1/36	1/36	1/36	0	0	0	1/6
4	0	0	0	1/36	1/36	1/36	1/36	1/36	1/36	0	0	1/6
5	0	0	0	0	1/36	1/36	1/36	1/36	1/36	1/36	0	1/6
6	0	0	0	0	0	1/36	1/36	1/36	1/36	1/36	1/36	1/6

Example 13. For continuous variables independence means you can factor the joint pdf or cdf as the product of a function of x and a function of y .

- (i) Suppose X has range $[0, 1/2]$, Y has range $[0, 1]$ and $f(x, y) = 96x^2y^3$ then X and Y are independent. The marginal densities are $f_X(x) = 24x^2$ and $f_Y(y) = 4y^3$.
- (ii) If $f(x, y) = 1.5(x^2 + y^2)$ over the unit square then X and Y are not independent because there is no way to factor $f(x, y)$ into a product $f_X(x)f_Y(y)$.
- (iii) If $F(x, y) = \frac{1}{2}(x^3y + xy^3)$ over the unit square then X and Y are not independent because the cdf does not factor into a product $F_X(x)F_Y(y)$.

Covariance and Correlation

Class 7, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Understand the meaning of covariance and correlation.
2. Be able to compute the covariance and correlation of two random variables.

2 Covariance

Covariance is a measure of how much two random variables vary together. For example, height and weight of giraffes have positive covariance because when one is big the other tends also to be big.

Definition: Suppose X and Y are random variables with means μ_X and μ_Y . The **covariance** of X and Y is defined as

$$\text{Cov}(X, Y) = E((X - \mu_X)(Y - \mu_Y)).$$

2.1 Properties of covariance

1. $\text{Cov}(aX + b, cY + d) = ac\text{Cov}(X, Y)$ for constants a, b, c, d .
2. $\text{Cov}(X_1 + X_2, Y) = \text{Cov}(X_1, Y) + \text{Cov}(X_2, Y)$.
3. $\text{Cov}(X, X) = \text{Var}(X)$
4. $\text{Cov}(X, Y) = E(XY) - \mu_X\mu_Y$.
5. $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$ for any X and Y .
6. If X and Y are independent then $\text{Cov}(X, Y) = 0$.

Warning: The converse is false: zero covariance does not always imply independence.

Notes. 1. Property 4 is like the similar property for variance. Indeed, if $X = Y$ it is exactly that property: $\text{Var}(X) = E(X^2) - \mu_X^2$.

By Property 5, the formula in Property 6 reduces to the earlier formula $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$ when X and Y are independent.

We give the proofs below. However, understanding and using these properties is more important than memorizing their proofs.

2.2 Sums and integrals for computing covariance

Since covariance is defined as an expected value we compute it in the usual way as a sum or integral.

Discrete case: If X and Y have joint pmf $p(x_i, y_j)$ then

$$\text{Cov}(X, Y) = \sum_{i=1}^n \sum_{j=1}^m p(x_i, y_j)(x_i - \mu_X)(y_j - \mu_Y) = \left(\sum_{i=1}^n \sum_{j=1}^m p(x_i, y_j)x_i y_j \right) - \mu_X \mu_Y.$$

Continuous case: If X and Y have joint pdf $f(x, y)$ over range $[a, b] \times [c, d]$ then

$$\text{Cov}(X, Y) = \int_c^d \int_a^b (x - \mu_x)(y - \mu_y)f(x, y) dx dy = \left(\int_c^d \int_a^b xy f(x, y) dx dy \right) - \mu_x \mu_y.$$

2.3 Examples

Example 1. Flip a fair coin 3 times. Let X be the number of heads in the first 2 flips and let Y be the number of heads on the last 2 flips (so there is overlap on the middle flip). Compute $\text{Cov}(X, Y)$.

answer: We'll do this twice, first using the joint probability table and the definition of covariance, and then using the properties of covariance.

With 3 tosses there are only 8 outcomes $\{\text{HHH}, \text{HHT}, \dots\}$, so we can create the joint probability table directly.

$X \setminus Y$	0	1	2	$p(x_i)$
0	1/8	1/8	0	1/4
1	1/8	2/8	1/8	1/2
2	0	1/8	1/8	1/4
$p(y_j)$	1/4	1/2	1/4	1

From the marginals we compute $E(X) = 1 = E(Y)$. Now we use [use the definition](#):

$$\text{Cov}(X, Y) = E((X - \mu_x)(Y - \mu_Y)) = \sum_{i,j} p(x_i, y_j)(x_i - 1)(y_j - 1)$$

We write out the sum leaving out all the terms that are 0, i.e. all the terms where $x_i = 1$ or $y_j = 1$ or the probability is 0.

$$\text{Cov}(X, Y) = \frac{1}{8}(0-1)(0-1) + \frac{1}{8}(2-1)(2-1) = \frac{1}{4}.$$

We could also have used property 4 to do the computation: From the full table we compute

$$E(XY) = 1 \cdot \frac{2}{8} + 2 \cdot \frac{1}{8} + 2 \cdot \frac{1}{8} + 4 \cdot \frac{1}{8} = \frac{5}{4}.$$

$$\text{So } \text{Cov}(XY) = E(XY) - \mu_X \mu_Y = \frac{5}{4} - 1 = \frac{1}{4}.$$

Next we redo the computation of $\text{Cov}(X, Y)$ using the properties of covariance. As usual, let X_i be the result of the i^{th} flip, so $X_i \sim \text{Bernoulli}(0.5)$. We have

$$X = X_1 + X_2 \quad \text{and} \quad Y = X_2 + X_3.$$

We know $E(X_i) = 1/2$ and $\text{Var}(X_i) = 1/4$. Therefore using Property 2 of covariance, we have

$$\text{Cov}(X, Y) = \text{Cov}(X_1 + X_2, X_2 + X_3) = \text{Cov}(X_1, X_2) + \text{Cov}(X_1, X_3) + \text{Cov}(X_2, X_2) + \text{Cov}(X_2, X_3).$$

Since the different tosses are independent we know

$$\text{Cov}(X_1, X_2) = \text{Cov}(X_1, X_3) = \text{Cov}(X_2, X_3) = 0.$$

Looking at the expression for $\text{Cov}(X, Y)$ there is only one non-zero term

$$\text{Cov}(X, Y) = \text{Cov}(X_2, X_2) = \text{Var}(X_2) = \boxed{\frac{1}{4}}.$$

Example 2. (Zero covariance does not imply independence.) Let X be a random variable that takes values $-2, -1, 0, 1, 2$; each with probability $1/5$. Let $Y = X^2$. Show that $\text{Cov}(X, Y) = 0$ but X and Y are not independent.

answer: We make a joint probability table:

$Y \setminus X$	-2	-1	0	1	2	$p(y_j)$
0	0	0	$1/5$	0	0	$1/5$
1	0	$1/5$	0	$1/5$	0	$2/5$
4	$1/5$	0	0	0	$1/5$	$2/5$
$p(x_i)$	$1/5$	$1/5$	$1/5$	$1/5$	$1/5$	1

Using the marginals we compute means $E(X) = 0$ and $E(Y) = 2$.

Next we show that X and Y are not independent. To do this all we have to do is find one place where the product rule fails, i.e. where $p(x_i, y_j) \neq p(x_i)p(y_j)$:

$$P(X = -2, Y = 0) = 0 \quad \text{but} \quad P(X = -2) \cdot P(Y = 0) = 1/25.$$

Since these are not equal X and Y are not independent. Finally we compute covariance using Property 4:

$$\text{Cov}(X, Y) = \frac{1}{5}(-8 - 1 + 1 + 8) - \mu_X \mu_Y = 0.$$

Discussion: This example shows that $\text{Cov}(X, Y) = 0$ does not imply that X and Y are independent. In fact, X and X^2 are as dependent as random variables can be: if you know the value of X then you know the value of X^2 with 100% certainty.

The key point is that $\text{Cov}(X, Y)$ measures the linear relationship between X and Y . In the above example X and X^2 have a quadratic relationship that is completely missed by $\text{Cov}(X, Y)$.

2.4 Proofs of the properties of covariance

1 and 2 follow from similar properties for expected value.

3. This is the definition of variance:

$$\text{Cov}(X, X) = E((X - \mu_X)(X - \mu_X)) = E((X - \mu_X)^2) = \text{Var}(X).$$

4. Recall that $E(X - \mu_x) = 0$. So

$$\begin{aligned} \text{Cov}(X, Y) &= E((X - \mu_X)(Y - \mu_Y)) \\ &= E(XY - \mu_X Y - \mu_Y X + \mu_X \mu_Y) \\ &= E(XY) - \mu_X E(Y) - \mu_Y E(X) + \mu_X \mu_Y \\ &= E(XY) - \mu_X \mu_Y - \mu_X \mu_Y + \mu_X \mu_Y \\ &= E(XY) - \mu_X \mu_Y. \end{aligned}$$

5. Using properties 3 and 2 we get

$$\text{Var}(X+Y) = \text{Cov}(X+Y, X+Y) = \text{Cov}(X, X) + 2\text{Cov}(X, Y) + \text{Cov}(Y, Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y).$$

6. If X and Y are independent then $f(x, y) = f_X(x)f_Y(y)$. Therefore

$$\begin{aligned} \text{Cov}(X, Y) &= \int \int (x - \mu_X)(y - \mu_Y) f_X(x)f_Y(y) dx dy \\ &= \int (x - \mu_X) f_X(x) dx \int (y - \mu_Y) f_Y(y) dy \\ &= E(X - \mu_X)E(Y - \mu_Y) \\ &= 0. \end{aligned}$$

3 Correlation

The units of covariance $\text{Cov}(X, Y)$ are ‘units of X times units of Y ’. This makes it hard to compare covariances: if we change scales then the covariance changes as well. Correlation is a way to remove the scale from the covariance.

Definition: The [correlation coefficient](#) between X and Y is defined by

$$\text{Cor}(X, Y) = \rho = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}.$$

3.1 Properties of correlation

1. ρ is the covariance of the standardizations of X and Y .

2. ρ is dimensionless (it’s a ratio!).

3. $-1 \leq \rho \leq 1$. Furthermore,

$$\rho = +1 \text{ if and only if } Y = aX + b \text{ with } a > 0,$$

$$\rho = -1 \text{ if and only if } Y = aX + b \text{ with } a < 0.$$

Property 3 shows that ρ measures the [linear](#) relationship between variables. If the correlation is positive then when X is large, Y will tend to be large as well. If the correlation is negative then when X is large, Y will tend to be small.

Example 2 above shows that correlation can completely miss higher order relationships.

3.2 Proof of Property 3 of correlation

(This is for the mathematically interested.)

$$0 \leq \text{Var}\left(\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y}\right) = \text{Var}\left(\frac{X}{\sigma_X}\right) + \text{Var}\left(\frac{Y}{\sigma_Y}\right) - 2\text{Cov}\left(\frac{X}{\sigma_X}, \frac{Y}{\sigma_Y}\right) = 2 - 2\rho$$

$$\Rightarrow \rho \leq 1$$

Likewise $0 \leq \text{Var}\left(\frac{X}{\sigma_X} + \frac{Y}{\sigma_Y}\right) \Rightarrow -1 \leq \rho$.

If $\rho = 1$ then $0 = \text{Var}\left(\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y}\right) \Rightarrow \frac{X}{\sigma_X} - \frac{Y}{\sigma_Y} = c$. ■

Example. We continue Example 1. To compute the correlation we divide the covariance by the standard deviations. In Example 1 we found $\text{Cov}(X, Y) = 1/4$ and $\text{Var}(X) = 2\text{Var}(X_j) = 1/2$. So, $\sigma_X = 1/\sqrt{2}$. Likewise $\sigma_Y = 1/\sqrt{2}$. Thus

$$\text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{1/4}{1/2} = \frac{1}{2}.$$

We see a positive correlation, which means that larger X tend to go with larger Y and smaller X with smaller Y . In Example 1 this happens because toss 2 is included in both X and Y , so it contributes to the size of both.

3.3 Bivariate normal distributions

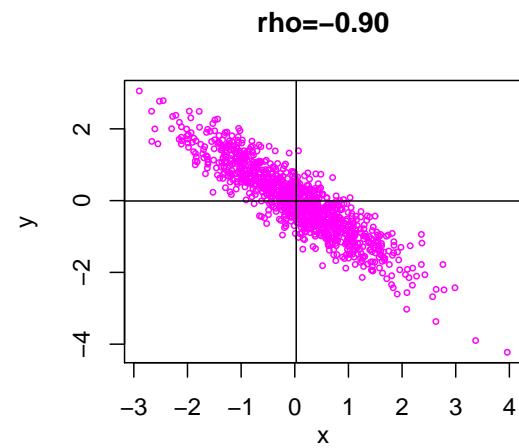
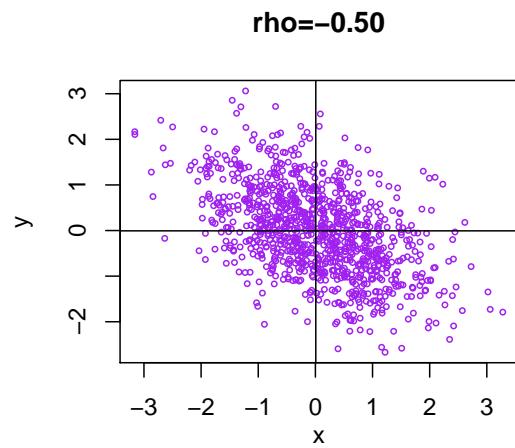
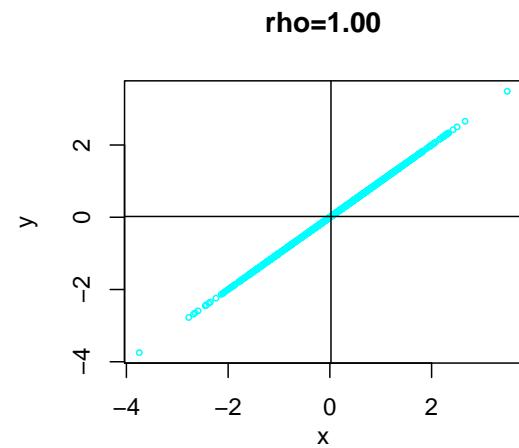
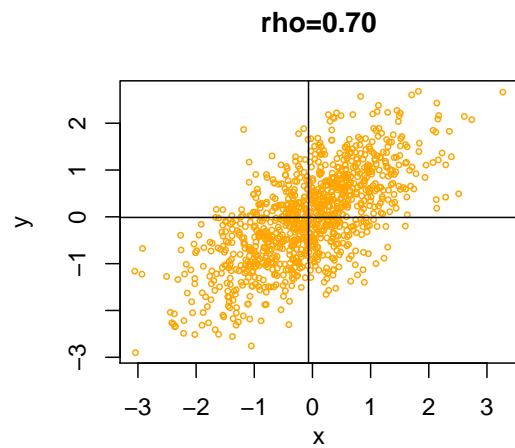
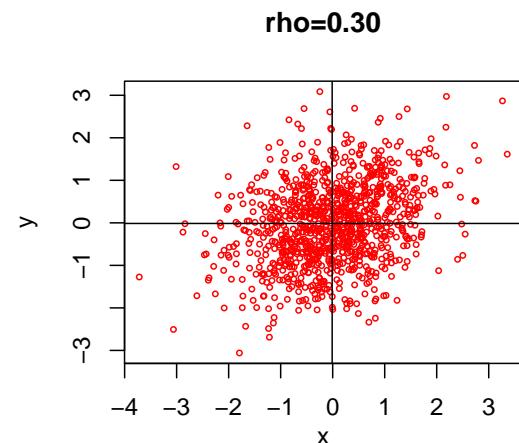
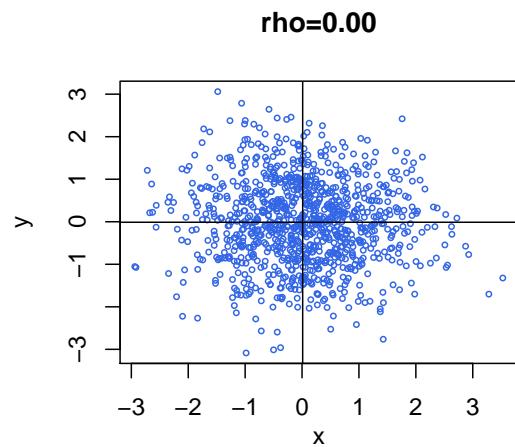
The [bivariate normal distribution](#) has density

$$f(x, y) = \frac{e^{\frac{-1}{2(1-\rho^2)} \left[\frac{(x-\mu_X)^2}{\sigma_X^2} + \frac{(y-\mu_Y)^2}{\sigma_Y^2} - \frac{2\rho(x-\mu_X)(y-\mu_Y)}{\sigma_X \sigma_Y} \right]}}{2\pi\sigma_X \sigma_Y \sqrt{1-\rho^2}}$$

For this distribution, the marginal distributions for X and Y are normal and the correlation between X and Y is ρ .

In the figures below we used R to simulate the distribution for various values of ρ . Individually X and Y are standard normal, i.e. $\mu_X = \mu_Y = 0$ and $\sigma_X = \sigma_Y = 1$. The figures show scatter plots of the results.

These plots and the next set show an important feature of correlation. We divide the data into quadrants by drawing a horizontal and a vertical line at the means of the y data and x data respectively. A positive correlation corresponds to the data tending to lie in the 1st and 3rd quadrants. A negative correlation corresponds to data tending to lie in the 2nd and 4th quadrants. You can see the data gathering about a line as ρ becomes closer to ± 1 .

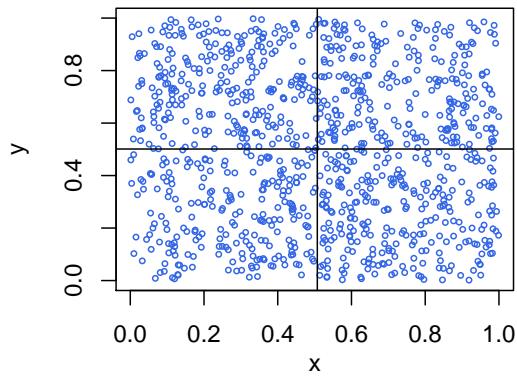


3.4 Overlapping uniform distributions

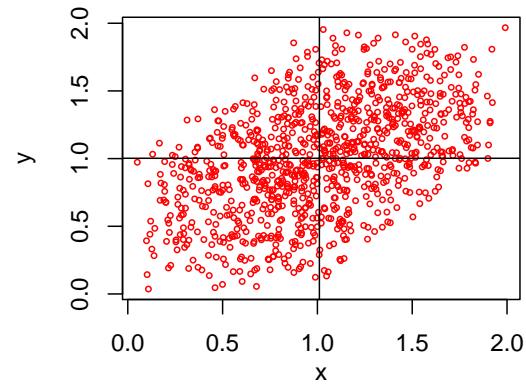
We ran simulations in R of the following scenario. X_1, X_2, \dots, X_{20} are i.i.d and follow a $U(0, 1)$ distribution. X and Y are both sums of the same number of X_i . We call the number of X_i common to both X and Y the overlap. The notation in the figures below indicates the number of X_i being summed and the number which overlap. For example, 5,3 indicates that X and Y were each the sum of 5 of the X_i and that 3 of the X_i were common to both sums. (The data was generated using `rand(1,1000)` ;)

Using the linearity of covariance it is easy to compute the theoretical correlation. For each plot we give both the theoretical correlation and the correlation of the data from the simulated sample.

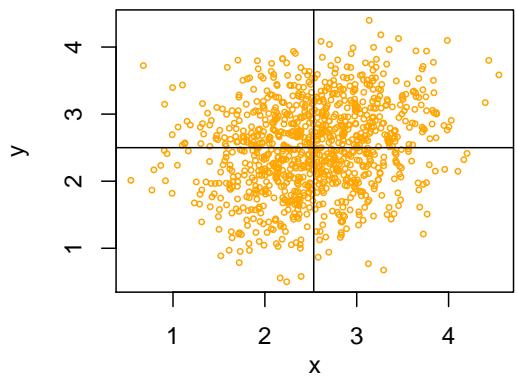
(1, 0) cor=0.00, sample_cor=-0.07



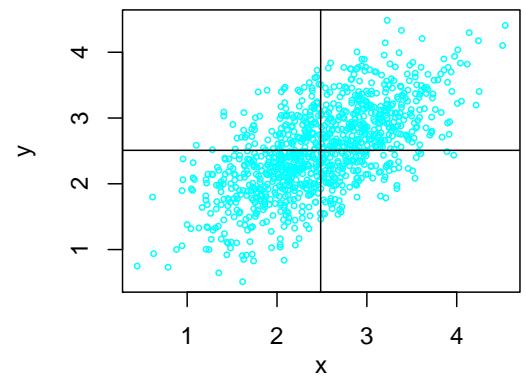
(2, 1) cor=0.50, sample_cor=0.48

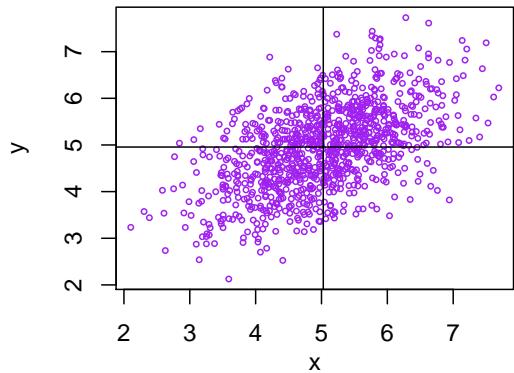
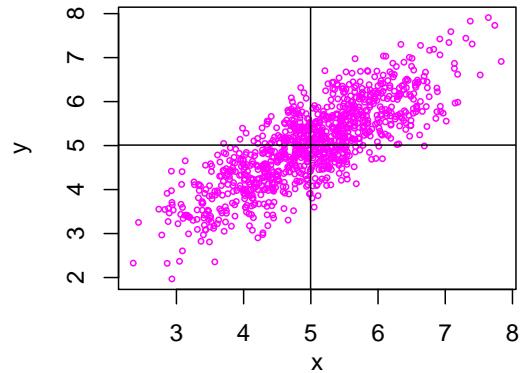


(5, 1) cor=0.20, sample_cor=0.21



(5, 3) cor=0.60, sample_cor=0.63



(10, 5) cor=0.50, sample_cor=0.53**(10, 8) cor=0.80, sample_cor=0.81**

Introduction to Statistics

Class 10, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Know the three overlapping “phases” of statistical practice.
2. Know what is meant by the term *statistic*.

2 Introduction to statistics

Statistics deals with data. Generally speaking, the goal of statistics is to make inferences based on data. We can divide this process into three phases: collecting data, describing data and analyzing data. This fits into the paradigm of the scientific method. We make hypotheses about what’s true, collect data in experiments, describe the results, and then infer from the results the **strength of the evidence** concerning our hypotheses.

2.1 Experimental design

The design of an experiment is crucial to making sure the collected data is useful. The adage ‘garbage in, garbage out’ applies here. A poorly designed experiment will produce poor quality data, from which it may be impossible to draw useful, valid inferences. To quote R.A. Fisher one of the founders of modern statistics,

To consult a statistician after an experiment is finished is often merely to ask him to conduct a post-mortem examination. He can perhaps say what the experiment died of.

2.2 Descriptive statistics

Raw data often takes the form of a massive list, array, or database of labels and numbers. To make sense of the data, we can calculate **summary statistics** like the mean, median, and interquartile range. We can also visualize the data using graphical devices like histograms, scatterplots, and the empirical cdf. These methods are useful for both communicating and exploring the data to gain insight into its structure, such as whether it might follow a familiar probability distribution.

2.3 Inferential statistics

Ultimately we want to draw inferences about the world. Often this takes the form of specifying a statistical model for the random process by which the data arises. For example, suppose the data takes the form of a series of measurements whose error we believe follows a normal distribution. (Note this is always an approximation since we know the error must

have some bound while a normal distribution has range $(-\infty, \infty)$.) We might then use the data to provide evidence for or against this hypothesis. Our focus in 18.05 will be on how to use data to draw inferences about model parameters. For example, assuming gestational length follows a $N(\mu, \sigma)$ distribution, we'll use the data of the gestational lengths of, say, 500 pregnancies to draw inferences about the values of the parameters μ and σ . Similarly, we may model the result of a two-candidate election by a $Bernoulli(p)$ distribution, and use poll data to draw inferences about the value of p .

We can rarely make definitive statements about such parameters because the data itself comes from a random process (such as choosing who to poll). Rather, our statistical evidence will always involve probability statements. Unfortunately, the media and public at large are wont to misunderstand the probabilistic meaning of statistical statements. In fact, researchers themselves often commit the same errors. In this course, we will emphasize the **meaning** of statistical statements alongside the **methods** which produce them.

Example 1. To study the effectiveness of new treatment for cancer, patients are recruited and then divided into an experimental group and a control group. The experimental group is given the new treatment and the control group receives the current standard of care. Data collected from the patients might include demographic information, medical history, initial state of cancer, progression of the cancer over time, treatment cost, and the effect of the treatment on tumor size, remission rates, longevity, and quality of life. The data will be used to make inferences about the effectiveness of the new treatment compared to the current standard of care.

Notice that this study will go through all three phases described above. The experimental design must specify the size of the study, who will be eligible to join, how the experimental and control groups will be chosen, how the treatments will be administered, whether or not the subjects or doctors know who is getting which treatment, and precisely what data will be collected, among other things. Once the data is collected it must be described and analyzed to determine whether it supports the hypothesis that the new treatment is more (or less) effective than the current one(s), and by how much. These statistical conclusions will be framed as precise statements involving probabilities.

As noted above, misinterpreting the exact meaning of statistical statements is a common source of error which has led to tragedy on more than one occasion.

Example 2. In 1999 in Great Britain, Sally Clark was convicted of murdering her two children after each child died weeks after birth (the first in 1996, the second in 1998). Her conviction was largely based on a faulty use of statistics to rule out sudden infant death syndrome. Though her conviction was overturned in 2003, she developed serious psychiatric problems during and after her imprisonment and died of alcohol poisoning in 2007. See http://en.wikipedia.org/wiki/Sally_Clark

This TED talk discusses the Sally Clark case and other instances of poor statistical intuition:
<http://www.youtube.com/watch?v=kLmzxRcUT0>

2.4 What is a statistic?

We give a simple definition whose meaning is best elucidated by examples.

Definition. A **statistic** is anything that can be computed from the collected data.

Example 3. Consider the data of 1000 rolls of a die. All of the following are statistics: the average of the 1000 rolls; the number of times a 6 was rolled; the sum of the squares of the rolls minus the number of even rolls. It's hard to imagine how we would use the last example, but it is a statistic. On the other hand, the probability of rolling a 6 is *not* a statistic, whether or not the die is truly fair. Rather this probability is a property of the die (and the way we roll it) which we can **estimate** using the data. Such an estimate is given by the statistic ‘proportion of the rolls that were 6’.

Example 4. Suppose we treat a group of cancer patients with a new procedure and collect data on how long they survive post-treatment. From the data we can compute the average survival time of patients in the group. We might employ this statistic as an estimate of the average survival time for future cancer patients following the new procedure. The actual survival is *not* a statistic.

Example 5. Suppose we ask 1000 residents whether or not they support the proposal to legalize marijuana in Massachusetts. The proportion of the 1000 who support the proposal is a statistic. The proportion of all Massachusetts residents who support the proposal is *not* a statistic since we have not queried every single one (note the word “collected” in the definition). Rather, we hope to draw a statistical conclusion about the state-wide proportion based on the data of our random sample.

The following are two general types of statistics we will use in 18.05.

1. **Point statistics:** a single value computed from data, such as the sample average \bar{x}_n or the sample standard deviation s_n .
2. **Interval statistics:** an interval $[a, b]$ computed from the data. This is really just a pair of point statistics, and will often be presented in the form $\bar{x} \pm s$.

3 Review of Bayes' theorem

We cannot stress strongly enough how important Bayes' theorem is to our view of inferential statistics. Recall that Bayes' theorem allows us to ‘invert’ conditional probabilities. That is, if H and D are events, then Bayes' theorem says

$$P(H|D) = \frac{P(D|H)P(H)}{P(D)}.$$

In scientific experiments we start with a hypothesis and collect data to test the hypothesis. We will often let H represent the event ‘our hypothesis is true’ and let D be the collected data. In these words Bayes' theorem says

$$P(\text{hypothesis is true} \mid \text{data}) = \frac{P(\text{data} \mid \text{hypothesis is true}) \cdot P(\text{hypothesis is true})}{P(\text{data})}$$

The left-hand term is the probability our hypothesis is true given the data we collected. This is precisely what we'd like to know. When all the probabilities on the right are known exactly, we can compute the probability on the left exactly. This will be our focus next week. Unfortunately, in practice we rarely know the exact values of all the terms on the

right. Statisticians have developed a number of ways to cope with this lack of knowledge and still make useful inferences. We will be exploring these methods for the rest of the course.

Example 6. Screening for a disease redux

Suppose a screening test for a disease has a 1% false positive rate and a 1% false negative rate. Suppose also that the rate of the disease in the population is 0.002. Finally suppose a randomly selected person tests positive. In the language of hypothesis and data we have:

Hypothesis: H = ‘the person has the disease’

Data: D = ‘the test was positive.’

What we want to know: $P(H|D) = P(\text{the person has the disease} \mid \text{a positive test})$

In this example all the probabilities on the right are known so we can use Bayes’ theorem to compute what we want to know.

$$\begin{aligned} P(\text{hypothesis} \mid \text{data}) &= P(\text{the person has the disease} \mid \text{a positive test}) \\ &= P(H|D) \\ &= \frac{P(D|H)P(H)}{P(D)} \\ &= \frac{.99 \cdot .002}{.99 \cdot .002 + .01 \cdot .998} \\ &= 0.166 \end{aligned}$$

Before the test we would have said the probability the person had the disease was 0.002. After the test we see the probability is 0.166. That is, the positive test provides some evidence that the person has the disease.

Maximum Likelihood Estimates

Class 10, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to define the likelihood function for a parametric model given data.
2. Be able to compute the maximum likelihood estimate of unknown parameter(s).

2 Introduction

Suppose we know we have data consisting of values x_1, \dots, x_n drawn from an exponential distribution. The question remains: which exponential distribution?!

We have casually referred to *the* exponential distribution or *the* binomial distribution or *the* normal distribution. In fact the exponential distribution $\exp(\lambda)$ is not a single distribution but rather a one-parameter family of distributions. Each value of λ defines a different distribution in the family, with pdf $f_\lambda(x) = \lambda e^{-\lambda x}$ on $[0, \infty)$. Similarly, a binomial distribution $\text{bin}(n, p)$ is determined by the two parameters n and p , and a normal distribution $N(\mu, \sigma^2)$ is determined by the two parameters μ and σ^2 (or equivalently, μ and σ). Parameterized families of distributions are often called [parametric distributions](#) or [parametric models](#).

We are often faced with the situation of having random data which we know (or believe) is drawn from a parametric model, whose parameters we do not know. For example, in an election between two candidates, polling data constitutes draws from a $\text{Bernoulli}(p)$ distribution with unknown parameter p . In this case we would like to use the data to estimate the value of the parameter p , as the latter predicts the result of the election. Similarly, assuming gestational length follows a normal distribution, we would like to use the data of the gestational lengths from a random sample of pregnancies to draw inferences about the values of the parameters μ and σ^2 .

Our focus so far has been on computing the [probability of data](#) arising from a parametric model with [known parameters](#). Statistical inference flips this on its head: we will estimate the [probability of parameters](#) given a parametric model and [observed data](#) drawn from it. In the coming weeks we will see how parameter values are naturally viewed as hypotheses, so we are in fact estimating the probability of various hypotheses given the data.

3 Maximum Likelihood Estimates

There are many methods for estimating unknown parameters from data. We will first consider the [maximum likelihood estimate](#) (MLE), which answers the question:

[For which parameter value does the observed data have the biggest probability?](#)

The MLE is an example of a [point estimate](#) because it gives a single value for the unknown parameter (later our estimates will involve intervals and probabilities). Two advantages of

the MLE are that it is often easy to compute and that it agrees with our intuition in simple examples. We will explain the MLE through a series of examples.

Example 1. A coin is flipped 100 times. Given that there were 55 heads, find the maximum likelihood estimate for the probability p of heads on a single toss.

Before actually solving the problem, let's establish some notation and terms.

We can think of counting the number of heads in 100 tosses as an experiment. For a given value of p , the probability of getting 55 heads in this experiment is the binomial probability

$$P(55 \text{ heads}) = \binom{100}{55} p^{55} (1-p)^{45}.$$

The probability of getting 55 heads depends on the value of p , so let's include p in by using the notation of conditional probability:

$$P(55 \text{ heads} | p) = \binom{100}{55} p^{55} (1-p)^{45}.$$

You should read $P(55 \text{ heads} | p)$ as:

‘the probability of 55 heads given p ,’

or more precisely as

‘the probability of 55 heads given that the probability of heads on a single toss is p .’

Here are some standard terms we will use as we do statistics.

- **Experiment:** Flip the coin 100 times and count the number of heads.
- **Data:** The data is the result of the experiment. In this case it is ‘55 heads’.
- **Parameter(s) of interest:** We are interested in the value of the unknown parameter p .
- **Likelihood, or likelihood function:** this is $P(\text{data} | p)$. Note it is a function of both the data and the parameter p . In this case the likelihood is

$$P(55 \text{ heads} | p) = \binom{100}{55} p^{55} (1-p)^{45}.$$

Notes: **1.** The likelihood $P(\text{data} | p)$ changes as the parameter of interest p changes.

2. Look carefully at the definition. One typical source of confusion is to mistake the likelihood $P(\text{data} | p)$ for $P(p | \text{data})$. We know from our earlier work with Bayes' theorem that $P(\text{data} | p)$ and $P(p | \text{data})$ are usually very different.

Definition: Given data the **maximum likelihood estimate (MLE)** for the parameter p is the value of p that maximizes the likelihood $P(\text{data} | p)$. That is, the MLE is the value of p for which the data is most likely.

answer: For the problem at hand, we saw above that the likelihood

$$P(55 \text{ heads} | p) = \binom{100}{55} p^{55} (1-p)^{45}.$$

We'll use the notation \hat{p} for the MLE. We use calculus to find it by taking the derivative of the likelihood function and setting it to 0.

$$\frac{d}{dp} P(\text{data} | p) = \binom{100}{55} (55p^{54}(1-p)^{45} - 45p^{55}(1-p)^{44}) = 0.$$

Solving this for p we get

$$\begin{aligned} 55p^{54}(1-p)^{45} &= 45p^{55}(1-p)^{44} \\ 55(1-p) &= 45p \\ 55 &= 100p \\ \text{the MLE is } \hat{p} &= .55 \end{aligned}$$

- Note:
- 1.** The MLE for p turned out to be exactly the fraction of heads we saw in our data.
 - 2.** The MLE is computed from the data. That is, it is a statistic.
 - 3.** Officially you should check that the critical point is indeed a maximum. You can do this with the second derivative test.

3.1 Log likelihood

If is often easier to work with the natural log of the likelihood function. For short this is simply called the [log likelihood](#). Since $\ln(x)$ is an increasing function, the maxima of the likelihood and log likelihood coincide.

Example 2. Redo the previous example using log likelihood.

answer: We had the likelihood $P(55 \text{ heads} | p) = \binom{100}{55} p^{55}(1-p)^{45}$. Therefore the log likelihood is

$$\ln(P(55 \text{ heads} | p)) = \ln\left(\binom{100}{55}\right) + 55 \ln(p) + 45 \ln(1-p).$$

Maximizing likelihood is the same as maximizing log likelihood. We check that calculus gives us the same answer as before:

$$\begin{aligned} \frac{d}{dp} (\text{log likelihood}) &= \frac{d}{dp} \left[\ln\left(\binom{100}{55}\right) + 55 \ln(p) + 45 \ln(1-p) \right] \\ &= \frac{55}{p} - \frac{45}{1-p} = 0 \\ \Rightarrow 55(1-p) &= 45p \\ \Rightarrow \hat{p} &= .55 \end{aligned}$$

3.2 Maximum likelihood for continuous distributions

For continuous distributions, we use the probability density function to define the likelihood. We show this in a few examples. In the next section we explain how this is analogous to what we did in the discrete case.

Example 3. Light bulbs

Suppose that the lifetime of *Badger* brand light bulbs is modeled by an exponential distribution with (unknown) parameter λ . We test 5 bulbs and find they have lifetimes of 2, 3, 1, 3, and 4 years, respectively. What is the MLE for λ ?

answer: We need to be careful with our notation. With five different values it is best to use subscripts. Let X_j be the lifetime of the i^{th} bulb and let x_i be the value X_i takes. Then each X_i has pdf $f_{X_i}(x_i) = \lambda e^{-\lambda x_i}$. We assume the lifetimes of the bulbs are independent, so the joint pdf is the product of the individual densities:

$$f(x_1, x_2, x_3, x_4, x_5 | \lambda) = (\lambda e^{-\lambda x_1})(\lambda e^{-\lambda x_2})(\lambda e^{-\lambda x_3})(\lambda e^{-\lambda x_4})(\lambda e^{-\lambda x_5}) = \lambda^5 e^{-\lambda(x_1+x_2+x_3+x_4+x_5)}.$$

Note that we write this as a conditional density, since it depends on λ . Viewing the data as fixed and λ as variable, this density is the likelihood function. Our data had values

$$x_1 = 2, x_2 = 3, x_3 = 1, x_4 = 3, x_5 = 4.$$

So the likelihood and log likelihood functions with this data are

$$f(2, 3, 1, 3, 4 | \lambda) = \lambda^5 e^{-13\lambda}, \quad \ln(f(2, 3, 1, 3, 4 | \lambda)) = 5 \ln(\lambda) - 13\lambda$$

Finally we use calculus to find the MLE:

$$\frac{d}{d\lambda}(\text{log likelihood}) = \frac{5}{\lambda} - 13 = 0 \Rightarrow \boxed{\hat{\lambda} = \frac{5}{13}}.$$

Note: **1.** In this example we used an uppercase letter for a random variable and the corresponding lowercase letter for the value it takes. This will be our usual practice.

2. The MLE for λ turned out to be the reciprocal of the sample mean \bar{x} , so $X \sim \exp(\hat{\lambda})$ satisfies $E(X) = \bar{x}$.

The following example illustrates how we can use the method of maximum likelihood to estimate multiple parameters at once.

Example 4. Normal distributions

Suppose the data x_1, x_2, \dots, x_n is drawn from a $N(\mu, \sigma^2)$ distribution, where μ and σ are unknown. Find the maximum likelihood estimate for the pair (μ, σ^2) .

answer: Let's be precise and phrase this in terms of random variables and densities. Let uppercase X_1, \dots, X_n be i.i.d. $N(\mu, \sigma^2)$ random variables, and let lowercase x_i be the value X_i takes. The density for each X_i is

$$f_{X_i}(x_i) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}.$$

Since the X_i are independent their joint pdf is the product of the individual pdf's:

$$f(x_1, \dots, x_n | \mu, \sigma) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n e^{-\sum_{i=1}^n \frac{(x_i-\mu)^2}{2\sigma^2}}.$$

For the fixed data x_1, \dots, x_n , the likelihood and log likelihood are

$$f(x_1, \dots, x_n | \mu, \sigma) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n e^{-\sum_{i=1}^n \frac{(x_i-\mu)^2}{2\sigma^2}}, \quad \ln(f(x_1, \dots, x_n | \mu, \sigma)) = -n \ln(\sqrt{2\pi}) - n \ln(\sigma) - \sum_{i=1}^n \frac{(x_i-\mu)^2}{2\sigma^2}.$$

Since $\ln(f(x_1, \dots, x_n | \mu, \sigma))$ is a function of the two variables μ, σ we use partial derivatives to find the MLE. The easy value to find is $\hat{\mu}$:

$$\frac{\partial f(x_1, \dots, x_n | \mu, \sigma)}{\partial \mu} = \sum_{i=1}^n \frac{(x_i - \mu)}{\sigma^2} = 0 \Rightarrow \sum_{i=1}^n x_i = n\mu \Rightarrow \hat{\mu} = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}.$$

To find $\hat{\sigma}$ we differentiate and solve for σ :

$$\frac{\partial f(x_1, \dots, x_n | \mu, \sigma)}{\partial \sigma} = -\frac{n}{\sigma} + \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^3} = 0 \Rightarrow \hat{\sigma}^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}.$$

We already know $\hat{\mu} = \bar{x}$, so we use that as the value for μ in the formula for $\hat{\sigma}$. We get the maximum likelihood estimates

$$\begin{aligned} \hat{\mu} &= \bar{x} &&= \text{the mean of the data} \\ \hat{\sigma}^2 &= \sum_{i=1}^n \frac{1}{n} (x_i - \hat{\mu})^2 = \sum_{i=1}^n \frac{1}{n} (x_i - \bar{x})^2 &&= \text{the variance of the data.} \end{aligned}$$

Example 5. Uniform distributions

Suppose our data x_1, \dots, x_n are independently drawn from a uniform distribution $U(a, b)$. Find the MLE estimate for a and b .

answer: This example is different from the previous ones in that we won't use calculus to find the MLE. The density for $U(a, b)$ is $\frac{1}{b-a}$ on $[a, b]$. Therefore our likelihood function is

$$f(x_1, \dots, x_n | a, b) = \begin{cases} \left(\frac{1}{b-a}\right)^n & \text{if all } x_i \text{ are in the interval } [a, b] \\ 0 & \text{otherwise.} \end{cases}$$

This is maximized by making $b - a$ as small as possible. The only restriction is that the interval $[a, b]$ must include all the data. Thus the MLE for the pair (a, b) is

$$\hat{a} = \min(x_1, \dots, x_n) \quad \hat{b} = \max(x_1, \dots, x_n).$$

Example 6. Capture/recapture method

The capture/recapture method is a way to estimate the size of a population in the wild. The method assumes that each animal in the population is equally likely to be captured by a trap.

Suppose 10 animals are captured, tagged and released. A few months later, 20 animals are captured, examined, and released. 4 of these 20 are found to be tagged. Estimate the size of the wild population using the MLE for the probability that a wild animal is tagged.

answer: Our unknown parameter n is the number of animals in the wild. Our data is that 4 out of 20 recaptured animals were tagged (and that there are 10 tagged animals). The likelihood function is

$$P(\text{data} | n \text{ animals}) = \frac{\binom{n-10}{16} \binom{10}{4}}{\binom{n}{20}}$$

(The numerator is the number of ways to choose 16 animals from among the $n-10$ untagged ones times the number of ways to choose 4 out of the 10 tagged animals. The denominator

is the number of ways to choose 20 animals from the entire population of n .) We can use R to compute that the likelihood function is maximized when $n = 50$. This should make some sense. It says our best estimate is that the fraction of all animals that are tagged is 10/50 which equals the fraction of recaptured animals which are tagged.

Example 7. Hardy-Weinberg. Suppose that a particular gene occurs as one of two alleles (A and a), where allele A has frequency θ in the population. That is, a random copy of the gene is A with probability θ and a with probability $1 - \theta$. Since a diploid genotype consists of two genes, the probability of each genotype is given by:

genotype	AA	Aa	aa
probability	θ^2	$2\theta(1 - \theta)$	$(1 - \theta)^2$

Suppose we test a random sample of people and find that k_1 are AA , k_2 are Aa , and k_3 are aa . Find the MLE of θ .

answer: The likelihood function is given by

$$P(k_1, k_2, k_3 | \theta) = \binom{k_1 + k_2 + k_3}{k_1} \binom{k_2 + k_3}{k_2} \binom{k_3}{k_3} \theta^{2k_1} (2\theta(1 - \theta))^{k_2} (1 - \theta)^{2k_3}.$$

So the log likelihood is given by

$$\text{constant} + 2k_1 \ln(\theta) + k_2 \ln(2\theta(1 - \theta)) + k_3 \ln(1 - \theta)$$

We set the derivative equal to zero:

$$\frac{2k_1 + k_2}{\theta} - \frac{k_2 + 2k_3}{1 - \theta} = 0$$

Solving for θ , we find the MLE is

$$\hat{\theta} = \frac{2k_1 + k_2}{2k_1 + 2k_2 + 2k_3},$$

which is simply the fraction of A alleles among all the genes in the sampled population.

4 Why we use the density to find the MLE for continuous distributions

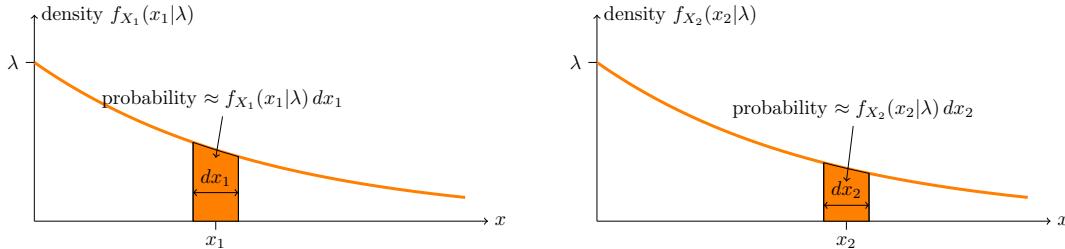
The idea for the maximum likelihood estimate is to find the value of the parameter(s) for which the data has the highest probability. In this section we'll see that we're doing this is really what we are doing with the densities. We will do this by considering a smaller version of the light bulb example.

Example 8. Suppose we have two light bulbs whose lifetimes follow an exponential(λ) distribution. Suppose also that we independently measure their lifetimes and get data $x_1 = 2$ years and $x_2 = 3$ years. Find the value of λ that maximizes the probability of this data.

answer: The main paradox to deal with is that for a continuous distribution the probability of a single value, say $x_1 = 2$, is zero. We resolve this paradox by remembering that a single

measurement really means a range of values, e.g. in this example we might check the light bulb once a day. So the data $x_1 = 2$ years really means x_1 is somewhere in a range of 1 day around 2 years.

If the range is small we call it dx_1 . The probability that X_1 is in the range is approximated by $f_{X_1}(x_1|\lambda) dx_1$. This is illustrated in the figure below. The data value x_2 is treated in exactly the same way.



The usual relationship between density and probability for small ranges.

Since the data is collected independently the joint probability is the product of the individual probabilities. Stated carefully

$$P(X_1 \text{ in range}, X_2 \text{ in range}|\lambda) \approx f_{X_1}(x_1|\lambda) dx_1 \cdot f_{X_2}(x_2|\lambda) dx_2$$

Finally, using the values $x_1 = 2$ and $x_2 = 3$ and the formula for an exponential pdf we have

$$P(X_1 \text{ in range}, X_2 \text{ in range}|\lambda) \approx \lambda e^{-2\lambda} dx_1 \cdot \lambda e^{-3\lambda} dx_2 = \lambda^2 e^{-5\lambda} dx_1 dx_2.$$

Now that we have a genuine probability we can look for the value of λ that maximizes it. Looking at the formula above we see that the factor $dx_1 dx_2$ will play no role in finding the maximum. So for the MLE we drop it and simply call the density the likelihood:

$$\text{likelihood} = f(x_1, x_2|\lambda) = \lambda^2 e^{-5\lambda}.$$

The value of λ that maximizes this is found just like in the example above. It is $\hat{\lambda} = 2/5$.

5 Appendix: Properties of the MLE

For the interested reader, we note several nice features of the MLE. These are quite technical and will not be on any exams.

The MLE behaves well under transformations. That is, if \hat{p} is the MLE for p and g is a one-to-one function, then $g(\hat{p})$ is the MLE for $g(p)$. For example, if $\hat{\sigma}$ is the MLE for the standard deviation σ then $(\hat{\sigma})^2$ is the MLE for the variance σ^2 .

Furthermore, the MLE is **asymptotically unbiased** and has **asymptotically minimal variance**. To explain these notions, note that the MLE is itself a random variable since the data is random and the MLE is computed from the data. Let x_1, x_2, \dots be an infinite sequence of samples from a distribution with parameter p . Let \hat{p}_n be the MLE for p based on the data x_1, \dots, x_n .

Asymptotically unbiased means that as the amount of data grows, the mean of the MLE converges to p . In symbols: $E(\hat{p}_n) \rightarrow p$ as $n \rightarrow \infty$. Of course, we would like the MLE to be

close to p with high probability, not just on average, so the smaller the variance of the MLE the better. Asymptotically minimal variance means that as the amount of data grows, the MLE has the minimal variance among all unbiased estimators of p . In symbols: for any unbiased estimator \tilde{p}_n and $\epsilon > 0$ we have that $\text{Var}(\tilde{p}_n) + \epsilon > \text{Var}(\hat{p}_n)$ as $n \rightarrow \infty$.

Bayesian Updating with Discrete Priors

Class 11, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to apply Bayes' theorem to compute probabilities.
2. Be able to define the and to identify the roles of prior probability, likelihood (Bayes term), posterior probability, data and hypothesis in the application of Bayes' Theorem.
3. Be able to use a Bayesian update table to compute posterior probabilities.

2 Review of Bayes' theorem

Recall that Bayes' theorem allows us to ‘invert’ conditional probabilities. If \mathcal{H} and \mathcal{D} are events, then:

$$P(\mathcal{H} | \mathcal{D}) = \frac{P(\mathcal{D} | \mathcal{H})P(\mathcal{H})}{P(\mathcal{D})}$$

Our view is that Bayes' theorem forms the foundation for inferential statistics. We will begin to justify this view today.

2.1 The base rate fallacy

When we first learned Bayes' theorem we worked an example about screening tests showing that $P(\mathcal{D}|\mathcal{H})$ can be very different from $P(\mathcal{H}|\mathcal{D})$. In the appendix we work a similar example. If you are not comfortable with Bayes' theorem you should read the example in the appendix now.

3 Terminology and Bayes' theorem in tabular form

We now use a coin tossing problem to introduce terminology and a tabular format for Bayes' theorem. This will provide a simple, uncluttered example that shows our main points.

Example 1. There are three types of coins which have different probabilities of landing heads when tossed.

- Type A coins are fair, with probability 0.5 of heads
- Type B coins are bent and have probability 0.6 of heads
- Type C coins are bent and have probability 0.9 of heads

Suppose I have a drawer containing 5 coins: 2 of type A , 2 of type B , and 1 of type C . I reach into the drawer and pick a coin at random. Without showing you the coin I flip it once and get heads. What is the probability it is type A ? Type B ? Type C ?

answer: Let A , B , and C be the event that the chosen coin was type A , type B , and type C . Let \mathcal{D} be the event that the toss is heads. The problem asks us to find

$$P(A|\mathcal{D}), \quad P(B|\mathcal{D}), \quad P(C|\mathcal{D}).$$

Before applying Bayes' theorem, let's introduce some terminology.

- **Experiment:** pick a coin from the drawer at random, flip it, and record the result.
- **Data:** the result of our experiment. In this case the event \mathcal{D} = ‘heads’. We think of \mathcal{D} as data that provides evidence for or against each hypothesis.
- **Hypotheses:** we are testing three hypotheses: the coin is type A , B or C .
- **Prior probability:** the probability of each hypothesis prior to tossing the coin (collecting data). Since the drawer has 2 coins of type A , 2 of type B and 1 of type C we have

$$P(A) = 0.4, \quad P(B) = 0.4, \quad P(C) = 0.2.$$

- **Likelihood:** (This is the same likelihood we used for the MLE.) The likelihood function is $P(\mathcal{D}|\mathcal{H})$, i.e., the probability of the data assuming that the hypothesis is true. Most often we will consider the data as fixed and let the hypothesis vary. For example, $P(\mathcal{D}|A)$ = probability of heads if the coin is type A . In our case the likelihoods are

$$P(\mathcal{D}|A) = 0.5, \quad P(\mathcal{D}|B) = 0.6, \quad P(\mathcal{D}|C) = 0.9.$$

The name likelihood is so well established in the literature that we have to teach it to you. However in colloquial language likelihood and probability are synonyms. This leads to the likelihood function often being confused with the probability of a hypothesis. Because of this we'd prefer to use the name Bayes' term. However since we are stuck with ‘likelihood’ we will try to use it very carefully and in a way that minimizes any confusion.

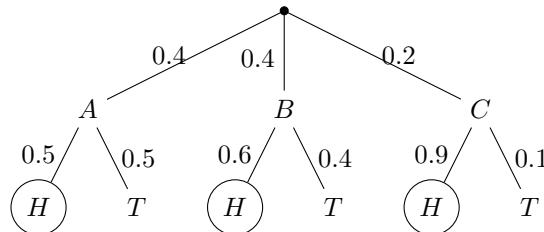
- **Posterior probability:** the probability (posterior to) of each hypothesis given the data from tossing the coin.

$$P(A|\mathcal{D}), \quad P(B|\mathcal{D}), \quad P(C|\mathcal{D}).$$

These posterior probabilities are what the problem asks us to find.

We now use Bayes' theorem to compute each of the posterior probabilities. We are going to write this out in complete detail so we can pick out each of the parts (Remember that the data \mathcal{D} is that the toss was heads.)

First we organize the probabilities into a tree:



Probability tree for choosing and tossing a coin.

Bayes' theorem says, e.g. $P(A|\mathcal{D}) = \frac{P(\mathcal{D}|A)P(A)}{P(\mathcal{D})}$. The denominator $P(\mathcal{D})$ is computed using the law of total probability:

$$P(\mathcal{D}) = P(\mathcal{D}|A)P(A) + P(\mathcal{D}|B)P(B) + P(\mathcal{D}|C)P(C) = 0.5 \cdot 0.4 + 0.6 \cdot 0.4 + 0.9 \cdot 0.2 = 0.62.$$

Now each of the three posterior probabilities can be computed:

$$\begin{aligned} P(A|\mathcal{D}) &= \frac{P(\mathcal{D}|A)P(A)}{P(\mathcal{D})} = \frac{0.5 \cdot 0.4}{0.62} = \frac{0.2}{0.62} \\ P(B|\mathcal{D}) &= \frac{P(\mathcal{D}|B)P(B)}{P(\mathcal{D})} = \frac{0.6 \cdot 0.4}{0.62} = \frac{0.24}{0.62} \\ P(C|\mathcal{D}) &= \frac{P(\mathcal{D}|C)P(C)}{P(\mathcal{D})} = \frac{0.9 \cdot 0.2}{0.62} = \frac{0.18}{0.62} \end{aligned}$$

Notice that the total probability $P(\mathcal{D})$ is the same in each of the denominators and that it is the sum of the three numerators. We can organize all of this very neatly in a [Bayesian update table](#):

Bayes				
hypothesis	prior	likelihood	numerator	posterior
\mathcal{H}	$P(\mathcal{H})$	$P(\mathcal{D} \mathcal{H})$	$P(\mathcal{D} \mathcal{H})P(\mathcal{H})$	$P(\mathcal{H} \mathcal{D})$
A	0.4	0.5	0.2	0.3226
B	0.4	0.6	0.24	0.3871
C	0.2	0.9	0.18	0.2903
total	1		0.62	1

The [Bayes numerator](#) is the product of the prior and the likelihood. We see in each of the Bayes' formula computations above that the posterior probability is obtained by dividing the Bayes numerator by $P(\mathcal{D}) = 0.625$. We also see that the law of total probability says that $P(\mathcal{D})$ is the sum of the entries in the Bayes numerator column.

Bayesian updating: The process of going from the prior probability $P(\mathcal{H})$ to the posterior $P(\mathcal{H}|\mathcal{D})$ is called [Bayesian updating](#). Bayesian updating uses the data to alter our understanding of the probability of each of the possible hypotheses.

3.1 Important things to notice

- There are two types of probabilities: Type one is the standard probability of data, e.g. the probability of heads is $p = 0.9$. Type two is the probability of the hypotheses, e.g. the probability the chosen coin is type A , B or C . This second type has prior (before the data) and posterior (after the data) values.
- The posterior (after the data) probabilities for each hypothesis are in the last column. We see that coin B is now the most probable, though its probability has decreased from a prior probability of 0.4 to a posterior probability of 0.39. Meanwhile, the probability of type C has increased from 0.2 to 0.29.
- The Bayes numerator column determines the posterior probability column. To compute the latter, we simply rescaled the Bayes numerator so that it sums to 1.

4. If all we care about is finding the most likely hypothesis, the Bayes numerator works as well as the normalized posterior.
5. The likelihood column does not sum to 1. The likelihood function is *not* a probability function.
6. The posterior probability represents the outcome of a ‘tug-of-war’ between the likelihood and the prior. When calculating the posterior, a large prior may be deflated by a small likelihood, and a small prior may be inflated by a large likelihood.
7. The maximum likelihood estimate (MLE) for Example 1 is hypothesis C , with a likelihood $P(\mathcal{D}|C) = 0.9$. The MLE is useful, but you can see in this example that it is not the entire story, since type B has the greatest posterior probability.

Terminology in hand, we can express Bayes’ theorem in various ways:

$$P(\mathcal{H}|\mathcal{D}) = \frac{P(\mathcal{D}|\mathcal{H})P(\mathcal{H})}{P(\mathcal{D})}$$

$$P(\text{hypothesis}|\text{data}) = \frac{P(\text{data}|\text{hypothesis})P(\text{hypothesis})}{P(\text{data})}$$

With the data fixed, the denominator $P(\mathcal{D})$ just serves to normalize the total posterior probability to 1. So we can also express Bayes’ theorem as a statement about the proportionality of two functions of \mathcal{H} (i.e., of the last two columns of the table).

$$P(\text{hypothesis}|\text{data}) \propto P(\text{data}|\text{hypothesis})P(\text{hypothesis})$$

This leads to the most elegant form of Bayes’ theorem in the context of Bayesian updating:

posterior \propto likelihood \times prior

3.2 Prior and posterior probability mass functions

Earlier in the course we saw that it is convenient to use random variables and probability mass functions. To do this we had to assign values to events (head is 1 and tails is 0). We will do the same thing in the context of Bayesian updating.

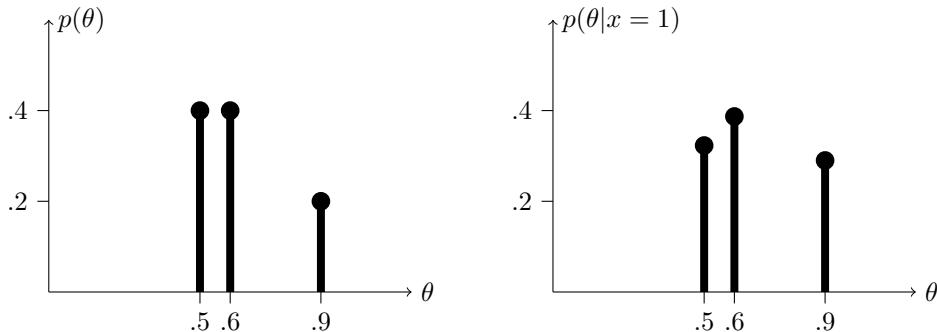
Our standard notations will be:

- θ is the **value of the hypothesis**.
- $p(\theta)$ is the **prior probability mass function of the hypothesis**.
- $p(\theta|\mathcal{D})$ is the **posterior probability mass function of the hypothesis given the data**.
- $p(\mathcal{D}|\theta)$ is the **likelihood function**. (This is not a pmf!)

In Example 1 we can represent the three hypotheses A , B , and C by $\theta = 0.5, 0.6, 0.9$. For the data we’ll let $x = 1$ mean heads and $x = 0$ mean tails. Then the prior and posterior probabilities in the table define the prior and posterior probability mass functions.

Hypothesis	θ	prior pmf $p(\theta)$	poster pmf $p(\theta x=1)$
A	0.5	$P(A) = p(0.5) = 0.4$	$P(A \mathcal{D}) = p(0.5 x=1) = 0.3226$
B	0.6	$P(B) = p(0.6) = 0.4$	$P(B \mathcal{D}) = p(0.6 x=1) = 0.3871$
C	0.9	$P(C) = p(0.9) = 0.2$	$P(C \mathcal{D}) = p(0.9 x=1) = 0.2903$

Here are plots of the prior and posterior pmf's from the example.



Prior pmf $p(\theta)$ and posterior pmf $p(\theta|x=1)$ for Example 1

If the data was different then the likelihood column in the Bayesian update table would be different. We can plan for different data by building the entire [likelihood table](#) ahead of time. In the coin example there are two possibilities for the data: the toss is heads or the toss is tails. So the full likelihood table has two likelihood columns:

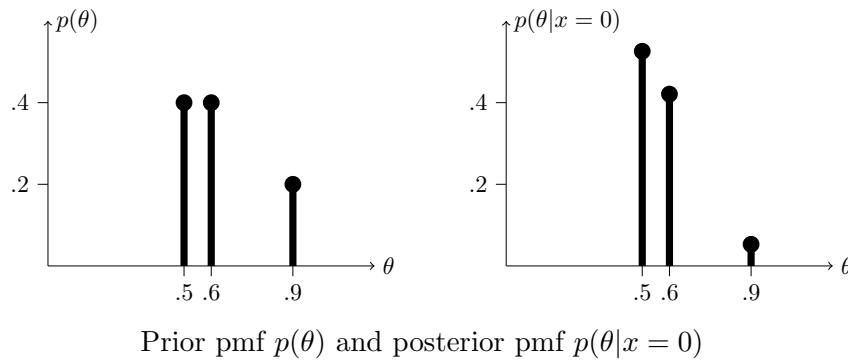
hypothesis	likelihood $p(x \theta)$	
θ	$p(x=0 \theta)$	$p(x=1 \theta)$
0.5	0.5	0.5
0.6	0.4	0.6
0.9	0.1	0.9

Example 2. Using the notation $p(\theta)$, etc., redo Example 1 assuming the flip was tails.

answer: Since the data has changed, the likelihood column in the Bayesian update table is now for $x = 0$. That is, we must take the $p(x=0|\theta)$ column from the likelihood table.

hypothesis	prior	likelihood	Bayes	
			numerator	posterior
θ	$p(\theta)$	$p(x=0 \theta)$	$p(x=0 \theta)p(\theta)$	$p(\theta x=0)$
0.5	0.4	0.5	0.2	0.5263
0.6	0.4	0.4	0.16	0.4211
0.9	0.2	0.1	0.02	0.0526
total	1		0.38	1

Now the probability of type A has increased from 0.4 to 0.5263, while the probability of type C has decreased from 0.2 to only 0.0526. Here are the corresponding plots:



3.3 Food for thought.

Suppose that in Example 1 you didn't know how many coins of each type were in the drawer. You picked one at random and got heads. How would you go about deciding which hypothesis (coin type) if any was most supported by the data?

4 Updating again and again

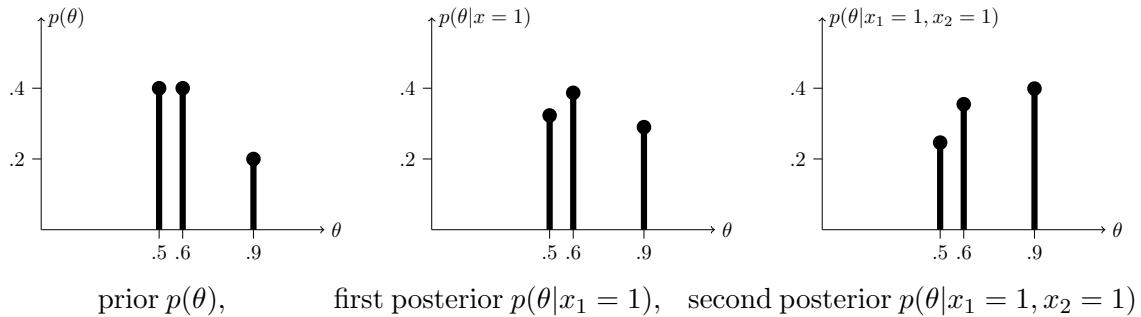
In life we are continually updating our beliefs with each new experience of the world. In Bayesian inference, after updating the prior to the posterior, we can take more data and update again! For the second update, the posterior from the first data becomes the prior for the second data.

Example 3. Suppose you have picked a coin as in Example 1. You flip it once and get heads. Then you flip the same coin and get heads again. What is the probability that the coin was type A? Type B? Type C?

answer: As we update several times the table gets big, so we use a smaller font to fit it in:

hypothesis	prior	Bayes			Bayes			posterior 2
		likelihood 1	numerator 1	likelihood 2	numerator 2			
θ	$p(\theta)$	$p(x_1 = 1 \theta)$	$p(x_1 = 1 \theta)p(\theta)$	$p(x_2 = 1 \theta)$	$p(x_2 = 1 \theta)p(x_1 = 1 \theta)p(\theta)$			
0.5	0.4	0.5	0.2	0.5	0.1			0.2463
0.6	0.4	0.6	0.24	0.6	0.144			0.3547
0.9	0.2	0.9	0.18	0.9	0.162			0.3990
total	1				0.406			1

Note that the second Bayes numerator is computed by multiplying the first Bayes numerator and the second likelihood; since we are only interested in the final posterior, there is no need to normalize until the last step. As shown in the last column and plot, after two heads the type C hypothesis has finally taken the lead!



5 Appendix: the base rate fallacy

Example 4. A screening test for a disease is both sensitive and specific. By that we mean it is usually positive when testing a person with the disease and usually negative when testing someone without the disease. Let's assume the true positive rate is 99% and the false positive rate is 2%. Suppose the prevalence of the disease in the general population is 0.5%. If a random person tests positive, what is the probability that they have the disease?

answer: As a review we first do the computation using trees. Next we will redo the computation using tables.

Let's use notation established above for hypotheses and data: let \mathcal{H}_+ be the hypothesis (event) that the person has the disease and let \mathcal{H}_- be the hypothesis they do not. Likewise, let \mathcal{T}_+ and \mathcal{T}_- represent the data of a positive and negative screening test respectively. We are asked to compute $P(\mathcal{H}_+|\mathcal{T}_+)$.

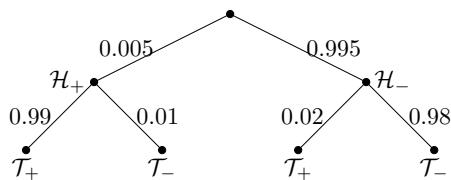
We are given

$$P(\mathcal{T}_+|\mathcal{H}_+) = 0.99, \quad P(\mathcal{T}_+|\mathcal{H}_-) = 0.02, \quad P(\mathcal{H}_+) = 0.005.$$

From these we can compute the false negative and true negative rates:

$$P(\mathcal{T}_-|\mathcal{H}_+) = 0.01, \quad P(\mathcal{T}_-|\mathcal{H}_-) = 0.98$$

All of these probabilities can be displayed quite nicely in a tree.



Bayes' theorem yields

$$P(\mathcal{H}_+|\mathcal{T}_+) = \frac{P(\mathcal{T}_+|\mathcal{H}_+)P(\mathcal{H}_+)}{P(\mathcal{T}_+)} = \frac{0.99 \cdot 0.005}{0.99 \cdot 0.005 + 0.02 \cdot 0.995} = 0.19920 \approx 20\%$$

Now we redo this calculation using a Bayesian update table:

hypothesis	prior	likelihood	Bayes numerator	posterior
\mathcal{H}	$P(\mathcal{H})$	$P(\mathcal{T}_+ \mathcal{H})$	$P(\mathcal{T}_+ \mathcal{H})P(\mathcal{H})$	$P(\mathcal{H} \mathcal{T}_+)$
\mathcal{H}_+	0.005	0.99	0.00495	0.19920
\mathcal{H}_-	0.995	0.02	0.01990	0.80080
total	1	NO SUM	0.02485	1

The table shows that the posterior probability $P(\mathcal{H}_+|\mathcal{T}_+)$ that a person with a positive test has the disease is about 20%. This is far less than the sensitivity of the test (99%) but much higher than the prevalence of the disease in the general population (0.5%).

Bayesian Updating: Probabilistic Prediction

Class 12, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to use the law of total probability to compute prior and posterior predictive probabilities.

2 Introduction

In the previous class we looked at updating the probability of hypotheses based on data. We can also use the data to update the probability of each possible outcome of a future experiment. In this class we will look at how this is done.

2.1 Probabilistic prediction; words of estimative probability (WEP)

There are many ways to word predictions:

- Prediction: “It will rain tomorrow.”
- Prediction using words of estimative probability (WEP): “It is likely to rain tomorrow.”
- Probabilistic prediction: “Tomorrow it will rain with probability 60% (and not rain with probability 40%).”

Each type of wording is appropriate at different times.

In this class we are going to focus on probabilistic prediction and precise quantitative statements. You can see http://en.wikipedia.org/wiki/Words_of_Estimative_Probability for an interesting discussion about the appropriate use of words of estimative probability. The article also contains a list of *weasel words* such as ‘might’, ‘cannot rule out’, ‘it’s conceivable’ that should be avoided as almost certain to cause confusion.

There are many places where we want to make a probabilistic prediction. Examples are

- Medical treatment outcomes
- Weather forecasting
- Climate change
- Sports betting
- Elections
- ...

These are all situations where there is uncertainty about the outcome and we would like as precise a description of what could happen as possible.

3 Predictive Probabilities

Probabilistic prediction simply means assigning a probability to each possible outcomes of an experiment.

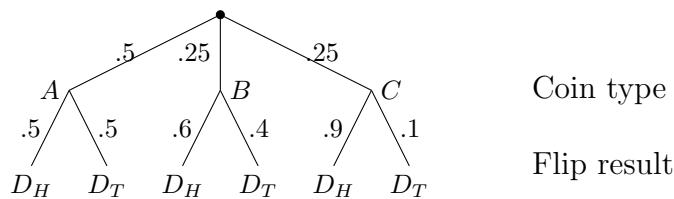
Recall the coin example from the previous class notes: there are three types of coins which are indistinguishable apart from their probability of landing heads when tossed.

- Type A coins are fair, with probability 0.5 of heads
 - Type B coins have probability 0.6 of heads
 - Type C coins have probability 0.9 of heads

You have a drawer containing 4 coins: 2 of type A , 1 of type B , and 1 of type C . You reach into the drawer and pick a coin at random. We let A stand for the event ‘the chosen coin is of type A ’. Likewise for B and C .

3.1 Prior predictive probabilities

Before taking data we can compute the probability that our chosen coin will land heads (or tails) if flipped. Let D_H be the event it lands heads and let D_T the event it lands tails. We can use the [law of total probability](#) to determine the probabilities of these events. Either by drawing a tree or directly proceeding to the algebra, we get:



$$\begin{aligned} P(D_H) &= P(D_H|A)P(A) + P(D_H|B)P(B) + P(D_H|C)P(C) \\ &= 0.5 \cdot 0.5 + 0.6 \cdot 0.25 + 0.9 \cdot 0.25 = 0.625 \end{aligned}$$

$$\begin{aligned} P(D_T) &= P(D_T|A)P(A) + P(D_T|B)P(B) + P(D_T|C)P(C) \\ &= 0.5 \cdot 0.5 + 0.4 \cdot 0.25 + 0.1 \cdot 0.25 = 0.375 \end{aligned}$$

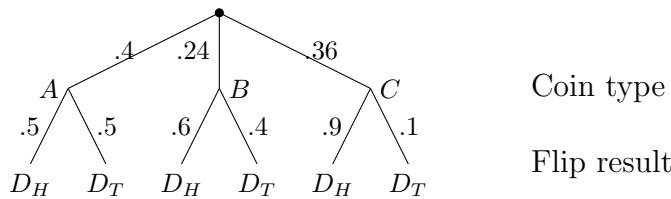
Definition: These probabilities give a (probabilistic) prediction of what will happen if the coin is tossed. Because they are computed before we collect any data they are called [prior predictive probabilities](#).

3.2 Posterior predictive probabilities

Suppose we flip the coin once and it lands heads. We now have data D , which we can use to update the prior probabilities of our hypotheses to posterior probabilities. Last class we learned to use a Bayes table to facilitate this computation:

hypothesis	prior	likelihood	Bayes numerator	posterior
H	$P(H)$	$P(D H)$	$P(D H)P(H)$	$P(H D)$
A	0.5	0.5	0.25	0.4
B	0.25	0.6	0.15	0.24
C	0.25	0.9	0.225	0.36
total	1		0.625	1

Having flipped the coin once and gotten heads, we can compute the probability that our chosen coin will land heads (or tails) if flipped a second time. We proceed just as before, but using the posterior probabilities $P(A|D)$, $P(B|D)$, $P(C|D)$ in place of the prior probabilities $P(A)$, $P(B)$, $P(C)$.



$$\begin{aligned} P(D_H|D) &= P(D_H|A)P(A|D) + P(D_H|B)P(B|D) + P(D_H|C)P(C|D) \\ &= 0.5 \cdot 0.4 + 0.6 \cdot 0.24 + 0.9 \cdot 0.36 = 0.668 \end{aligned}$$

$$\begin{aligned} P(D_T|D) &= P(D_T|A)P(A|D) + P(D_T|B)P(B|D) + P(D_T|C)P(C|D) \\ &= 0.5 \cdot 0.4 + 0.4 \cdot 0.24 + 0.1 \cdot 0.36 = 0.332 \end{aligned}$$

Definition: These probabilities give a (probabilistic) prediction of what will happen if the coin is tossed again. Because they are computed after collecting data and updating the prior to the posterior, they are called **posterior predictive probabilities**.

Note that heads on the first toss increases the probability of heads on the second toss.

3.3 Review

Here's a succinct description of the preceding sections that may be helpful:

Each hypothesis gives a different probability of heads, so the total probability of heads is a weighted average. For the prior predictive probability of heads, the weights are given by the prior probabilities of the hypotheses. For the posterior predictive probability of heads, the weights are given by the posterior probabilities of the hypotheses.

Remember: Prior and posterior probabilities are for hypotheses. Prior predictive and posterior predictive probabilities are for data. To keep this straight, remember that the latter **predict** future data.

Bayesian Updating: Odds

Class 12, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to convert between odds and probability.
2. Be able to update prior odds to posterior odds using Bayes factors.
3. Understand how Bayes factors measure the extent to which data provides evidence for or against a hypothesis.

2 Odds

When comparing two events, it common to phrase probability statements in terms of odds.

Definition The **odds** of event E versus event E' are the ratio of their probabilities $P(E)/P(E')$. If unspecified, the second event is assumed to be the complement E^c . So the **odds** of E are:

$$O(E) = \frac{P(E)}{P(E^c)}.$$

For example, $O(\text{rain}) = 2$ means that the probability of rain is twice the probability of no rain ($2/3$ versus $1/3$). We might say ‘the odds of rain are 2 to 1.’

Example. For a fair coin, $O(\text{heads}) = \frac{1/2}{1/2} = 1$. We might say the odds of heads are **1 to 1** or **fifty-fifty**.

Example. For a standard die, the odds of rolling a 4 are $\frac{1/6}{5/6} = \frac{1}{5}$. We might say the odds are ‘1 to 5 for’ or ‘5 to 1 against’ rolling a 4.

Example. The probability of a pair in a five card poker hand is 0.42257. So the odds of a pair are $0.42257/(1-0.42257) = 0.73181$.

We can go back and forth between probability and odds as follows.

Conversion formulas: if $P(E) = p$ then $O(E) = \frac{p}{1-p}$. If $O(E) = q$ then $P(E) = \frac{q}{1+q}$.

Notes:

1. The second formula simply solves $q = p/(1-p)$ for p .
2. Probabilities are between 0 and 1, while odds are between 0 to ∞ .
3. The property $P(E^c) = 1 - P(E)$ becomes $O(E^c) = 1/O(E)$.

Example. Let F be the event that a five card poker hand is a full house. Then $P(F) = 0.00145214$ so $O(F) = 0.0014521/(1 - 0.0014521) = 0.0014542$.

The odds not having a full house are $O(F^c) = (1 - 0.0014521)/0.0014521 = 687 = 1/O(F)$.

4. If $P(E)$ or $O(E)$ is small then $O(E) \approx P(E)$. This follows from the conversion formulas.

Example. In the poker example where F = ‘full house’ we saw that $P(F)$ and $O(F)$ differ only in the fourth significant digit.

3 Updating odds

3.1 Introduction

In Bayesian updating, we used the likelihood of data to update prior probabilities of hypotheses to posterior probabilities. In the language of odds, we will update [prior odds](#) to [posterior odds](#). One of our key points will be that the data can provide evidence supporting or negating a hypothesis depending on whether its posterior odds are greater or less than its prior odds.

3.2 Example: Marfan syndrome

Marfan syndrome is a genetic disease of connective tissue that occurs in 1 of every 15000 people. The main ocular features of Marfan syndrome include bilateral ectopia lentis (lens dislocation), myopia and retinal detachment. About 70% of people with Marfan syndrome have at least one of these ocular features; only 7% of people without Marfan syndrome do. (We don’t guarantee the accuracy of these numbers, but they will work perfectly well for our example.)

If a person has at least one of these ocular features, what are the odds that they have Marfan syndrome?

answer: This is a standard Bayesian updating problem. Our hypotheses are:

M = ‘the person has Marfan syndrome’

M^c = ‘the person does not have Marfan syndrome’

The data is:

F = ‘the person has at least one ocular feature’.

We are given the prior probability of M and the likelihoods of F given M or M^c :

$$P(M) = 1/15000, \quad P(F|M) = 0.7, \quad P(F|M^c) = 0.07.$$

As before, we can compute the posterior probabilities using a table:

hypothesis	prior	likelihood	numerator	posterior
H	$P(H)$	$P(F H)$	$P(F H)P(H)$	$P(H F)$
M	0.000067	0.7	0.0000467	0.00066
M^c	0.999933	0.07	0.069995	0.99933
total	1		0.07004	1

First we find the prior odds:

$$O(M) = \frac{P(M)}{P(M^c)} = \frac{1/15000}{14999/15000} = \frac{1}{14999} \approx 0.000067.$$

The posterior odds are given by the ratio of the posterior probabilities or the Bayes numerators, since the normalizing factor will be the same in both numerator and denominator.

$$O(M|F) = \frac{P(M|F)}{P(M^c|F)} = \frac{P(F|M)P(M)}{P(F|M^c)P(M^c)} = 0.000667.$$

The posterior odds are a factor of 10 larger than the prior odds. In that sense, having an ocular feature is strong evidence in favor of the hypothesis M . However, because the prior odds are so small, it is still highly unlikely the person has Marfan syndrome.

4 Bayes factors and strength of evidence

The factor of 10 in the previous example is called a Bayes factor. The exact definition is the following.

Definition: For a hypothesis H and data D , the **Bayes factor** is the ratio of the likelihoods:

$$\text{Bayes factor} = \frac{P(D|H)}{P(D|H^c)}.$$

Let's see exactly where the Bayes factor arises in updating odds. We have

$$\begin{aligned} O(H|D) &= \frac{P(H|D)}{P(H^c|D)} \\ &= \frac{P(D|H) P(H)}{P(D|H^c) P(H^c)} \\ &= \frac{P(D|H)}{P(D|H^c)} \cdot \frac{P(H)}{P(H^c)} \\ &= \frac{P(D|H)}{P(D|H^c)} \cdot O(H) \end{aligned}$$

$$\text{posterior odds} = \text{Bayes factor} \times \text{prior odds}$$

From this formula, we see that the Bayes' factor (BF) tells us whether the data provides evidence for or against the hypothesis.

- If $BF > 1$ then the posterior odds are greater than the prior odds. So the data provides evidence for the hypothesis.
- If $BF < 1$ then the posterior odds are less than the prior odds. So the data provides evidence against the hypothesis.
- If $BF = 1$ then the prior and posterior odds are equal. So the data provides no evidence either way.

The following example is taken from the textbook *Information Theory, Inference, and Learning Algorithms* by David J. C. MacKay, who has this to say regarding trial evidence.

In my view, a jury's task should generally be to multiply together carefully evaluated likelihood ratios from each independent piece of admissible evidence with an equally carefully reasoned prior probability. This view is shared by many statisticians but learned British appeal judges recently disagreed and actually overturned the verdict of a trial because the jurors *had* been taught to use Bayes' theorem to handle complicated DNA evidence.

Example 1. Two people have left traces of their own blood at the scene of a crime. A suspect , Oliver, is tested and found to have type 'O' blood. The blood groups of the two traces are found to be of type 'O' (a common type in the local population, having frequency 60%) and type 'AB' (a rare type, with frequency 1%). Does this data (type 'O' and 'AB' blood were found at the scene) give evidence in favor of the proposition that Oliver was one of the two people present at the scene of the crime?"

answer: There are two hypotheses:

S = 'Oliver and another unknown person were at the scene of the crime'

S^c = 'two unknown people were at the scene of the crime'

The data is:

D = 'type 'O' and 'AB' blood were found'

The Bayes factor for Oliver's presence is $BF_{\text{Oliver}} = \frac{P(D|S)}{P(D|S^c)}$. We compute the numerator and denominator of this separately.

The data says that both type O and type AB blood were found. If Oliver was at the scene then 'type O' blood would be there. So $P(D|S)$ is the probability that the other person had type AB blood. We are told this is .01, so $P(D|S) = 0.01$.

If Oliver was not at the scene then there were two random people one with type O and one with type AB blood. The probability of this is $2 \cdot 0.6 \cdot 0.01$. The factor of 2 is because there are two ways this can happen –the first person is type O and the second is type AB or vice versa.*

Thus the Bayes factor for Oliver's presence is

$$BF_{\text{Oliver}} = \frac{P(D|S)}{P(D|S^c)} = \frac{0.01}{2 \cdot 0.6 \cdot 0.01} = 0.83.$$

Since $BF_{\text{Oliver}} < 1$, the data provides (weak) evidence against Oliver being at the scene.

*We have assumed the blood types of the two people are independent. This is not precisely true, but for a large population it is close enough. The exact probability is $\frac{2 \cdot N_O \cdot N_{AB}}{N \cdot (N - 1)}$ where N_O is the number of people with type O blood, N_{AB} the number with type AB blood and N the size of the population. We have $\frac{N_O}{N} = 0.6$. For large N we have $N \approx N - 1$, so $\frac{N_{AB}}{N - 1} \approx 0.01$. This shows the probability is approximately $2 \cdot 0.6 \cdot 0.01$ as claimed.

Example 2. Another suspect Alberto is found to have type 'AB' blood. Do the same data give evidence in favor of the proposition that Alberto was one of the two people present at the crime?

answer: Reusing the above notation with Alberto in place of Oliver we have:

$$BF_{\text{Alberto}} = \frac{P(D|S)}{P(D|S^c)} = \frac{0.6}{2 \cdot 0.6 \cdot 0.01} = 50.$$

Since $BF_{\text{Alberto}} \gg 1$, the data provides strong evidence in favor of Alberto being at the scene.

Notes:

1. In both examples, we have only computed the Bayes factor, not the posterior odds. To compute the latter, we would need to know the prior odds that Oliver (or Alberto) was at the scene based on other evidence.
2. Note that if 50% of the population had type O blood instead of 60%, then the Oliver's Bayes factor would be 1 (neither for nor against). More generally, the break-even point for blood type evidence is when the proportion of the suspect's blood type in the general population equals the proportion of the suspect's blood type among those who left blood at the scene.

4.1 Updating again and again

Suppose we collect data in two stages, first D_1 , then D_2 . We have seen in our dice and coin examples that the final posterior can be computed all at once or in two stages where we first update the prior using the likelihoods for D_1 and then update the resulting posterior using the likelihoods for D_2 . The latter approach works whenever likelihoods multiply:

$$P(D_1, D_2 | H) = P(D_1 | H)P(D_2 | H).$$

Since likelihoods are conditioned on hypotheses, we say that D_1 and D_2 are **conditionally independent** if the above equation holds for every hypothesis H .

Example. There are five dice in a drawer, with 4, 6, 8, 12, and 20 sides (these are the hypotheses). I pick a die at random and roll it twice. The first roll gives 7. The second roll gives 11. Are these results conditionally independent? Are they independent?

answer: These results are conditionally independent. For example, for the hypothesis of the 8-sided die we have:

$$\begin{aligned} P(7 \text{ on roll 1} | 8\text{-sided die}) &= 1/8 \\ P(11 \text{ on roll 2} | 8\text{-sided die}) &= 0 \\ P(7 \text{ on roll 1, } 11 \text{ on roll 2} | 8\text{-sided die}) &= 0 \end{aligned}$$

For the hypothesis of the 20-sided die we have:

$$\begin{aligned} P(7 \text{ on roll 1} | 20\text{-sided die}) &= 1/20 \\ P(11 \text{ on roll 2} | 20\text{-sided die}) &= 1/20 \\ P(7 \text{ on roll 1, } 11 \text{ on roll 2} | 20\text{-sided die}) &= (1/20)^2 \end{aligned}$$

However, the results of the rolls are *not* independent. That is:

$$P(7 \text{ on roll 1, } 11 \text{ on roll 2}) \neq P(7 \text{ on roll 1})P(11 \text{ on roll 2}).$$

Intuitively, this is because a 7 on the roll 1 allows us to rule out the 4- and 6-sided dice, making an 11 on roll 2 more likely. Let's check this intuition by computing both sides precisely. On the righthand side we have:

$$\begin{aligned} P(7 \text{ on roll 1}) &= \frac{1}{5} \cdot \frac{1}{8} + \frac{1}{5} \cdot \frac{1}{12} + \frac{1}{5} \cdot \frac{1}{20} = \frac{31}{600} \\ P(11 \text{ on roll 2}) &= \frac{1}{5} \cdot \frac{1}{12} + \frac{1}{5} \cdot \frac{1}{20} = \frac{2}{75} \end{aligned}$$

On the lefthand side we have:

$$\begin{aligned} P(7 \text{ on roll 1}, 11 \text{ on roll 2}) &= P(11 \text{ on roll 2} | 7 \text{ on roll 1})P(7 \text{ on roll 1}) \\ &= \left(\frac{30}{93} \cdot \frac{1}{12} + \frac{6}{31} \cdot \frac{1}{20} \right) \cdot \frac{31}{600} \\ &= \frac{17}{465} \cdot \frac{31}{600} = \frac{17}{9000} \end{aligned}$$

Here $\frac{30}{93}$ and $\frac{6}{31}$ are the posterior probabilities of the 12- and 20-sided dice given a 7 on roll 1. We conclude that, without conditioning on hypotheses, the rolls are not independent.

Returning to general setup, if D_1 and D_2 are conditionally independent for H and H^c then it makes sense to consider each Bayes factor independently:

$$BF_i = \frac{P(D_i|H)}{P(D_i|H^c)}.$$

The prior odds of H are $O(H)$. The posterior odds after D_1 are

$$O(H|D_1) = BF_1 \cdot O(H).$$

And the posterior odds after D_1 and D_2 are

$$\begin{aligned} O(H|D_1, D_2) &= BF_2 \cdot O(H|D_1) \\ &= BF_2 \cdot BF_1 \cdot O(H) \end{aligned}$$

We have the beautifully simple notion that updating with new data just amounts to multiplying the current posterior odds by the Bayes factor of the new data.

Example 3. Other symptoms of Marfan Syndrome

Recall from the earlier example that the Bayes factor for at least one ocular feature (F) is

$$BF_F = \frac{P(F|M)}{P(F|M^c)} = \frac{0.7}{0.07} = 10.$$

The wrist sign (W) is the ability to wrap one hand around your other wrist to cover your pinky nail with your thumb. Assume 10% of the population have the wrist sign, while 90% of people with Marfan's have it. Therefore the Bayes factor for the wrist sign is

$$BF_W = \frac{P(W|M)}{P(W|M^c)} = \frac{0.9}{0.1} = 9.$$

We will assume that F and W are conditionally independent symptoms. That is, among people with Marfan syndrome, ocular features and the wrist sign are independent, and among people without Marfan syndrome, ocular features and the wrist sign are independent. Given this assumption, the posterior odds of Marfan syndrome for someone with both an ocular feature and the wrist sign are

$$O(M|F, W) = BF_W \cdot BF_F \cdot O(M) = 9 \cdot 10 \cdot \frac{1}{14999} \approx \frac{6}{1000}.$$

We can convert the posterior odds back to probability, but since the odds are so small the result is nearly the same:

$$P(M|F, W) \approx \frac{6}{1000 + 6} \approx 0.596\%.$$

So ocular features and the wrist sign are both strong evidence in favor of the hypothesis M , and taken together they are very strong evidence. Again, because the prior odds are so small, it is still unlikely that the person has Marfan syndrome, but at this point it might be worth undergoing further testing given potentially fatal consequences of the disease (such as aortic aneurysm or dissection).

Note also that if a person has exactly one of the two symptoms, then the product of the Bayes factors is near 1 (either 9/10 or 10/9). So the two pieces of data essentially cancel each other out with regard to the evidence they provide for Marfan's syndrome.

5 Log odds

In practice, people often find it convenient to work with the natural log of the odds in place of odds. Naturally enough these are called the **log odds**. The Bayesian update formula

$$O(H|D_1, D_2) = BF_2 \cdot BF_1 \cdot O(H)$$

becomes

$$\ln(O(H|D_1, D_2)) = \ln(BF_2) + \ln(BF_1) + \ln(O(H)).$$

We can interpret the above formula for the posterior log odds as the sum of the prior log odds and all the evidence $\ln(BF_i)$ provided by the data. Note that by taking logs, evidence in favor ($BF_i > 1$) is positive and evidence against ($BF_i < 1$) is negative.

To avoid lengthier computations, we will work with odds rather than log odds in this course. Log odds are nice because sums are often more intuitive than products. Log odds also play a central role in logistic regression, an important statistical model related to linear regression.

Bayesian Updating with Continuous Priors

Class 13, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Understand a parameterized family of distributions as representing a continuous range of hypotheses for the observed data.
2. Be able to state Bayes' theorem and the law of total probability for continuous densities.
3. Be able to apply Bayes' theorem to update a prior probability density function to a posterior pdf given data and a likelihood function.
4. Be able to interpret and compute posterior predictive probabilities.

2 Introduction

Up to now we have only done Bayesian updating when we had a finite number of hypothesis, e.g. our dice example had five hypotheses (4, 6, 8, 12 or 20 sides). Now we will study Bayesian updating when there is a [continuous range of hypotheses](#). The Bayesian update process will be essentially the same as in the discrete case. As usual when moving from discrete to continuous we will need to replace the probability mass function by a probability density function, and sums by integrals.

The first few sections of this note are devoted to working with pdfs. In particular we will cover the law of total probability and Bayes' theorem. We encourage you to focus on how these are essentially identical to the discrete versions. After that, we will apply Bayes' theorem and the law of total probability to Bayesian updating.

3 Examples with continuous ranges of hypotheses

Here are three standard examples with continuous ranges of hypotheses.

Example 1. Suppose you have a system that can succeed or fail with probability p . Then we can hypothesize that p is anywhere in the range $[0, 1]$. That is, we have a continuous range of hypotheses. We will often model this example with a ‘bent’ coin with unknown probability p of heads.

Example 2. The lifetime of a certain isotope is modeled by an exponential distribution $\exp(\lambda)$. In principle, the mean lifetime $1/\lambda$ can be any real number in $(0, \infty)$.

Example 3. We are not restricted to a single parameter. In principle, the parameters μ and σ of a normal distribution can be any real numbers in $(-\infty, \infty)$ and $(0, \infty)$, respectively. If we model gestational length for single births by a normal distribution, then from millions of data points we know that μ is about 40 weeks and σ is about one week.

In all of these examples we modeled the random process giving rise to the data by a distribution with parameters –called a **parametrized distribution**. Every possible choice of the parameter(s) is a hypothesis, e.g. we can hypothesize that the probability of success in Example 1 is $p = 0.7313$. We have a continuous set of hypotheses because we could take any value between 0 and 1.

4 Notational conventions

4.1 Parametrized models

As in the examples above our hypotheses often take the form a certain parameter has value θ . We will often use the letter θ to stand for an arbitrary hypothesis. This will leave symbols like p , f , and x to take their usual meanings as pmf, pdf, and data. Also, rather than saying ‘the hypothesis that the parameter of interest has value θ ’ we will simply say the hypothesis θ .

4.2 Big and little letters

We have two parallel notations for outcomes and probability:

1. (Big letters) Event A , probability function $P(A)$.
2. (Little letters) Value x , pmf $p(x)$ or pdf $f(x)$.

These notations are related by $P(X = x) = p(x)$, where x is a value the discrete random variable X and ‘ $X = x$ ’ is the corresponding event.

We carry these notations over to the probabilities used in Bayesian updating.

1. (Big letters) From hypotheses \mathcal{H} and data \mathcal{D} we compute several associated probabilities

$$P(\mathcal{H}), P(\mathcal{D}), P(\mathcal{H}|\mathcal{D}), P(\mathcal{D}|\mathcal{H}).$$

In the coin example we might have \mathcal{H} = ‘the chosen coin has probability 0.6 of heads’, \mathcal{D} = ‘the flip was heads’, and $P(\mathcal{D}|\mathcal{H}) = 0.6$

2. (Small letters) Hypothesis values θ and data values x both have probabilities or probability densities:

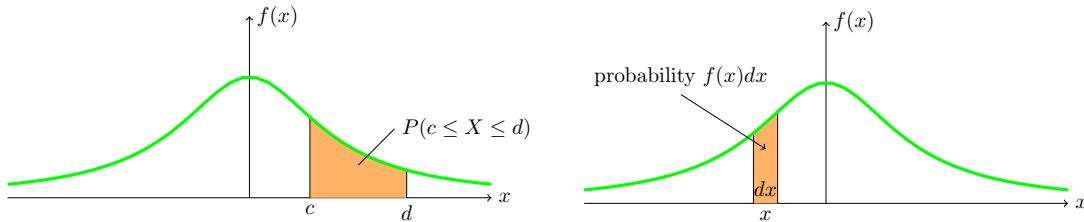
$$\begin{array}{cccc} p(\theta) & p(x) & p(\theta|x) & p(x|\theta) \\ f(\theta) & f(x) & f(\theta|x) & f(x|\theta) \end{array}$$

In the coin example we might have $\theta = 0.6$ and $x = 1$, so $p(x|\theta) = 0.6$. We might also write $p(x = 1|\theta = 0.6)$ to emphasize the values of x and θ , but we will never just write $p(1|0.6)$ because it is unclear which value is x and which is θ .

Although we will still use both types of notation, from now on we will mostly use the small letter notation involving pmfs and pdfs. Hypotheses will usually be parameters represented by Greek letters ($\theta, \lambda, \mu, \sigma, \dots$) while data values will usually be represented by English letters (x, x_i, y, \dots).

5 Quick review of pdf and probability

Suppose X is a random variable with pdf $f(x)$. Recall $f(x)$ is a density; its units are probability/(units of x).



The probability that the value of X is in $[c, d]$ is given by

$$\int_c^d f(x) dx.$$

The probability that X is in an infinitesimal range dx around x is $f(x) dx$. In fact, the integral formula is just the ‘sum’ of these infinitesimal probabilities. We can visualize these probabilities by viewing the integral as area under the graph of $f(x)$.

In order to manipulate probabilities instead of densities in what follows, we will make frequent use of the notion that $f(x) dx$ is the probability that X is in an infinitesimal range around x of width dx . Please make sure that you fully understand this notion.

6 Continuous priors, discrete likelihoods

In the Bayesian framework we have probabilities of hypotheses –called prior and posterior probabilities– and probabilities of data given a hypothesis –called likelihoods. In earlier classes both the hypotheses and the data had discrete ranges of values. We saw in the introduction that we might have a continuous range of hypotheses. The same is true for the data, but for today we will assume that our data can only take a discrete set of values. In this case, the likelihood of data x given hypothesis θ is written using a pmf: $p(x|\theta)$.

We will use the following coin example to explain these notions. We will carry this example through in each of the succeeding sections.

Example 4. Suppose we have a bent coin with unknown probability θ of heads. The value of θ is random and could be anywhere between 0 and 1. For this and the examples that follow we’ll suppose that the value of θ follows a distribution with **continuous prior probability density** $f(\theta) = 2\theta$. We have a **discrete likelihood** because tossing a coin has only two outcomes, $x = 1$ for heads and $x = 0$ for tails.

$$p(x = 1|\theta) = \theta, \quad p(x = 0|\theta) = 1 - \theta.$$

Think: This can be tricky to wrap your mind around. We have a coin with an unknown probability θ of heads. The value of the parameter θ is itself random and has a prior pdf $f(\theta)$. It may help to see that the discrete examples we did in previous classes are similar. For example, we had a coin that might have probability of heads 0.5, 0.6, or 0.9. So,

we called our hypotheses $H_{0.5}$, $H_{0.6}$, $H_{0.9}$ and these had prior probabilities $P(H_{0.5})$ etc. In other words, we had a coin with an unknown probability of heads, we had hypotheses about that probability and each of these hypotheses had a prior probability.

7 The law of total probability

The law of total probability for continuous probability distributions is essentially the same as for discrete distributions. We replace the prior pmf by a prior pdf and the sum by an integral. We start by reviewing the law for the discrete case.

Recall that for a discrete set of hypotheses $\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_n$ the law of total probability says

$$P(\mathcal{D}) = \sum_{i=1}^n P(\mathcal{D}|\mathcal{H}_i)P(\mathcal{H}_i). \quad (1)$$

This is the total **prior probability** of \mathcal{D} because we used the prior probabilities $P(\mathcal{H}_i)$

In the little letter notation with $\theta_1, \theta_2, \dots, \theta_n$ for hypotheses and x for data the law of total probability is written

$$p(x) = \sum_{i=1}^n p(x|\theta_i)p(\theta_i). \quad (2)$$

We also called this the **prior predictive probability** of the outcome x to distinguish it from the prior probability of the hypothesis θ .

Likewise, there is a law of total probability for continuous pdfs. We state it as a theorem using little letter notation.

Theorem. Law of total probability. Suppose we have a continuous parameter θ in the range $[a, b]$, and discrete random data x . Assume θ is itself random with density $f(\theta)$ and that x and θ have likelihood $p(x|\theta)$. In this case, the total probability of x is given by the formula.

$$p(x) = \int_a^b p(x|\theta)f(\theta) d\theta \quad (3)$$

Proof. Our proof will be by analogy to the discrete version: The probability term $p(x|\theta)f(\theta) d\theta$ is perfectly analogous to the term $p(x|\theta_i)p(\theta_i)$ in Equation 2 (or the term $P(\mathcal{D}|\mathcal{H}_i)P(\mathcal{H}_i)$ in Equation 1). Continuing the analogy: the sum in Equation 2 becomes the integral in Equation 3

As in the discrete case, when we think of θ as a hypothesis explaining the probability of the data we call $p(x)$ the **prior predictive probability for x** .

Example 5. (Law of total probability.) Continuing with Example 4. We have a bent coin with probability θ of heads. The value of θ is random with prior pdf $f(\theta) = 2\theta$ on $[0, 1]$.

Suppose I flip the coin once. What is the total probability of heads?

answer: In Example 4 we noted that the likelihoods are $p(x = 1|\theta) = \theta$ and $p(x = 0|\theta) = 1 - \theta$. So the total probability of $x = 1$ is

$$p(x = 1) = \int_0^1 p(x = 1|\theta) f(\theta) d\theta = \int_0^1 \theta \cdot 2\theta d\theta = \int_0^1 2\theta^2 d\theta = \frac{2}{3}.$$

Since the prior is weighted towards higher probabilities of heads, so is the total probability.

8 Bayes' theorem for continuous probability densities

The statement of Bayes' theorem for continuous pdfs is essentially identical to the statement for pmfs. We state it including $d\theta$ so we have genuine probabilities:

Theorem. Bayes' Theorem. Use the same assumptions as in the law of total probability, i.e. θ is a continuous parameter with pdf $f(\theta)$ and range $[a, b]$; x is random discrete data; together they have likelihood $p(x|\theta)$. With these assumptions:

$$f(\theta|x) d\theta = \frac{p(x|\theta)f(\theta) d\theta}{p(x)} = \frac{p(x|\theta)f(\theta) d\theta}{\int_a^b p(x|\theta)f(\theta) d\theta}. \quad (4)$$

Proof. Since this is a statement about probabilities it is just the usual statement of Bayes' theorem. This is important enough to warrant spelling it out in words: Let Θ be the random variable that produces the value θ . Consider the events

$$H = \text{'}\Theta\text{ is in an interval of width }d\theta\text{ around the value }\theta\text{'}$$

and

$$D = \text{'the value of the data is }x\text{'}$$

Then $P(H) = f(\theta) d\theta$, $P(D) = p(x)$, and $P(D|H) = p(x|\theta)$. Now our usual form of Bayes' theorem becomes

$$f(\theta|x) d\theta = P(H|D) = \frac{P(D|H)P(H)}{P(D)} = \frac{p(x|\theta)f(\theta) d\theta}{p(x)}$$

Looking at the first and last terms in this equation we see the new form of Bayes' theorem.

Finally, we firmly believe that is is more conducive to careful thinking about probability to keep the factor of $d\theta$ in the statement of Bayes' theorem. But because it appears in the numerator on both sides of Equation 4 many people drop the $d\theta$ and write Bayes' theorem in terms of densities as

$$f(\theta|x) = \frac{p(x|\theta)f(\theta)}{p(x)} = \frac{p(x|\theta)f(\theta)}{\int_a^b p(x|\theta)f(\theta) d\theta}.$$

9 Bayesian updating with continuous priors

Now that we have Bayes' theorem and the law of total probability we can finally get to Bayesian updating. Before continuing with Example 4, we point out two features of the Bayesian updating table that appears in the next example:

1. The table for continuous priors is very simple: since we cannot have a row for each of an infinite number of hypotheses we'll have just **one row which uses a variable to stand for all hypotheses θ** .
2. By including $d\theta$, all the entries in the table are probabilities and all our usual probability rules apply.

Example 6. (**Bayesian updating.**) Continuing Examples 4 and 5. We have a bent coin with unknown probability θ of heads. The value of θ is random with prior pdf $f(\theta) = 2\theta$. Suppose we flip the coin once and get heads. Compute the posterior pdf for θ .

answer: We make an update table with the usual columns. Since this is our first example the first row is the abstract version of Bayesian updating in general and the second row is Bayesian updating for this particular example.

hypothesis	prior	likelihood	Bayes numerator	posterior
θ	$f(\theta) d\theta$	$p(x = 1 \theta)$	$p(x = 1 \theta)f(\theta) d\theta$	$f(\theta x = 1) d\theta$
θ	$2\theta d\theta$	θ	$2\theta^2 d\theta$	$3\theta^2 d\theta$
total	$\int_a^b f(\theta) d\theta = 1$		$p(x = 1) = \int_0^1 2\theta^2 d\theta = 2/3$	1

Therefore the posterior pdf (after seeing 1 heads) is $f(\theta|x) = 3\theta^2$.

We have a number of comments:

1. Since we used the prior probability $f(\theta) d\theta$, the hypothesis should have been:
'the unknown parameter is in an interval of width $d\theta$ around θ '.
Even for us that is too much to write, so you will have to think it everytime we write that the hypothesis is θ .
2. The **posterior pdf** for θ is found by removing the $d\theta$ from the posterior probability in the table.
$$f(\theta|x) = 3\theta^2.$$
3. (i) As always $p(x)$ is the **total probability**. Since we have a continuous distribution instead of a sum we compute an integral.
(ii) Notice that by including $d\theta$ in the table, it is clear what integral we need to compute to find the total probability $p(x)$.
4. The table organizes the continuous version of Bayes' theorem. Namely, the posterior pdf is related to the prior pdf and likelihood function via:

$$f(\theta|x)d\theta = \frac{p(x|\theta) f(\theta)d\theta}{\int_a^b p(x|\theta)f(\theta) d\theta} = \frac{p(x|\theta) f(\theta)}{p(x)}$$

Removing the $d\theta$ in the numerator of both sides we have the statement in terms of densities.

5. Regarding both sides as functions of θ , we can again express Bayes' theorem in the form:
$$f(\theta|x) \propto p(x|\theta) \cdot f(\theta)$$

posterior \propto likelihood \times prior.

9.1 Flat priors

One important prior is called a **flat or uniform prior**. A flat prior assumes that every hypothesis is equally probable. For example, if θ has range $[0, 1]$ then $f(\theta) = 1$ is a flat prior.

Example 7. (Flat priors.) We have a bent coin with unknown probability θ of heads. Suppose we toss it once and get tails. Assume a flat prior and find the posterior probability for θ .

answer: This is the just Example 6 with a change of prior and likelihood.

hypothesis	prior	likelihood	Bayes numerator	posterior
θ	$f(\theta) d\theta$	$p(x = 0 \theta)$		$f(\theta x = 0) d\theta$
θ	$1 \cdot d\theta$	$1 - \theta$	$(1 - \theta) d\theta$	$2(1 - \theta) d\theta$
total	$\int_a^b f(\theta) d\theta = 1$		$p(x = 0) = \int_0^1 (1 - \theta) d\theta = 1/2$	1

9.2 Using the posterior pdf

Example 8. In the previous example the prior probability was flat. First show that this means that a priori the coin is equally like to be biased towards heads or tails. Then, after observing one heads, what is the (posterior) probability that the coin is biased towards heads?

answer: Since the parameter θ is the probability the coin lands heads, the first part of the problem asks us to show $P(\theta > .5) = 0.5$ and the second part asks for $P(\theta > .5 | x = 1)$. These are easily computed from the prior and posterior pdfs respectively.

The prior probability that the coin is biased towards heads is

$$P(\theta > .5) = \int_{.5}^1 f(\theta) d\theta = \int_{.5}^1 1 \cdot d\theta = \theta|_{.5}^1 = \frac{1}{2}.$$

The probability of $1/2$ means the coin is equally likely to be biased toward heads or tails. The posterior probability that it's biased towards heads is

$$P(\theta > .5 | x = 1) = \int_{.5}^1 f(\theta|x = 1) d\theta = \int_{.5}^1 2\theta d\theta = \theta^2|_{.5}^1 = \frac{3}{4}.$$

We see that observing one heads has increased the probability that the coin is biased towards heads from $1/2$ to $3/4$.

10 Predictive probabilities

Just as in the discrete case we are also interested in using the posterior probabilities of the hypotheses to make predictions for what will happen next.

Example 9. (Prior and posterior prediction.) Continuing Examples 4, 5, 6: we have a coin with unknown probability θ of heads and the value of θ has prior pdf $f(\theta) = 2\theta$. Find the prior predictive probability of heads. Then suppose the first flip was heads and find the posterior predictive probabilities of both heads and tails on the second flip.

answer: For notation let x_1 be the result of the first flip and let x_2 be the result of the second flip. The prior predictive probability is exactly the total probability computed in Examples 5 and 6.

$$p(x_1 = 1) = \int_0^1 p(x_1 = 1|\theta)f(\theta) d\theta = \int_0^1 2\theta^2 d\theta = \frac{2}{3}.$$

The posterior predictive probabilities are the total probabilities computed using the posterior pdf. From Example 6 we know the posterior pdf is $f(\theta|x_1 = 1) = 3\theta^2$. So the posterior predictive probabilities are

$$p(x_2 = 1|x_1 = 1) = \int_0^1 p(x_2 = 1|\theta, x_1 = 1)f(\theta|x_1 = 1) d\theta = \int_0^1 \theta \cdot 3\theta^2 d\theta = 3/4$$

$$p(x_2 = 0|x_1 = 1) = \int_0^1 p(x_2 = 0|\theta, x_1 = 1)f(\theta|x_1 = 1) d\theta = \int_0^1 (1 - \theta) \cdot 3\theta^2 d\theta = 1/4$$

(More simply, we could have computed $p(x_2 = 0|x_1 = 1) = 1 - p(x_2 = 1|x_1 = 1) = 1/4$.)

11 From discrete to continuous Bayesian updating

To develop intuition for the transition from discrete to continuous Bayesian updating, we'll walk a familiar road from calculus. Namely we will:

- (i) approximate the continuous range of hypotheses by a finite number.
- (ii) create the discrete updating table for the finite number of hypotheses.
- (iii) consider how the table changes as the number of hypotheses goes to infinity.

In this way, we will see the prior and posterior pmf's converge to the prior and posterior pdf's.

Example 10. To keep things concrete, we will work with the ‘bent’ coin with a flat prior $f(\theta) = 1$ from Example 7. Our goal is to go from discrete to continuous by increasing the number of hypotheses

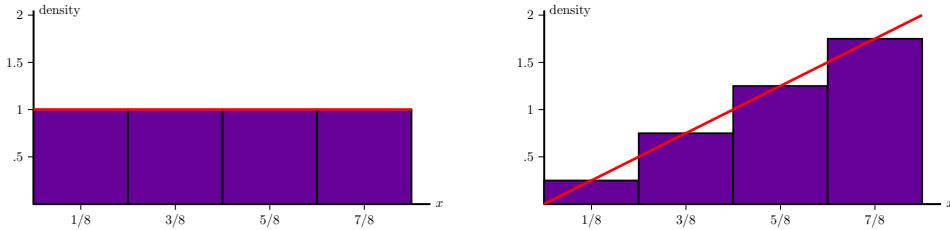
4 hypotheses. We slice $[0, 1]$ into 4 equal intervals: $[0, 1/4]$, $[1/4, 1/2]$, $[1/2, 3/4]$, $[3/4, 1]$. Each slice has width $\Delta\theta = 1/4$. We put our 4 hypotheses θ_i at the centers of the four slices:

$$\theta_1: \theta = 1/8, \quad \theta_2: \theta = 3/8, \quad \theta_3: \theta = 5/8, \quad \theta_4: \theta = 7/8.$$

The flat prior gives each hypothesis a probability of $1/4 = 1 \cdot \Delta\theta$. We have the table:

hypothesis	prior	likelihood	Bayes num.	posterior
$\theta = 1/8$	1/4	1/8	$(1/4) \times (1/8)$	1/16
$\theta = 3/8$	1/4	3/8	$(1/4) \times (3/8)$	3/16
$\theta = 5/8$	1/4	5/8	$(1/4) \times (5/8)$	5/16
$\theta = 7/8$	1/4	7/8	$(1/4) \times (7/8)$	7/16
Total	1	—	$\sum_{i=1}^n \theta_i \Delta\theta$	1

Here are the density histograms of the prior and posterior pmf. The prior and posterior pdfs from Example 7 are superimposed on the histograms in red.

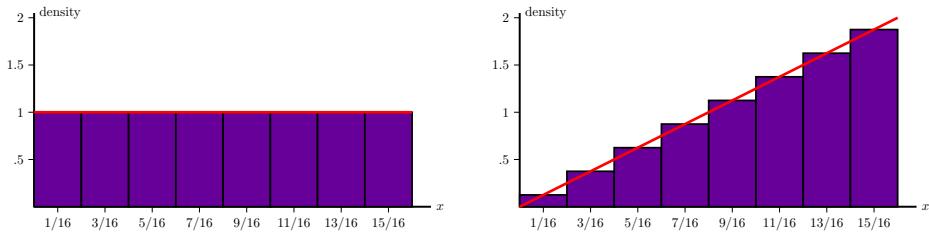


8 hypotheses. Next we slice $[0,1]$ into 8 intervals each of width $\Delta\theta = 1/8$ and use the center of each slice for our 8 hypotheses θ_i .

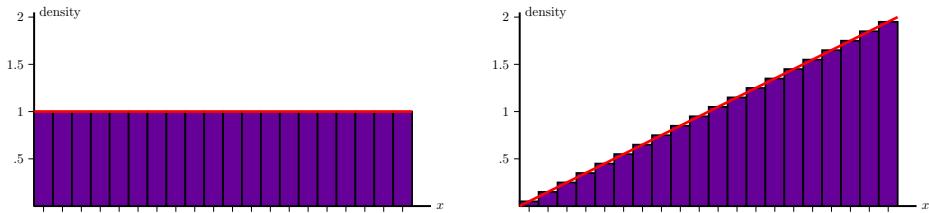
$$\begin{aligned}\theta_1: \quad & \theta = 1/16, \quad \theta_2: \quad \theta = 3/16, \quad \theta_3: \quad \theta = 5/16, \quad \theta_4: \quad \theta = 7/16 \\ \theta_5: \quad & \theta = 9/16, \quad \theta_6: \quad \theta = 11/16, \quad \theta_7: \quad \theta = 13/16, \quad \theta_8: \quad \theta = 15/16\end{aligned}$$

The flat prior gives each hypothesis the probability $1/8 = 1 \cdot \Delta\theta$. Here are the table and density histograms.

hypothesis	prior	likelihood	Bayes num.	posterior
$\theta = 1/16$	$1/8$	$1/16$	$(1/8) \times (1/16)$	$1/64$
$\theta = 3/16$	$1/8$	$3/16$	$(1/8) \times (3/16)$	$3/64$
$\theta = 5/16$	$1/8$	$5/16$	$(1/8) \times (5/16)$	$5/64$
$\theta = 7/16$	$1/8$	$7/16$	$(1/8) \times (7/16)$	$7/64$
$\theta = 9/16$	$1/8$	$9/16$	$(1/8) \times (9/16)$	$9/64$
$\theta = 11/16$	$1/8$	$11/16$	$(1/8) \times (11/16)$	$11/64$
$\theta = 13/16$	$1/8$	$13/16$	$(1/8) \times (13/16)$	$13/64$
$\theta = 15/16$	$1/8$	$15/16$	$(1/8) \times (15/16)$	$15/64$
Total	1	—	$\sum_{i=1}^n \theta_i \Delta\theta$	1



20 hypotheses. Finally we slice $[0,1]$ into 20 pieces. This is essentially identical to the previous two cases. Let's skip right to the density histograms.



Looking at the sequence of plots we see how the prior and posterior density histograms converge to the prior and posterior probability density functions.

Notational conventions
Class 13, 18.05
Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to work with the various notations and terms we use to describe probabilities and likelihood.

2 Introduction

We've introduced a number of different notations for probability, hypotheses and data. We collect them here, to have them in one place.

3 Notation and terminology for data and hypotheses

The problem of labeling data and hypotheses is a tricky one. When we started the course we talked about outcomes, e.g. heads or tails. Then when we introduced random variables we gave outcomes numerical values, e.g. 1 for heads and 0 for tails. This allowed us to do things like compute means and variances. We need to do something similar now. Recall our notational conventions:

- Events are labeled with capital letters, e.g. A, B, C .
- A random variable is capital X and takes values small x .
- The connection between values and events: ' $X = x$ ' is the event that X takes the value x .
- The probability of an event is capital $P(A)$.
- A discrete random variable has a probability mass function small $p(x)$ The connection between P and p is that $P(X = x) = p(x)$.
- A continuous random variable has a probability density function $f(x)$ The connection between P and f is that $P(a \leq X \leq b) = \int_a^b f(x) dx$.
- For a continuous random variable X the probability that X is in an infinitesimal interval of width dx round x is $f(x) dx$.

In the context of Bayesian updating we have similar conventions.

- We use capital letters, especially \mathcal{H} , to indicate a hypothesis, e.g. \mathcal{H} = 'the coin is fair'.

- We use lower case letters, especially θ , to indicate the hypothesized value of a model parameter, e.g. the probability the coin lands heads is $\theta = 0.5$.
- We use upper case letters, especially \mathcal{D} , when talking about data as events. For example, $\mathcal{D} = \text{'the sequence of tosses was HTH'}$.
- We use lower case letters, especially x , when talking about data as values. For example, the sequence of data was $x_1, x_2, x_3 = 1, 0, 1$.
- When the set of hypotheses is discrete we can use the probability of individual hypotheses, e.g. $p(\theta)$. When the set is continuous we need to use the probability for an infinitesimal range of hypotheses, e.g. $f(\theta) d\theta$.

The following table summarizes this for discrete θ and continuous θ . In both cases we are assuming a discrete set of possible outcomes (data) x . Tomorrow we will deal with a continuous set of outcomes.

		Bayes			
	hypothesis	prior	likelihood	numerator	posterior
	\mathcal{H}	$P(\mathcal{H})$	$P(\mathcal{D} \mathcal{H})$	$P(\mathcal{D} \mathcal{H})P(\mathcal{H})$	$P(\mathcal{H} \mathcal{D})$
Discrete θ :	θ	$p(\theta)$	$p(x \theta)$	$p(x \theta)p(\theta)$	$p(\theta x)$
Continuous θ :	θ	$f(\theta) d\theta$	$p(x \theta)$	$p(x \theta)f(\theta) d\theta$	$f(\theta x) d\theta$

Remember the continuous hypothesis θ is really a shorthand for ‘the parameter θ is in an interval of width $d\theta$ around θ ’.

Beta Distributions

Class 14, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be familiar with the 2-parameter family of beta distributions and its normalization.
2. Be able to update a beta prior to a beta posterior in the case of a binomial likelihood.

2 Beta distribution

The [beta distribution](#) $\text{beta}(a, b)$ is a [two-parameter](#) distribution with range $[0, 1]$ and pdf

$$f(\theta) = \frac{(a+b-1)!}{(a-1)!(b-1)!} \theta^{a-1} (1-\theta)^{b-1}$$

We have made an applet so you can explore the shape of the Beta distribution as you vary the parameters:

<http://mathlets.org/mathlets/beta-distribution/>.

As you can see in the applet, the beta distribution may be defined for any real numbers $a > 0$ and $b > 0$. In 18.05 we will stick to integers a and b , but you can get the full story here: http://en.wikipedia.org/wiki/Beta_distribution

In the context of Bayesian updating, a and b are often called [hyperparameters](#) to distinguish them from the unknown parameter θ representing our hypotheses. In a sense, a and b are ‘one level up’ from θ since they parameterize its pdf.

2.1 A simple but important observation!

If a pdf $f(\theta)$ has the form $c\theta^{a-1}(1-\theta)^{b-1}$ then $f(\theta)$ is a $\text{beta}(a, b)$ distribution and the normalizing constant must be

$$c = \frac{(a+b-1)!}{(a-1)!(b-1)!}.$$

This follows because the constant c must normalize the pdf to have total probability 1. There is only one such constant and it is given in the formula for the beta distribution.

A similar observation holds for normal distributions, exponential distributions, and so on.

2.2 Beta priors and posteriors for binomial random variables

Example 1. Suppose we have a bent coin with unknown probability θ of heads. We toss it 12 times and get 8 heads and 4 tails. Starting with a flat prior, show that the posterior pdf is a $\text{beta}(9, 5)$ distribution.

answer: This is nearly identical to examples from the previous class. We'll call the data from all 12 tosses x_1 . In the following table we call the leading constant factor in the posterior column c_2 . Our simple observation will tell us that it has to be the constant factor from the beta pdf.

The data is 8 heads and 4 tails. Since this comes from a binomial(12, θ) distribution, the likelihood $p(x_1|\theta) = \binom{12}{8} \theta^8 (1-\theta)^4$. Thus the Bayesian update table is

		Bayes		
hypothesis	prior	likelihood	numerator	posterior
θ	$1 \cdot d\theta$	$\binom{12}{8} \theta^8 (1-\theta)^4$	$\binom{12}{8} \theta^8 (1-\theta)^4 d\theta$	$c_2 \theta^8 (1-\theta)^4 d\theta$
total	1		$T = \binom{12}{8} \int_0^1 \theta^8 (1-\theta)^4 d\theta$	1

Our simple observation above holds with $a = 9$ and $b = 5$. Therefore the posterior pdf

$$f(\theta|x_1) = c_2 \theta^8 (1-\theta)^4$$

follows a beta(9, 5) distribution and the normalizing constant c_2 must be

$$c_2 = \frac{13!}{8! 4!}.$$

Note: We explicitly included the binomial coefficient $\binom{12}{8}$ in the likelihood. We could just as easily have given it a name, say c_1 and not bothered making its value explicit.

Example 2. Now suppose we toss the same coin again, getting n heads and m tails. Using the posterior pdf of the previous example as our new prior pdf, show that the new posterior pdf is that of a beta($9 + n, 5 + m$) distribution.

answer: It's all in the table. We'll call the data of these $n + m$ additional tosses x_2 . This time we won't make the binomial coefficient explicit. Instead we'll just call it c_3 . Whenever we need a new label we will simply use c with a new subscript.

		Bayes		
hyp.	prior	likelihood	posterior	numerator
θ	$c_2 \theta^8 (1-\theta)^4 d\theta$	$c_3 \theta^n (1-\theta)^m$	$c_2 c_3 \theta^{n+8} (1-\theta)^{m+4} d\theta$	$c_4 \theta^{n+8} (1-\theta)^{m+4} d\theta$
total	1		$T = \int_0^1 c_2 c_3 \theta^{n+8} (1-\theta)^{m+4} d\theta$	1

Again our simple observation holds and therefore the posterior pdf

$$f(\theta|x_1, x_2) = c_4 \theta^{n+8} (1-\theta)^{m+4}$$

follows a beta($n + 9, m + 5$) distribution.

Note: **Flat beta.** The beta(1, 1) distribution is the same as the uniform distribution on $[0, 1]$, which we have also called the flat prior on θ . This follows by plugging $a = 1$ and $b = 1$ into the definition of the beta distribution, giving $f(\theta) = 1$.

Summary: If the probability of heads is θ , the number of heads in $n + m$ tosses follows a binomial($n + m, \theta$) distribution. We have seen that if the prior on θ is a beta distribution then so is the posterior; only the parameters a, b of the beta distribution change! We summarize precisely how they change in a table. We assume the data is n heads in $n + m$ tosses.

hypothesis	data	prior	likelihood	posterior
θ	$x = n$	$\text{beta}(a, b)$	$\text{binomial}(n + m, \theta)$	$\text{beta}(a + n, b + m)$
θ	$x = n$	$c_1 \theta^{a-1} (1 - \theta)^{b-1} d\theta$	$c_2 \theta^n (1 - \theta)^m$	$c_3 \theta^{a+n-1} (1 - \theta)^{b+m-1} d\theta$

2.3 Conjugate priors

In the literature you'll see that the beta distribution is called a [conjugate prior](#) for the binomial distribution. This means that if the likelihood function is binomial, then a beta prior gives a beta posterior. In fact, the beta distribution is a conjugate prior for the Bernoulli and geometric distributions as well.

We will soon see another important example: the normal distribution is its own conjugate prior. In particular, if the likelihood function is normal with known variance, then a normal prior gives a normal posterior.

Conjugate priors are useful because they reduce Bayesian updating to modifying the parameters of the prior distribution (so-called hyperparameters) rather than computing integrals. We saw this for the beta distribution in the last table. For many more examples see:

http://en.wikipedia.org/wiki/Conjugate_prior_distribution

Continuous Data with Continuous Priors

Class 14, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to construct a Bayesian update table for continuous hypotheses and continuous data.
2. Be able to recognize the pdf of a normal distribution and determine its mean and variance.

2 Introduction

We are now ready to do Bayesian updating when both the hypotheses and the data take continuous values. The pattern is the same as what we've done before, so let's first review the previous two cases.

3 Previous cases

1. Discrete hypotheses, discrete data

Notation

- Hypotheses \mathcal{H}
- Data x
- Prior $P(\mathcal{H})$
- Likelihood $p(x | \mathcal{H})$
- Posterior $P(\mathcal{H} | x)$.

Example 1. Suppose we have data x and three possible explanations (hypotheses) for the data that we'll call A, B, C . Suppose also that the data can take two possible values, -1 and 1.

In order to use the data to help estimate the probabilities of the different hypotheses we need a prior pmf and a likelihood table. Assume the prior and likelihoods are given in the following table. (For this example we are only concerned with the formal process of Bayesian updating. So we just made up the prior and likelihoods.)

hypothesis \mathcal{H}	prior $P(\mathcal{H})$	likelihood $p(x \mathcal{H})$	
		$x = -1$	$x = 1$
A	0.1	0.2	0.8
B	0.3	0.5	0.5
C	0.6	0.7	0.3

Prior probabilities Likelihoods

Naturally, each entry in the likelihood table is a likelihood $p(x | \mathcal{H})$. For instance the 0.2 row A and column $x = -1$ is the likelihood $p(x = -1 | A)$.

Question: Suppose we run one trial and obtain the data $x_1 = 1$. Use this to find the posterior probabilities for the hypotheses.

answer: The data picks out one column from the likelihood table which we then use in our Bayesian update table.

hypothesis \mathcal{H}	prior $P(\mathcal{H})$	likelihood $p(x = 1 \mathcal{H})$	Bayes	
			numerator $p(x \mathcal{H})P(\mathcal{H})$	posterior $P(\mathcal{H} x) = \frac{p(x \mathcal{H})P(\mathcal{H})}{p(x)}$
A	0.1	0.8	0.08	0.195
B	0.3	0.5	0.15	0.366
C	0.6	0.3	0.18	0.439
total	1		$p(x) = 0.41$	1

To summarize: the prior probabilities of hypotheses and the likelihoods of data given hypothesis were given; the Bayes numerator is the product of the prior and likelihood; the total probability $p(x)$ is the sum of the probabilities in the Bayes numerator column; and we divide by $p(x)$ to normalize the Bayes numerator.

2. Continuous hypotheses, discrete data

Now suppose that we have data x that can take a discrete set of values and a continuous parameter θ that determines the distribution the data is drawn from.

Notation

- Hypotheses θ
- Data x
- Prior $f(\theta) d\theta$
- Likelihood $p(x | \theta)$
- Posterior $f(\theta | x) d\theta$.

Note: Here we multiplied by $d\theta$ to express the prior and posterior as probabilities. As densities, we have the prior pdf $f(\theta)$ and the posterior pdf $f(\theta | x)$.

Example 2. Assume that $x \sim \text{Binomial}(5, \theta)$. So θ is in the range $[0, 1]$ and the data x can take six possible values, $0, 1, \dots, 5$.

Since there is a continuous range of values we use a pdf to describe the prior on θ . Let's suppose the prior is $f(\theta) = 2\theta$. We can still make a likelihood table, though it only has one row representing an arbitrary hypothesis θ .

hypothesis	likelihood $p(x \theta)$					
	$x = 0$	$x = 1$	$x = 2$	$x = 3$	$x = 4$	$x = 5$
θ	$\binom{5}{0}(1 - \theta)^5$	$\binom{5}{1}\theta(1 - \theta)^4$	$\binom{5}{2}\theta^2(1 - \theta)^3$	$\binom{5}{3}\theta^3(1 - \theta)^2$	$\binom{5}{4}\theta^4(1 - \theta)$	$\binom{5}{5}\theta^5$

Likelihoods

Question: Suppose we run one trial and obtain the data $x_1 = 2$. Use this to find the posterior pdf for the parameter (hypotheses) θ .

answer: As before, the data picks out one column from the likelihood table which we can use in our Bayesian update table. Since we want to work with probabilities we write $f(\theta)d\theta$ and $f(\theta | x_1)d\theta$ for the pdf's.

hypothesis	prior	likelihood	Bayes numerator	posterior
θ	$f(\theta)d\theta$	$p(x = 2 \theta)$	$p(x \theta)f(\theta)d\theta$	$f(\theta x)d\theta = \frac{p(x \theta)f(\theta)d\theta}{p(x)}$
θ	$2\theta d\theta$	$\binom{5}{2}\theta^2(1 - \theta)^3$	$2\binom{5}{2}\theta^3(1 - \theta)^3 d\theta$	$f(\theta x)d\theta = \frac{3! 3!}{7!} \theta^3(1 - \theta)^3 d\theta$
total	1		$p(x) = \int_0^1 2\binom{5}{2}\theta^2(1 - \theta)^3 d\theta = 2\binom{5}{2} \frac{3! 3!}{7!}$	1

To summarize: the prior probabilities of hypotheses and the likelihoods of data given hypothesis were given; the Bayes numerator is the product of the prior and likelihood; the total probability $p(x)$ is the integral of the probabilities in the Bayes numerator column; and we divide by $p(x)$ to normalize the Bayes numerator.

4 Continuous hypotheses and continuous data

When both data and hypotheses are continuous, the only change to the previous example is that the likelihood function uses a pdf $f(x | \theta)$ instead of a pmf $p(x | \theta)$. The general shape of the Bayesian update table is the same.

Notation

- Hypotheses θ
- Data x
- Prior $f(\theta)d\theta$

- Likelihood $f(x | \theta) dx$
- Posterior $f(\theta | x) d\theta$.

Simplifying the notation. In the previous cases we included $d\theta$ so that we were working with probabilities instead of densities. When both data and hypotheses are continuous we will need both $d\theta$ and dx . This makes things conceptually simpler, but notationally cumbersome. To simplify the notation we will allow ourselves to dx in our tables. This is fine because the data x is a fixed. We keep the $d\theta$ because the hypothesis θ is allowed to vary.

For comparison, we first show the general table in simplified notation followed immediately afterward by the table showing the infinitesimals.

hypoth.	prior	likelihood	Bayes numerator	posterior
θ	$f(\theta) d\theta$	$f(x \theta)$	$f(x \theta)f(\theta) d\theta$	$f(\theta x) = \frac{f(x \theta)f(\theta) d\theta}{f(x)}$
total	1		$f(x) = \int f(x \theta)f(\theta) d\theta$	1

Bayesian update table without dx

hypoth.	prior	likelihood	Bayes numerator	posterior
θ	$f(\theta) d\theta$	$f(x \theta) dx$	$f(x \theta)f(\theta) d\theta dx$	$f(\theta x) d\theta = \frac{f(x \theta)f(\theta) d\theta dx}{f(x) dx} = \frac{f(x \theta)f(\theta) d\theta}{f(x)}$
total	1		$f(x) dx = (\int f(x \theta)f(\theta) d\theta) dx$	1

Bayesian update table with $d\theta$ and dx

To summarize: the prior probabilities of hypotheses and the likelihoods of data given hypothesis were given; the Bayes numerator is the product of the prior and likelihood; the total probability $f(x) dx$ is the integral of the probabilities in the Bayes numerator column; we divide by $f(x) dx$ to normalize the Bayes numerator.

5 Normal hypothesis, normal data

A standard example of continuous hypotheses and continuous data assumes that both the data and prior follow normal distributions. The following example assumes that the variance of the data is known.

Example 3. Suppose we have data $x = 5$ which was drawn from a normal distribution

with unknown mean θ and standard deviation 1.

$$x \sim N(\theta, 1)$$

Suppose further that our prior distribution for θ is $\theta \sim N(2, 1)$.

Let x represent an arbitrary data value.

- (a) Make a Bayesian table with prior, likelihood, and Bayes numerator.
- (b) Show that the posterior distribution for θ is normal as well.
- (c) Find the mean and variance of the posterior distribution.

answer: As we did with the tables above, a good compromise on the notation is to include $d\theta$ but not dx . The reason for this is that the total probability is computed by integrating over θ and the $d\theta$ reminds of us that.

Our prior pdf is

$$f(\theta) = \frac{1}{\sqrt{2\pi}} e^{-(\theta-2)^2/2}.$$

The likelihood function is

$$f(x = 5 | \theta) = \frac{1}{\sqrt{2\pi}} e^{-(5-\theta)^2/2}.$$

We know we are going to multiply the prior and the likelihood, so we carry out that algebra first. In the very last step we simplify the constant factor into one constant we call c_1 .

$$\begin{aligned} \text{prior} \cdot \text{likelihood} &= \frac{1}{\sqrt{2\pi}} e^{-(\theta-2)^2/2} \cdot \frac{1}{\sqrt{2\pi}} e^{-(5-\theta)^2/2} \\ &= \frac{1}{2\pi} e^{-(2\theta^2 - 14\theta + 29)/2} \\ &= \frac{1}{2\pi} e^{-(\theta^2 - 7\theta + 29/2)} \quad (\text{complete the square}) \\ &= \frac{1}{2\pi} e^{-((\theta - 7/2)^2 + 9/4)} \\ &= \frac{e^{-9/4}}{2\pi} e^{-(\theta - 7/2)^2} \\ &= c_1 e^{-(\theta - 7/2)^2} \end{aligned}$$

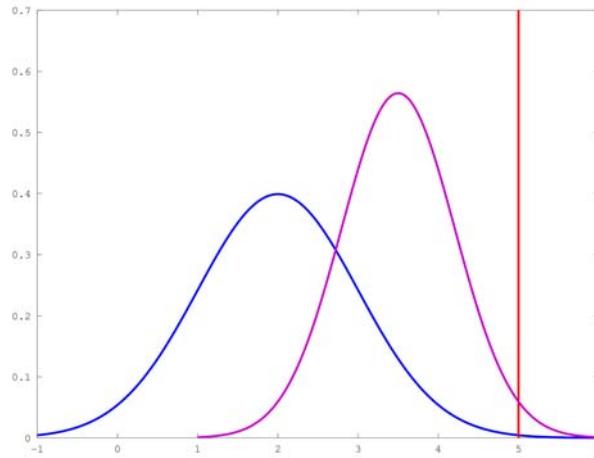
In the last step we replaced the complicated constant factor by the simpler expression c_1 .

hypothesis	prior	likelihood	Bayes numerator	posterior
θ	$f(\theta) d\theta$	$f(x = 5 \theta)$	$f(x = 5 \theta) f(\theta) d\theta$	$\frac{f(x = 5 \theta) f(\theta) d\theta}{f(x = 5)}$
θ	$\frac{1}{\sqrt{2\pi}} e^{-(\theta-2)^2/2} d\theta$	$\frac{1}{\sqrt{2\pi}} e^{-(5-\theta)^2/2}$	$c_1 e^{-(\theta - 7/2)^2}$	$c_2 e^{-(\theta - 7/2)^2}$
total	1	$f(x = 5) = \int f(x = 5 \theta) f(\theta) d\theta$		1

We can see by the form of the posterior pdf that it is a normal distribution. Because the exponential for a normal distribution is $e^{-(\theta-\mu)^2/2\sigma^2}$ we have mean $\mu = 7/2$ and $2\sigma^2 = 1$, so variance $\sigma^2 = 1/2$.

We don't need to bother computing the total probability; it is just used for normalization and we already know the normalization constant $\frac{1}{\sigma\sqrt{2\pi}}$ for a normal distribution.

Here is the graph of the prior and the posterior pdf's for this example. Note how the data 'pulls' the prior towards the data.



prior = blue; posterior = purple; data = red

Now we'll repeat the previous example for general x . When reading this if you mentally substitute 5 for x you will understand the algebra.

Example 4. Suppose our data x is drawn from a normal distribution with unknown mean θ and standard deviation 1.

$$x \sim N(\theta, 1)$$

answer: As before, we show the algebra used to simplify the Bayes numerator: The prior pdf and likelihood function are

$$f(\theta) = \frac{1}{\sqrt{2\pi}} e^{-(\theta-2)^2/2} \quad f(x | \theta) = \frac{1}{\sqrt{2\pi}} e^{-(x-\theta)^2/2}.$$

The Bayes numerator is the product of the prior and the likelihood:

$$\begin{aligned} \text{prior} \cdot \text{likelihood} &= \frac{1}{\sqrt{2\pi}} e^{-(\theta-2)^2/2} \cdot \frac{1}{\sqrt{2\pi}} e^{-(x-\theta)^2/2} \\ &= \frac{1}{2\pi} e^{-(2\theta^2 - (4+2x)\theta + 4+x^2)/2} \\ &= \frac{1}{2\pi} e^{-(\theta^2 - (2+x)\theta + (4+x^2)/2)} \quad (\text{complete the square}) \\ &= \frac{1}{2\pi} e^{-((\theta - (1+x/2))^2 - (1+x/2)^2 + (4+x^2)/2)} \\ &= c_1 e^{-(\theta - (1+x/2))^2} \end{aligned}$$

Just as in the previous example, in the last step we replaced all the constants, including the exponentials that just involve x , by the simple constant c_1 .

Now the Bayesian update table becomes

hypothesis	prior	likelihood	Bayes numerator	posterior $f(\theta x) d\theta$
θ	$f(\theta) d\theta$	$f(x \theta)$	$f(x \theta)f(\theta) d\theta$	$\frac{f(x \theta)f(\theta) d\theta}{f(x)}$
θ	$\frac{1}{\sqrt{2\pi}}e^{-(\theta-2)^2/2} d\theta$	$\frac{1}{\sqrt{2\pi}}e^{-(x-\theta)^2/2}$	$c_1 e^{-(\theta-(1+x/2))^2}$	$c_2 e^{-(\theta-(1+x/2))^2}$
total	1		$f(x) = f(x \theta)f(\theta) d\theta$	1

As in the previous example we can see by the form of the posterior that it must be a normal distribution with mean $1 + x/2$ and variance $1/2$. (Compare this with the case $x = 5$ in the previous example.)

6 Predictive probabilities

Since the data x is continuous it has prior and posterior predictive pdfs. The [prior predictive pdf](#) is the total probability density computed at the bottom of the Bayes numerator column:

$$f(x) = \int f(x|\theta)f(\theta) d\theta,$$

where the integral is computed over the entire range of θ .

The [posterior predictive pdf](#) has the same form as the prior predictive pdf, except it uses the posterior probabilities for θ :

$$f(x_2|x_1) = \int f(x_2|\theta, x_1)f(\theta|x_1) d\theta,$$

As usual, we usually assume x_1 and x_2 are [conditionally independent](#). That is,

$$f(x_2|\theta, x_1) = f(x_2|\theta).$$

In this case the formula for the posterior predictive pdf is a little simpler:

$$f(x_2|x_1) = \int f(x_2|\theta)f(\theta|x_1) d\theta,$$

Conjugate priors: Beta and normal
Class 15, 18.05
Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Understand the benefits of conjugate priors.
2. Be able to update a beta prior given a Bernoulli, binomial, or geometric likelihood.
3. Understand and be able to use the formula for updating a normal prior given a normal likelihood with known variance.

2 Introduction and definition

In this reading, we will elaborate on the notion of a conjugate prior for a likelihood function. With a conjugate prior the posterior is of the same type, e.g. for binomial likelihood the beta prior becomes a beta posterior. Conjugate priors are useful because they reduce Bayesian updating to modifying the parameters of the prior distribution (so-called hyperparameters) rather than computing integrals.

Our focus in 18.05 will be on two important examples of conjugate priors: beta and normal. For a far more comprehensive list, see the tables herein:

http://en.wikipedia.org/wiki/Conjugate_prior_distribution

We now give a definition of conjugate prior. It is best understood through the examples in the subsequent sections.

Definition. Suppose we have data with likelihood function $f(x|\theta)$ depending on a hypothesized parameter. Also suppose the prior distribution for θ is one of a family of parametrized distributions. If the posterior distribution for θ is in this family then we say the prior is a [conjugate prior](#) for the likelihood.

3 Beta distribution

In this section, we will show that the beta distribution is a conjugate prior for binomial, Bernoulli, and geometric likelihoods.

3.1 Binomial likelihood

We saw last time that the [beta distribution is a conjugate prior for the binomial distribution](#). This means that if the likelihood function is binomial and the prior distribution is beta then the posterior is also beta.

More specifically, suppose that the likelihood follows a binomial(N, θ) distribution where N is known and θ is the (unknown) parameter of interest. We also have that the data x from one trial is an integer between 0 and N . Then for a beta prior we have the following table:

hypothesis	data	prior	likelihood	posterior
θ	x	$\text{beta}(a, b)$	$\text{binomial}(N, \theta)$	$\text{beta}(a + x, b + N - x)$
θ	x	$c_1 \theta^{a-1} (1-\theta)^{b-1}$	$c_2 \theta^x (1-\theta)^{N-x}$	$c_3 \theta^{a+x-1} (1-\theta)^{b+N-x-1}$

The table is simplified by writing the normalizing coefficient as c_1 , c_2 and c_3 respectively. If needed, we can recover the values of the c_1 and c_2 by recalling (or looking up) the normalizations of the beta and binomial distributions.

$$c_1 = \frac{(a+b-1)!}{(a-1)! (b-1)!} \quad c_2 = \binom{N}{x} = \frac{N!}{x! (N-x)!} \quad c_3 = \frac{(a+b+N-1)!}{(a+x-1)! (b+N-x-1)!}$$

3.2 Bernoulli likelihood

The [beta distribution is a conjugate prior for the Bernoulli distribution](#). This is actually a special case of the binomial distribution, since $\text{Bernoulli}(\theta)$ is the same as $\text{binomial}(1, \theta)$. We do it separately because it is slightly simpler and of special importance. In the table below, we show the updates corresponding to success ($x = 1$) and failure ($x = 0$) on separate rows.

hypothesis	data	prior	likelihood	posterior
θ	x	$\text{beta}(a, b)$	$\text{Bernoulli}(\theta)$	$\text{beta}(a+1, b)$ or $\text{beta}(a, b+1)$
θ	$x = 1$	$c_1 \theta^{a-1} (1-\theta)^{b-1}$	θ	$c_3 \theta^a (1-\theta)^{b-1}$
θ	$x = 0$	$c_1 \theta^{a-1} (1-\theta)^{b-1}$	$1-\theta$	$c_3 \theta^{a-1} (1-\theta)^b$

The constants c_1 and c_3 have the same formulas as in the previous (binomial likelihood case) with $N = 1$.

3.3 Geometric likelihood

Recall that the $\text{geometric}(\theta)$ distribution describes the probability of x successes before the first failure, where the probability of success on any single independent trial is θ . The corresponding pmf is given by $p(x) = \theta^x (1-\theta)$.

Now suppose that we have a data point x , and our hypothesis θ is that x is drawn from a $\text{geometric}(\theta)$ distribution. From the table we see that the [beta distribution is a conjugate prior for a geometric likelihood](#) as well:

hypothesis	data	prior	likelihood	posterior
θ	x	$\text{beta}(a, b)$	$\text{geometric}(\theta)$	$\text{beta}(a+x, b+1)$
θ	x	$c_1 \theta^{a-1} (1-\theta)^{b-1}$	$\theta^x (1-\theta)$	$c_3 \theta^{a+x-1} (1-\theta)^b$

At first it may seem strange that the beta distribution is a conjugate prior for both the binomial and geometric distributions. The key reason is that the binomial and geometric likelihoods are proportional as functions of θ . Let's illustrate this in a concrete example.

Example 1. While traveling through the Mushroom Kingdom, Mario and Luigi find some rather unusual coins. They agree on a prior of $f(\theta) \sim \text{beta}(5, 5)$ for the probability of heads,

though they disagree on what experiment to run to investigate θ further.

a) Mario decides to flip a coin 5 times. He gets four heads in five flips.

b) Luigi decides to flip a coin until the first tails. He gets four heads before the first tail.

Show that Mario and Luigi will arrive at the same posterior on θ , and calculate this posterior.

answer: We will show that both Mario and Luigi find the posterior pdf for θ is a beta(9, 6) distribution.

Mario's table

hypothesis	data	prior	likelihood	posterior
θ	$x = 4$	beta(5, 5)	binomial(5, θ)	???
θ	$x = 4$	$c_1 \theta^4 (1 - \theta)^4$	$\binom{5}{4} \theta^4 (1 - \theta)$	$c_3 \theta^8 (1 - \theta)^5$

Luigi's table

hypothesis	data	prior	likelihood	posterior
θ	$x = 4$	beta(5, 5)	geometric(θ)	???
θ	$x = 4$	$c_1 \theta^4 (1 - \theta)^4$	$\theta^4 (1 - \theta)$	$c_3 \theta^8 (1 - \theta)^5$

Since both Mario and Luigi's posterior has the form of a beta(9, 6) distribution that's what they both must be. The normalizing factor is the same in both cases because it's determined by requiring the total probability to be 1.

4 Normal begets normal

We now turn to another important example: [the normal distribution is its own conjugate prior](#). In particular, if the likelihood function is normal with known variance, then a normal prior gives a normal posterior. Now both the hypotheses and the data are continuous.

Suppose we have a measurement $x \sim N(\theta, \sigma^2)$ where the variance σ^2 is known. That is, the mean θ is our unknown parameter of interest and we are given that the likelihood comes from a normal distribution with variance σ^2 . If we choose a normal prior pdf

$$f(\theta) \sim N(\mu_{\text{prior}}, \sigma_{\text{prior}}^2)$$

then the posterior pdf is also normal: $f(\theta|x) \sim N(\mu_{\text{post}}, \sigma_{\text{post}}^2)$ where

$$\frac{\mu_{\text{post}}}{\sigma_{\text{post}}^2} = \frac{\mu_{\text{prior}}}{\sigma_{\text{prior}}^2} + \frac{x}{\sigma^2}, \quad \frac{1}{\sigma_{\text{post}}^2} = \frac{1}{\sigma_{\text{prior}}^2} + \frac{1}{\sigma^2} \quad (1)$$

The following form of these formulas is easier to read and shows that μ_{post} is a weighted average between μ_{prior} and the data x .

$$a = \frac{1}{\sigma_{\text{prior}}^2}, \quad b = \frac{1}{\sigma^2}, \quad \mu_{\text{post}} = \frac{a\mu_{\text{prior}} + bx}{a + b}, \quad \sigma_{\text{post}}^2 = \frac{1}{a + b}. \quad (2)$$

With these formulas in mind, we can express the update via the table:

hypothesis	data	prior	likelihood	posterior
θ	x	$f(\theta) \sim N(\mu_{\text{prior}}, \sigma_{\text{prior}}^2)$	$f(x \theta) \sim N(\theta, \sigma^2)$	$f(\theta x) \sim N(\mu_{\text{post}}, \sigma_{\text{post}}^2)$
θ	x	$c_1 \exp\left(\frac{-(\theta - \mu_{\text{prior}})^2}{2\sigma_{\text{prior}}^2}\right)$	$c_2 \exp\left(\frac{-(x - \theta)^2}{2\sigma^2}\right)$	$c_3 \exp\left(\frac{-(\theta - \mu_{\text{post}})^2}{2\sigma_{\text{post}}^2}\right)$

We leave the proof of the general formulas to the problem set. It is an involved algebraic manipulation which is essentially the same as the following numerical example.

Example 2. Suppose we have prior $\theta \sim N(4, 8)$, and likelihood function likelihood $x \sim N(\theta, 5)$. Suppose also that we have one measurement $x_1 = 3$. Show the posterior distribution is normal.

answer: We will show this by grinding through the algebra which involves completing the square.

$$\text{prior: } f(\theta) = c_1 e^{-(\theta-4)^2/16}; \quad \text{likelihood: } f(x_1|\theta) = c_2 e^{-(x_1-\theta)^2/10} = c_2 e^{-(3-\theta)^2/10}$$

We multiply the prior and likelihood to get the posterior:

$$\begin{aligned} f(\theta|x_1) &= c_3 e^{-(\theta-4)^2/16} e^{-(3-\theta)^2/10} \\ &= c_3 \exp\left(-\frac{(\theta-4)^2}{16} - \frac{(3-\theta)^2}{10}\right) \end{aligned}$$

We complete the square in the exponent

$$\begin{aligned} -\frac{(\theta-4)^2}{16} - \frac{(3-\theta)^2}{10} &= -\frac{5(\theta-4)^2 + 8(3-\theta)^2}{80} \\ &= -\frac{13\theta^2 - 88\theta + 152}{80} \\ &= -\frac{\theta^2 - \frac{88}{13}\theta + \frac{152}{13}}{80/13} \\ &= -\frac{(\theta - 44/13)^2 + 152/13 - (44/13)^2}{80/13}. \end{aligned}$$

Therefore the posterior is

$$f(\theta|x_1) = c_3 e^{-\frac{(\theta-44/13)^2+152/13-(44/13)^2}{80/13}} = c_4 e^{-\frac{(\theta-44/13)^2}{80/13}}.$$

This has the form of the pdf for $N(44/13, 40/13)$. QED

For practice we check this against the formulas (2).

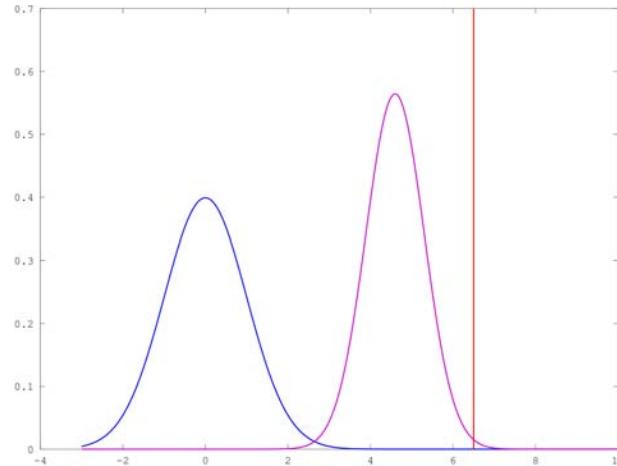
$$\mu_{\text{prior}} = 4, \quad \sigma_{\text{prior}}^2 = 8, \quad \sigma^2 = 5 \Rightarrow a = \frac{1}{8}, \quad b = \frac{1}{5}.$$

Therefore

$$\begin{aligned} \mu_{\text{post}} &= \frac{a\mu_{\text{prior}} + bx}{a+b} = \frac{44}{13} = 3.38 \\ \sigma_{\text{post}}^2 &= \frac{1}{a+b} = \frac{40}{13} = 3.08. \end{aligned}$$

Example 3. Suppose that we know the data $x \sim N(\theta, 1)$ and we have prior $N(0, 1)$. We get one data value $x = 6.5$. Describe the changes to the pdf for θ in updating from the prior to the posterior.

answer: Here is a graph of the prior pdf with the data point marked by a red line.



Prior in blue, posterior in magenta, data in red

The posterior mean will be a weighted average of the prior mean and the data. So the peak of the posterior pdf will be between the peak of the prior and the red line. A little algebra with the formula shows

$$\sigma_{\text{post}}^2 = \frac{1}{1/\sigma_{\text{prior}}^2 + 1/\sigma^2} = \sigma_{\text{prior}}^2 \cdot \frac{\sigma^2}{\sigma_{\text{prior}}^2 + \sigma^2} < \sigma_{\text{prior}}^2$$

That is the posterior has smaller variance than the prior, i.e. data makes us more certain about where in its range θ lies.

4.1 More than one data point

Example 4. Suppose we have data x_1, x_2, x_3 . Use the formulas (1) to update sequentially.

answer: Let's label the prior mean and variance as μ_0 and σ_0^2 . The updated means and variances will be μ_i and σ_i^2 . In sequence we have

$$\begin{aligned} \frac{1}{\sigma_1^2} &= \frac{1}{\sigma_0^2} + \frac{1}{\sigma^2}; & \frac{\mu_1}{\sigma_1^2} &= \frac{\mu_0}{\sigma_0^2} + \frac{x_1}{\sigma^2} \\ \frac{1}{\sigma_2^2} &= \frac{1}{\sigma_1^2} + \frac{1}{\sigma^2} = \frac{1}{\sigma_0^2} + \frac{2}{\sigma^2}; & \frac{\mu_2}{\sigma_2^2} &= \frac{\mu_1}{\sigma_1^2} + \frac{x_2}{\sigma^2} = \frac{\mu_0}{\sigma_0^2} + \frac{x_1 + x_2}{\sigma^2} \\ \frac{1}{\sigma_3^2} &= \frac{1}{\sigma_2^2} + \frac{1}{\sigma^2} = \frac{1}{\sigma_0^2} + \frac{3}{\sigma^2}; & \frac{\mu_3}{\sigma_3^2} &= \frac{\mu_2}{\sigma_2^2} + \frac{x_3}{\sigma^2} = \frac{\mu_0}{\sigma_0^2} + \frac{x_1 + x_2 + x_3}{\sigma^2} \end{aligned}$$

The example generalizes to n data values x_1, \dots, x_n :

Normal-normal update formulas for n data points

$$\frac{\mu_{\text{post}}}{\sigma_{\text{post}}^2} = \frac{\mu_{\text{prior}}}{\sigma_{\text{prior}}^2} + \frac{n\bar{x}}{\sigma^2}, \quad \frac{1}{\sigma_{\text{post}}^2} = \frac{1}{\sigma_{\text{prior}}^2} + \frac{n}{\sigma^2}, \quad \bar{x} = \frac{x_1 + \dots + x_n}{n}. \quad (3)$$

Again we give the easier to read form, showing μ_{post} is a weighted average of μ_{prior} and the sample average \bar{x} :

$$a = \frac{1}{\sigma_{\text{prior}}^2}, \quad b = \frac{n}{\sigma^2}, \quad \mu_{\text{post}} = \frac{a\mu_{\text{prior}} + b\bar{x}}{a + b}, \quad \sigma_{\text{post}}^2 = \frac{1}{a + b}. \quad (4)$$

Interpretation: μ_{post} is a weighted average of μ_{prior} and \bar{x} . If the number of data points is large then the weight b is large and \bar{x} will have a strong influence on the posterior. If σ_{prior}^2 is small then the weight a is large and μ_{prior} will have a strong influence on the posterior. To summarize:

1. Lots of data has a big influence on the posterior.
2. High certainty (low variance) in the prior has a big influence on the posterior.

The actual posterior is a balance of these two influences.

Choosing priors
Class 15, 18.05
Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Learn that the choice of prior affects the posterior.
2. See that too rigid a prior can make it difficult to learn from the data.
3. See that more data lessens the dependence of the posterior on the prior.
4. Be able to make a reasonable choice of prior, based on prior understanding of the system under consideration.

2 Introduction

Up to now we have always been handed a prior pdf. In this case, statistical inference from data is essentially an application of Bayes' theorem. When the prior is known there is no controversy on how to proceed. The art of statistics starts when the prior is not known with certainty. There are two main schools on how to proceed in this case: **Bayesian** and **frequentist**. For now we are following the Bayesian approach. Starting next week we will learn the frequentist approach.

Recall that given data D and a hypothesis H we used Bayes' theorem to write

$$P(H|D) = \frac{P(D|H) \cdot P(H)}{P(D)}$$

posterior \propto likelihood · prior.

Bayesian: Bayesians make inferences using the posterior $P(H|D)$, and therefore always need a prior $P(H)$. If a prior is not known with certainty the Bayesian must try to make a reasonable choice. There are many ways to do this and reasonable people might make different choices. In general it is good practice to justify your choices and to explore a range of priors to see if they all point to the same conclusion.

Frequentist: Very briefly, frequentists do not try to create a prior. Instead, they make inferences using the likelihood $P(D|H)$.

We will compare the two approaches in detail once we have more experience with each. For now we simply list two benefits of the Bayesian approach.

1. The posterior probability $P(H|D)$ for the hypothesis given the evidence is usually exactly what we'd like to know. The Bayesian can say something like 'the parameter of interest has probability 0.95 of being between 0.49 and 0.51.'
2. The assumptions that go into choosing the prior can be clearly spelled out.

More good data: It is always the case that **more good data** allows for stronger conclusions and lessens the influence of the prior. The emphasis should be as much on good data (quality) as on more data (quantity).

3 Example: Dice

Suppose we have a drawer full of dice, each of which has either 4, 6, 8, 12, or 20 sides. This time, we do not know how many of each type are in the drawer. A die is picked at random from the drawer and rolled 5 times. The results in order are 4, 2, 4, 7, and 5.

3.1 Uniform prior

Suppose we have no idea what the distribution of dice in the drawer might be. In this case it's reasonable to use a flat prior. Here is the update table for the posterior probabilities that result from updating after each roll. In order to fit all the columns, we leave out the unnormalized posteriors.

hyp.	prior	lik ₁	post ₁	lik ₂	post ₂	lik ₃	post ₃	lik ₄	post ₄	lik ₅	post ₅
H_4	1/5	1/4	0.370	1/4	0.542	1/4	0.682	0	0.000	0	0.000
H_6	1/5	1/6	0.247	1/6	0.241	1/6	0.202	0	0.000	1/6	0.000
H_8	1/5	1/8	0.185	1/8	0.135	1/8	0.085	1/8	0.818	1/8	0.876
H_{12}	1/5	1/12	0.123	1/12	0.060	1/12	0.025	1/12	0.161	1/12	0.115
H_{20}	1/5	1/20	0.074	1/20	0.022	1/20	0.005	1/20	0.021	1/20	0.009

This should look familiar. Given the data the final posterior is heavily weighted towards hypothesis H_8 that the 8-sided die was picked.

3.2 Other priors

To see how much the above posterior depended on our choice of prior, let's try some other priors. Suppose we have reason to believe that there are ten times as many 20-sided dice in the drawer as there are each of the other types. The table becomes:

hyp.	prior	lik ₁	post ₁	lik ₂	post ₂	lik ₃	post ₃	lik ₄	post ₄	lik ₅	post ₅
H_4	0.071	1/4	0.222	1/4	0.453	1/4	0.650	0	0.000	0	0.000
H_6	0.071	1/6	0.148	1/6	0.202	1/6	0.193	0	0.000	1/6	0.000
H_8	0.071	1/8	0.111	1/8	0.113	1/8	0.081	1/8	0.688	1/8	0.810
H_{12}	0.071	1/12	0.074	1/12	0.050	1/12	0.024	1/12	0.136	1/12	0.107
H_{20}	0.714	1/20	0.444	1/20	0.181	1/20	0.052	1/20	0.176	1/20	0.083

Even here the final posterior is heavily weighted to the hypothesis H_8 .

What if the 20-sided die is 100 times more likely than each of the others?

hyp.	prior	lik ₁	post ₁	lik ₂	post ₂	lik ₃	post ₃	lik ₄	post ₄	lik ₅	post ₅
H_4	0.0096	1/4	0.044	1/4	0.172	1/4	0.443	0	0.000	0	0.000
H_6	0.0096	1/6	0.030	1/6	0.077	1/6	0.131	0	0.000	1/6	0.000
H_8	0.0096	1/8	0.022	1/8	0.043	1/8	0.055	1/8	0.266	1/8	0.464
H_{12}	0.0096	1/12	0.015	1/12	0.019	1/12	0.016	1/12	0.053	1/12	0.061
H_{20}	0.9615	1/20	0.889	1/20	0.689	1/20	0.354	1/20	0.681	1/20	0.475

With such a strong prior belief in the 20-sided die, the final posterior gives a lot of weight to the theory that the data arose from a 20-sided die, even though it extremely unlikely the

20-sided die would produce a maximum of 7 in 5 roles. The posterior now gives roughly even odds that an 8-sided die versus a 20-sided die was picked.

3.3 Rigid priors

Mild cognitive dissonance. Too rigid a prior belief can overwhelm any amount of data. Suppose I've got it in my head that the die has to be 20-sided. So I set my prior to $P(H_{20}) = 1$ with the other 4 hypotheses having probability 0. Look what happens in the update table.

hyp.	prior	lik ₁	post ₁	lik ₂	post ₂	lik ₃	post ₃	lik ₄	post ₄	lik ₅	post ₅
H_4	0	1/4	0	1/4	0	1/4	0	0	0	0	0
H_6	0	1/6	0	1/6	0	1/6	0	0	0	1/6	0
H_8	0	1/8	0	1/8	0	1/8	0	1/8	0	1/8	0
H_{12}	0	1/12	0	1/12	0	1/12	0	1/12	0	1/12	0
H_{20}	1	1/20	1	1/20	1	1/20	1	1/20	1	1/20	1

No matter what the data, a hypothesis with prior probability 0 will have posterior probability 0. In this case I'll never get away from the hypothesis H_{20} , although I might experience some mild cognitive dissonance.

Severe cognitive dissonance. Rigid priors can also lead to absurdities. Suppose I now have it in my head that the die must be 4-sided. So I set $P(H_4) = 1$ and the other prior probabilities to 0. With the given data on the fourth roll I reach an impasse. A roll of 7 can't possibly come from a 4-sided die. Yet this is the only hypothesis I'll allow. My unnormalized posterior is a column of all zeros which cannot be normalized.

hyp.	prior	lik ₁	post ₁	lik ₂	post ₂	lik ₃	post ₃	lik ₄	unnorm.	post ₄	post ₄
H_4	1	1/4	1	1/4	1	1/4	1	0	0	???	???
H_6	0	1/6	0	1/6	0	1/6	0	0	0	???	???
H_8	0	1/8	0	1/8	0	1/8	0	1/8	0	???	???
H_{12}	0	1/12	0	1/12	0	1/12	0	1/12	0	???	???
H_{20}	0	1/20	0	1/20	0	1/20	0	1/20	0	???	???

I must adjust my belief about what is possible or, more likely, I'll suspect you of accidentally or deliberately messing up the data.

4 Example: Malaria

Here is a real example adapted from *Statistics, A Bayesian Perspective* by Donald Berry:

By the 1950's scientists had begun to formulate the hypothesis that carriers of the sickle-cell gene were more resistant to malaria than noncarriers. There was a fair amount of circumstantial evidence for this hypothesis. It also helped explain the persistence of an otherwise deleterious gene in the population. In one experiment scientists injected 30 African volunteers with malaria. Fifteen of the volunteers carried one copy of the sickle-cell gene and the other 15 were noncarriers. Fourteen out of 15 noncarriers developed malaria while only 2

out of 15 carriers did. Does this small sample support the hypothesis that the sickle-cell gene protects against malaria?

Let S represent a carrier of the sickle-cell gene and N represent a non-carrier. Let $D+$ indicate developing malaria and $D-$ indicate not developing malaria. The data can be put in a table.

	$D+$	$D-$	
S	2	13	15
N	14	1	15
	16	14	30

Before analysing the data we should say a few words about the experiment and experimental design. First, it is clearly unethical: to gain some information they infected 16 people with malaria. We also need to worry about bias. How did they choose the test subjects. Is it possible the noncarriers were weaker and thus more susceptible to malaria than the carriers? Berry points out that it is reasonable to assume that an injection is similar to a mosquito bite, but it is not guaranteed. This last point means that if the experiment shows a relation between sickle-cell and protection against injected malaria, we need to consider the hypothesis that the protection from mosquito transmitted malaria is weaker or non-existent. Finally, we will frame our hypothesis as 'sickle-cell protects against malaria', but really all we can hope to say from a study like this is that 'sickle-cell is correlated with protection against malaria'.

Model. For our model let θ_S be the probability that an injected carrier S develops malaria and likewise let θ_N be the probability that an injected noncarrier N develops malaria. We assume independence between all the experimental subjects. With this model, the likelihood is a function of both θ_S and θ_N :

$$P(\text{data}|\theta_S, \theta_N) = c \theta_S^2 (1 - \theta_S)^{13} \theta_N^{14} (1 - \theta_N).$$

As usual we leave the constant factor c as a letter. (It is a product of two binomial coefficients: $c = \binom{15}{2} \binom{15}{14}$.)

Hypotheses. Each hypothesis consists of a pair (θ_N, θ_S) . To keep things simple we will only consider a finite number of values for these probabilities. We could easily consider many more values or even a continuous range of hypotheses. Assume θ_S and θ_N are each one of 0, 0.2, 0.4, 0.6, 0.8, 1. This leads to two-dimensional tables.

First is a table of hypotheses. The color coding indicates the following:

1. Light orange squares along the diagonal are where $\theta_S = \theta_N$, i.e. sickle-cell makes no difference one way or the other.
2. Pink and red squares above the diagonal are where $\theta_N > \theta_S$, i.e. sickle-cell provides some protection against malaria.
3. In the red squares $\theta_N - \theta_S \geq 0.6$, i.e. sickle-cell provides a lot of protection.
4. White squares below diagonal are where $\theta_S > \theta_N$, i.e. sickle-cell actually increases the probability of developing malaria.

$\theta_N \setminus \theta_S$	0	0.2	0.4	0.6	0.8	1
1	(0,1)	(.2,1)	(.4,1)	(.6,1)	(.8,1)	(1,1)
0.8	(0,.8)	(.2,.8)	(.4,.8)	(.6,.8)	(.8,.8)	(1,.8)
0.6	(0,.6)	(.2,.6)	(.4,.6)	(.6,.6)	(.8,.6)	(1,.6)
0.4	(0,.4)	(.2,.4)	(.4,.4)	(.6,.4)	(.8,.4)	(1,.4)
0.2	(0,.2)	(.2,.2)	(.4,.2)	(.6,.2)	(.8,.2)	(1,.2)
0	(0,0)	(.2,0)	(.4,0)	(.6,0)	(.8,0)	(1,0)

Hypotheses on level of protection due to S :
red = strong; pink = some; orange = none; white = negative.

Next is the table of likelihoods. (Actually we've taken advantage of our indifference to scale and scaled all the likelihoods by $100000/c$ to make the table more presentable.) Notice that, to the precision of the table, many of the likelihoods are 0. The color coding is the same as in the hypothesis table. We've highlighted the biggest likelihoods with a blue border.

$\theta_N \setminus \theta_S$	0	0.2	0.4	0.6	0.8	1
1	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
0.8	0.00000	1.93428	0.18381	0.00213	0.00000	0.00000
0.6	0.00000	0.06893	0.00655	0.00008	0.00000	0.00000
0.4	0.00000	0.00035	0.00003	0.00000	0.00000	0.00000
0.2	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
0	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000

Likelihoods $p(\text{data}|\theta_S, \theta_N)$ scaled by $100000/c$

4.1 Flat prior

Suppose we have no opinion whatsoever on whether and to what degree sickle-cell protects against malaria. In this case it is reasonable to use a flat prior. Since there are 36 hypotheses each one gets a prior probability of $1/36$. This is given in the table below. Remember each square in the table represents one hypothesis. Because it is a probability table we include the marginal pmf.

$\theta_N \setminus \theta_S$	0	0.2	0.4	0.6	0.8	1	$p(\theta_N)$
1	1/36	1/36	1/36	1/36	1/36	1/36	1/6
0.8	1/36	1/36	1/36	1/36	1/36	1/36	1/6
0.6	1/36	1/36	1/36	1/36	1/36	1/36	1/6
0.4	1/36	1/36	1/36	1/36	1/36	1/36	1/6
0.2	1/36	1/36	1/36	1/36	1/36	1/36	1/6
0	1/36	1/36	1/36	1/36	1/36	1/36	1/6
$p(\theta_S)$	1/6	1/6	1/6	1/6	1/6	1/6	1

Flat prior $p(\theta_S, \theta_N)$: every hypothesis (square) has equal probability

To compute the posterior we simply multiply the likelihood table by the prior table and

normalize. Normalization means making sure the entire table sums to 1.

$\theta_N \setminus \theta_S$	0	0.2	0.4	0.6	0.8	1	$p(\theta_N \text{data})$
0	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
0.2	0.00000	0.88075	0.08370	0.00097	0.00000	0.00000	0.96542
0.4	0.00000	0.03139	0.00298	0.00003	0.00000	0.00000	0.03440
0.6	0.00000	0.00016	0.00002	0.00000	0.00000	0.00000	0.00018
0.8	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
1	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
$p(\theta_S \text{data})$	0.00000	0.91230	0.08670	0.00100	0.00000	0.00000	1.00000

Posterior to flat prior: $p(\theta_S, \theta_N | \text{data})$

To decide whether S confers protection against malaria, we compute the posterior probabilities of ‘some protection’ and of ‘strong protection’. These are computed by summing the corresponding squares in the posterior table.

Some protection: $P(\theta_N > \theta_S) = \text{sum of pink and red} = .99995$

Strong protection: $P(\theta_N - \theta_S > .6) = \text{sum of red} = .88075$

Working from the flat prior, it is effectively certain that sickle-cell provides some protection and very probable that it provides strong protection.

4.2 Informed prior

The experiment was not run without prior information. There was a lot of circumstantial evidence that the sickle-cell gene offered some protection against malaria. For example it was reported that a greater percentage of carriers survived to adulthood.

Here’s one way to build an informed prior. We’ll reserve a reasonable amount of probability for the hypotheses that S gives no protection. Let’s say 24% split evenly among the 6 (orange) cells where $\theta_N = \theta_S$. We know we shouldn’t set any prior probabilities to 0, so let’s spread 6% of the probability evenly among the 15 white cells below the diagonal. That leaves 70% of the probability for the 15 pink and red squares above the diagonal.

$\theta_N \setminus \theta_S$	0	0.2	0.4	0.6	0.8	1	$p(\theta_N)$
0	0.04667	0.04667	0.04667	0.04667	0.04667	0.04000	0.27333
0.2	0.04667	0.04667	0.04667	0.04667	0.04000	0.00400	0.23067
0.4	0.04667	0.04667	0.04667	0.04000	0.00400	0.00400	0.18800
0.6	0.04667	0.04667	0.04667	0.04000	0.00400	0.00400	0.14533
0.8	0.04667	0.04000	0.00400	0.00400	0.00400	0.00400	0.10267
1	0.04000	0.00400	0.00400	0.00400	0.00400	0.00400	0.06000
$p(\theta_S)$	0.27333	0.23067	0.18800	0.14533	0.10267	0.06000	1.0

Informed prior $p(\theta_S, \theta_N)$: makes use of prior information that sickle-cell is protective.

We then compute the posterior pmf.

$\theta_N \setminus \theta_S$	0	0.2	0.4	0.6	0.8	1	$p(\theta_N \text{data})$
$p(\theta_S \text{data})$	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
1	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
0.8	0.00000	0.88076	0.08370	0.00097	0.00000	0.00000	0.96543
0.6	0.00000	0.03139	0.00298	0.00003	0.00000	0.00000	0.03440
0.4	0.00000	0.00016	0.00001	0.00000	0.00000	0.00000	0.00017
0.2	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
0	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
	0.00000	0.91231	0.08669	0.00100	0.00000	0.00000	1.00000

Posterior to informed prior: $p(\theta_S, \theta_N | \text{data})$

We again compute the posterior probabilities of ‘some protection’ and ‘strong protection’.

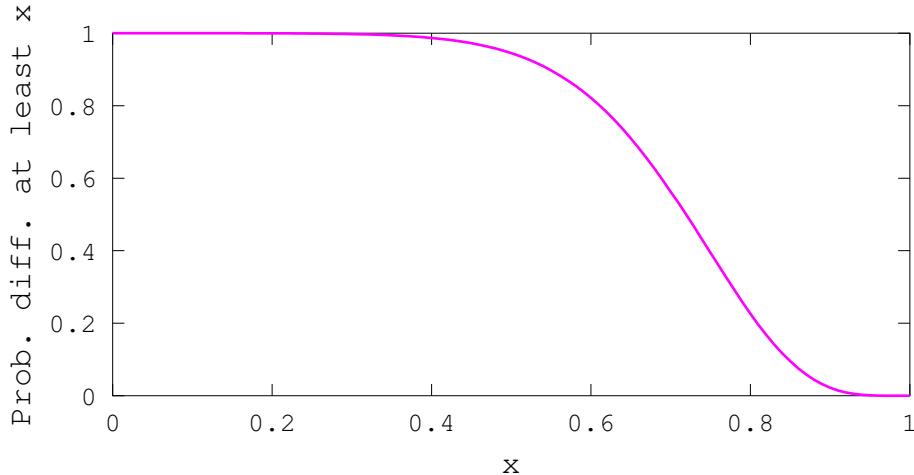
Some protection: $P(\theta_N > \theta_S) = \text{sum of pink and red} = .99996$

Strong protection: $P(\theta_N - \theta_S > .6) = \text{sum of red} = .88076$

Note that the informed posterior is nearly identical to the flat posterior.

4.3 PDALX

The following plot is based on the flat prior. For each x , it gives the probability that $\theta_N - \theta_S \geq x$. To make it smooth we used many more hypotheses.



Probability the difference $\theta_N - \theta_S$ is at least x (PDALX).

Notice that it is virtually certain that the difference is at least .4.

Probability intervals

Class 16, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to find probability intervals given a pmf or pdf.
2. Understand how probability intervals summarize belief in Bayesian updating.
3. Be able to use subjective probability intervals to construct reasonable priors.
4. Be able to construct subjective probability intervals by systematically estimating quantiles.

2 Probability intervals

Suppose we have a pmf $p(\theta)$ or pdf $f(\theta)$ describing our belief about the value of an unknown parameter of interest θ .

Definition: A *p*-probability interval for θ is an interval $[a, b]$ with $P(a \leq \theta \leq b) = p$.

Notes.

1. In the discrete case with pmf $p(\theta)$, this means $\sum_{a \leq \theta_i \leq b} p(\theta_i) = p$.
2. In the continuous case with pdf $f(\theta)$, this means $\int_a^b f(\theta) d\theta = p$.
3. We may say *90%-probability interval* to mean 0.9-probability interval. Probability intervals are also called *credible intervals* to contrast them with confidence intervals, which we'll introduce in the frequentist unit.

Example 1. Between the 0.05 and 0.55 quantiles is a 0.5 probability interval. There are many 50% probability intervals, e.g. the interval from the 0.25 to the 0.75 quantiles.

In particular, notice that the *p*-probability interval for θ is *not unique*.

Q-notation. We can phrase probability intervals in terms of **quantiles**. Recall that the s -quantile for θ is the value q_s with $P(\theta \leq q_s) = s$. So for $s \leq t$, the amount of probability between the s -quantile and the t -quantile is just $t - s$. In these terms, a *p*-probability interval is any interval $[q_s, q_t]$ with $t - s = p$.

Example 2. We have 0.5 probability intervals $[q_{0.25}, q_{0.75}]$ and $[q_{0.05}, q_{0.55}]$.

Symmetric probability intervals.

The interval $[q_{0.25}, q_{0.75}]$ is *symmetric* because the amount of probability remaining on either side of the interval is the same, namely 0.25. If the pdf is not too skewed, the symmetric interval is usually a good default choice.

More notes.

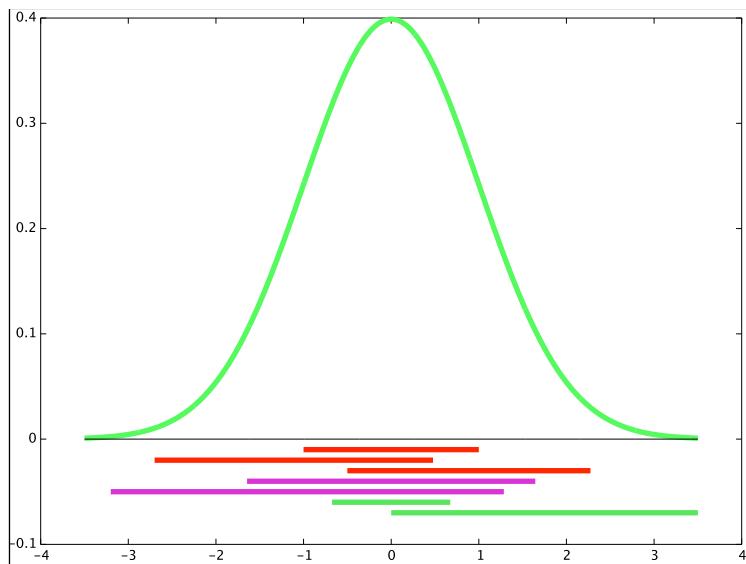
1. Different *p*-probability intervals for θ may have different widths. We can make the width

smaller by centering the interval under the highest part of the pdf. Such an interval is usually a good choice since it contains the most likely values. See the examples below for normal and beta distributions.

- 2.** Since the width can vary for fixed p , a larger p does not always mean a larger width. Here's what is true: if a p_1 -probability interval is fully contained in a p_2 -probability interval, then p_1 is bigger than p_2 .

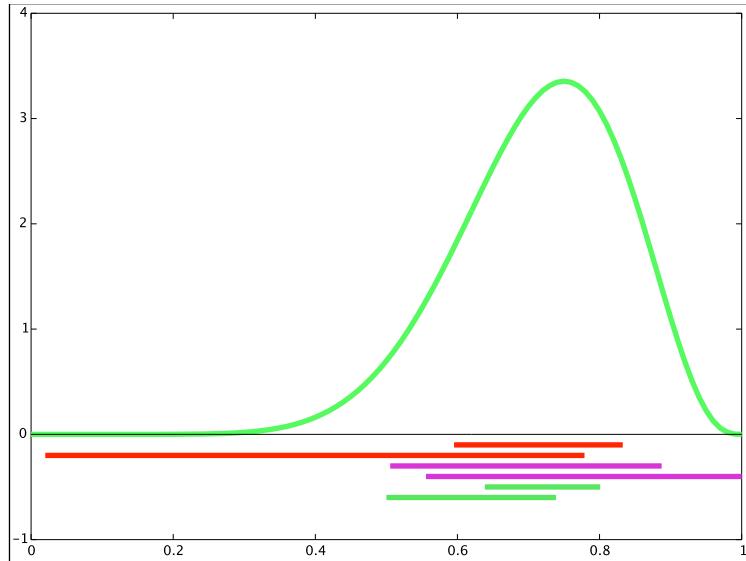
Probability intervals for a normal distribution. The figure shows a number of probability intervals for the standard normal.

1. All of the red bars span a 0.68-probability interval. Notice that the smallest red bar runs between -1 and 1. This runs from the 16th percentile to the 84th percentile so it is a symmetric interval.
2. All the magenta bars span a 0.9-probability interval. They are longer than the red bars because they include more probability. Note again that the shortest magenta bar is symmetric.



$$\text{red} = 0.68, \text{ magenta} = 0.9, \text{ green} = 0.5$$

Probability intervals for a beta distribution. The following figure shows probability intervals for a beta distribution. Notice how the two red bars have very different lengths yet cover the same probability $p = 0.68$.



red = 0.68, magenta = 0.9, green = 0.5

3 Uses of probability intervals

3.1 Summarizing and communicating your beliefs

Probability intervals are an intuitive and effective way to summarize and communicate your beliefs. It's hard to describe an entire function $f(\theta)$ to a friend in words. If the function isn't from a parameterized family then it's especially hard. Even with a beta distribution, it's easier to interpret "I think θ is between 0.45 and 0.65 with 50% probability" than "I think θ follows a beta(8,6) distribution". An exception to this rule of communication might be the normal distribution, but only if the recipient is also comfortable with standard deviation. Of course, what we gain in clarity we lose in precision, since the function contains more information than the probability interval.

Probability intervals also play well with Bayesian updating. If we update from the prior $f(\theta)$ to the posterior $f(\theta|x)$, then the p -probability interval for the posterior will tend to be shorter than the p -probability interval for the prior. In this sense, the data has made us more certain. See for example the election example below.

4 Constructing a prior using subjective probability intervals

Probability intervals are also useful when we do not have a pmf or pdf at hand. In this case, [subjective probability intervals](#) give us a method for constructing a reasonable prior for θ "from scratch". The thought process is to ask yourself a series of questions, e.g., 'what is my expected value for θ ?'; 'my 0.5-probability interval?'; 'my 0.9-probability interval?' Then build a prior that is consistent with these intervals.

4.1 Estimating the intervals directly

Example 3. Building priors

In 2013 there was a special election for a congressional seat in a district in South Carolina. The election pitted Republican Mark Sanford against Democrat Elizabeth Colbert Busch. Let θ be the fraction of the population who favored Busch. Our goal in this example is to build a subjective prior for θ . We'll use the following prior evidence.

- Sanford is a former S. Carolina Congressman and Governor
- He had famously resigned after having an affair in Argentina while he claimed to be hiking the Appalachian trail.
- In 2013 Sanford won the Republican primary over 15 primary opponents.
- In the district in the 2012 presidential election the Republican Romney beat the Democrat Obama 58% to 40%.
- The Colbert bump: Elizabeth Colbert Busch is the sister of well-known comedian Stephen Colbert.

Our strategy will be to use our intuition to construct some probability intervals and then find a beta distribution that approximately matches these intervals. This is subjective so someone else might give a different answer.

Step 1. Use the evidence to construct 0.5 and 0.9 probability intervals for θ .

We'll start by thinking about the 90% interval. The single strongest prior evidence is the 58% to 40% of Romney over Obama. Given the negatives for Sanford we don't expect he'll win much more than 58% of the vote. So we'll put the top of the 0.9 interval at 0.65. With all of Sanford's negatives he could lose big. So we'll put the bottom at 0.3.

$$0.9 \text{ interval: } [0.3, 0.65]$$

For the 0.5 interval we'll pull these endpoints in. It really seems unlikely Sanford will get more votes than Romney, so we can leave 0.25 probability that he'll get above 57%. The lower limit seems harder to predict. So we'll leave 0.25 probability that he'll get under 42%.

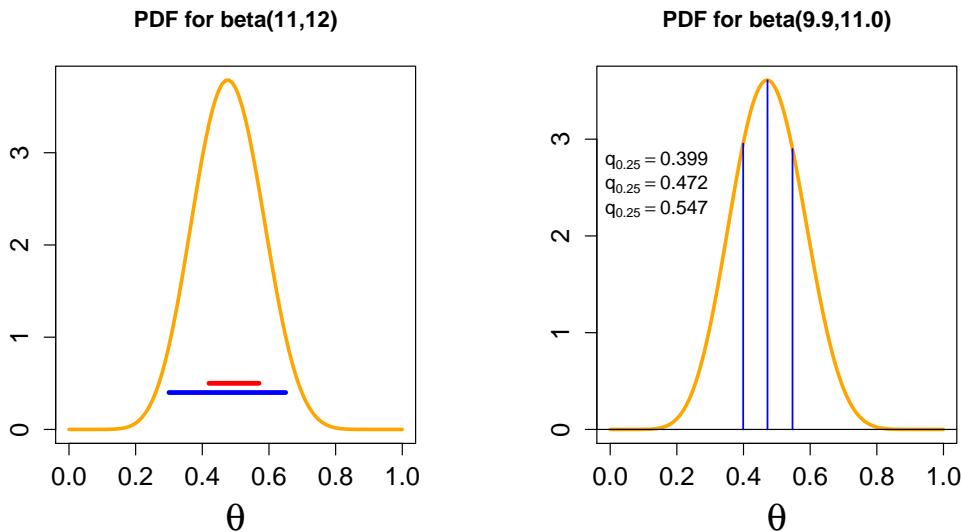
$$0.5 \text{ interval: } [0.42, 0.57]$$

Step 2. Use our 0.5 and 0.9 probability intervals to pick a beta distribution that approximates these intervals. We used the R function `pbeta` and a little trial and error to choose `beta(11,12)`. Here is our R code.

```
a = 11
b = 12
pbeta(0.65, a, b) - pbeta(0.3, a, b)
pbeta(0.57, a, b) - pbeta(0.42, a, b)
```

This computed $P([0.3, 0.65]) = 0.91$ and $P([0.42, 0.57]) = 0.52$. So our intervals are actually 0.91 and 0.52-probability intervals. This is pretty close to what we wanted!

At right is a graph of the density of beta(11,12). The red line shows our interval [0.42, 0.57] and the blue line shows our interval [0.3, 0.65].



beta(11,12) found using probability intervals and beta(9.9,11.0) found using quantiles

4.2 Constructing a prior by estimating quantiles

The method in Example 3 gives a good feel for building priors from probability intervals. Here we illustrate a slightly different way of building a prior by estimating quantiles. The basic strategy is to first estimate the median, then divide and conquer to estimate the first and third quartiles. Finally you choose a prior distribution that fits these estimates.

Example 4. Redo the Sanford vs. Colbert-Busch election example using quantiles.

answer: We start by estimating the median. Just as before the single strongest evidence is the 58% to 40% victory of Romney over Obama. However, given Sanford's negatives and Busch's Colbert bump we'll estimate the median at 0.47.

In a district that went 58 to 40 for the Republican Romney it's hard to imagine Sanford's vote going a lot below 40%. So we'll estimate Sanford 25th percentile as 0.40. Likewise, given his negatives it's hard to imagine him going above 58%, so we'll estimate his 75th percentile as 0.55.

We used R to search through values of a and b for the beta distribution that matches these quartiles the best. Since the beta distribution does not require a and b to be integers we looked for the best fit to 1 decimal place. We found beta(9.9, 11.0). Above is a plot of beta(9.9,11.0) with its actual quartiles shown. These match the desired quartiles pretty well.

Historic note. In the election Sanford won 54% of the vote and Busch won 45.2%. (Source: <http://elections.huffingtonpost.com/2013/mark-sanford-vs-elizabeth-colbert-busch-sc1>

The Frequentist School of Statistics
Class 17, 18.05
Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to explain the difference between the frequentist and Bayesian approaches to statistics.
2. Know our working definition of a statistic and be able to distinguish a statistic from a non-statistic.

2 Introduction

After much foreshadowing, the time has finally come to switch from Bayesian statistics to frequentist statistics. For much of the twentieth century, frequentist statistics has been the dominant school. If you've ever encountered confidence intervals, p -values, t -tests, or χ^2 -tests, you've seen frequentist statistics. With the rise of high-speed computing and big data, Bayesian methods are becoming more common. After we've studied frequentist methods we will compare the strengths and weaknesses of the two approaches.

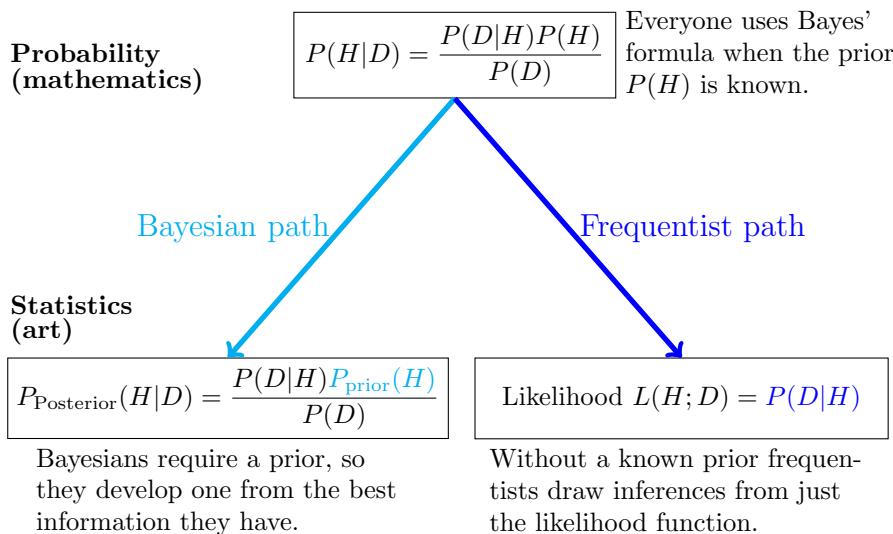
2.1 The fork in the road

Both schools of statistics start with probability. In particular both know and love Bayes' theorem:

$$P(H|D) = \frac{P(D|H) P(H)}{P(D)}.$$

When the prior is known exactly all statisticians will use this formula. For Bayesian inference we take H to be a hypothesis and D some data. Over the last few weeks we have seen that, given a prior and a likelihood model, Bayes' theorem is a complete recipe for updating our beliefs in the face of new data. This works perfectly when the prior was known perfectly. We saw this in our dice examples. We also saw examples of a disease with a known frequency in the general population and a screening test of known accuracy.

In practice we saw that there is usually no universally-accepted prior – different people will have different [a priori beliefs](#) – but we would still like to make useful inferences from data. Bayesians and frequentists take fundamentally different approaches to this challenge, as summarized in the figure below.



The reasons for this split are both practical (ease of implementation and computation) and philosophical (subjectivity versus objectivity and the nature of probability).

2.2 What is probability?

The main philosophical difference concerns the meaning of probability. The term **frequentist** refers to the idea that probabilities represent longterm frequencies of repeatable random experiments. For example, ‘a coin has probability 1/2 of heads’ means that the relative frequency of heads (number of heads out of number of flips) goes to 1/2 as the number of flips goes to infinity. This means the frequentist finds it non-sensical to specify a probability distribution for a parameter with a fixed value. While Bayesians are happy to use probability to describe their incomplete knowledge of a fixed parameter, frequentists reject the use of probability to quantify degree of belief in hypotheses.

Example 1. Suppose I have a bent coin with unknown probability θ of heads. The value of θ may be unknown, but it is a fixed value. Thus, to the frequentist there can be no prior pdf $f(\theta)$. By comparison the Bayesian may agree that θ has a fixed value, but interprets $f(\theta)$ as representing **uncertainty** about that value. Both the Bayesian and the frequentist are perfectly happy with $p(\text{heads}|\theta) = \theta$, since the longterm frequency of heads given θ is θ .

In short, Bayesians put probability distributions on everything (hypotheses and data), while frequentists put probability distributions on (random, repeatable, experimental) data given a hypothesis. For the frequentist when dealing with data from an unknown distribution only the likelihood has meaning. The prior and posterior do not.

3 Working definition of a statistic

Our view of statistics is that it is the art of drawing conclusions (making inferences) from data. With that in mind we can make a simple working definition of a statistic. There is a more formal definition, but we don’t need to introduce it at this point.

Statistic. A [statistic](#) is anything that can be computed from data. Sometimes to be more precise we'll say a statistic is a [rule](#) for computing something from data and the [value](#) of the statistic is what is computed. This can include computing likelihoods where we hypothesize values of the model parameters. But it does not include anything that requires we know the true value of a model parameter with unknown value.

Examples. 1. The mean of data is a statistic. It is a rule that says given data x_1, \dots, x_n compute $\frac{x_1 + \dots + x_n}{n}$.

2. The maximum of data is a statistic. It is a rule that says to pick the maximum value of the data x_1, \dots, x_n .

3. Suppose $x \sim N(\mu, 9)$ where μ is unknown. Then the likelihood

$$p(x|\mu = 7) = \frac{1}{3\sqrt{2\pi}} e^{-\frac{(x-7)^2}{18}}$$

is a statistic. However, the distance of x from the true mean μ is [not](#) a statistic since we cannot compute it without knowing μ

Point statistic. A [point statistic](#) is a single value computed from data. For example, the mean and the maximum are both point statistics. The maximum likelihood estimate is also a point statistic since it is computed directly from the data based on a likelihood model.

Interval statistic. An [interval statistic](#) is an interval computed from data. For example, the range from the minimum to maximum of x_1, \dots, x_n is an interval statistic, e.g. the data 0.5, 1.0, 0.2, 3.0, 5.0 has range [0.2, 5.0].

Set statistic. A [set statistic](#) is a set computed from data.

Example. Suppose we have five dice: 4, 6, 8, 12 and 20-sided. We pick one at random and roll it once. The value of the roll is the data. The set of dice for which this roll is possible is a set statistic. For example, if the roll is a 10 then the value of this set statistic is {12, 20}. If the roll is a 7 then this set statistic has value {8, 12, 20}.

It's important to remember that a statistic is itself a random variable since it is computed from random data. For example, if data is drawn from $N(\mu, \sigma^2)$ then the mean of n data points follows $N(\mu, \sigma^2/n)$.

Sampling distribution. The probability distribution of a statistic is called its [sampling distribution](#).

Point estimate. We can use statistics to make a [point estimate](#) of a parameter θ . For example, if the parameter θ represents the true mean then the data mean \bar{x} is a point estimate of θ .

Null Hypothesis Significance Testing I

Class 17, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Know the definitions of the significance testing terms: NHST, null hypothesis, alternative hypothesis, simple hypothesis, composite hypothesis, significance level, power.
2. Be able to design and run a significance test for Bernoulli or binomial data.
3. Be able to compute a p -value for a normal hypothesis and use it in a significance test.

2 Introduction

Frequentist statistics is often applied in the framework of null hypothesis significance testing (NHST). We will look at the [Neyman-Pearson](#) paradigm which focuses on one hypothesis called the [null hypothesis](#). There are other paradigms for hypothesis testing, but Neyman-Pearson is the most common. Stated simply, this method asks if the data is well outside the region where we would expect to see it under the null hypothesis. If so, then we reject the null hypothesis in favor of a second hypothesis called the alternative hypothesis.

The computations done here all involve the likelihood function. There are two main differences between what we'll do here and what we did in Bayesian updating.

1. The evidence of the data will be considered purely through the likelihood function it will not be weighted by our prior beliefs.
2. We will need a notion of extreme data, e.g. 95 out of 100 heads in a coin toss or a Mayfly that lives for a month.

2.1 Motivating examples

Example 1. Suppose you want to decide whether a coin is fair. If you toss it 100 times and get 85 heads, would you think the coin is likely to be unfair? What about 60 heads? Or 52 heads? Most people would guess that 85 heads is strong evidence that the coin is unfair, whereas 52 heads is no evidence at all. Sixty heads is less clear. [Null hypothesis significance testing \(NHST\)](#) is a frequentist approach to thinking quantitatively about these questions.

Example 2. Suppose you want to compare a new medical treatment to a placebo or the current standard of care. What sort of evidence would convince you that the new treatment is better than the placebo or the current standard? Again, NHST is a quantitative framework for answering these questions.

3 Significance testing

We'll start by listing the ingredients for NHST. Formally they are pretty simple. There is an art to choosing good ingredients. We will explore the art in examples. If you have never seen NHST before just scan this list now and come back to it after reading through the examples and explanations given below.

3.1 Ingredients

- H_0 : the **null hypothesis**. This is the default assumption for the model generating the data.
- H_A : the **alternative hypothesis**. If we reject the null hypothesis we accept this alternative as the best explanation for the data.
- X : the **test statistic**. We compute this from the data.
- **Null distribution**: the probability distribution of X assuming H_0 .
- **Rejection region**: if X is in the rejection region we reject H_0 in favor of H_A .
- **Non-rejection region**: the complement to the rejection region. If X is in this region we do not reject H_0 . Note that we say ‘do not reject’ rather than ‘accept’ because usually the best we can say is that the data does not support rejecting H_0 .

The null hypothesis H_0 and the alternative hypothesis H_A play different roles. Typically we choose H_0 to be either a simple hypothesis or the default which we'll only reject if we have enough evidence against it. The examples below will clarify this.

4 NHST Terminology

In this section we will use one extended example to introduce and explore the terminology used in null hypothesis significance testing (NHST).

Example 3. To test whether a coin is fair we flip it 10 times. If we get an unexpectedly large or small number of heads we'll suspect the coin is unfair. To make this precise in the language of NHST we set up the ingredients as follows. Let θ be the probability that the coin lands heads when flipped.

1. Null hypothesis: $H_0 = \text{‘the coin is fair’}$, i.e. $\theta = 0.5$.
2. Alternative hypothesis: $H_A = \text{‘the coin is not fair’}$, i.e. $\theta \neq .5$
3. Test statistic: $X = \text{number of heads in 10 flips}$
4. Null distribution: This is the probability function based on the null hypothesis

$$p(x | \theta = 0.5) \sim \text{binomial}(10, 0.5).$$

Here is the probability table for the null distribution.

x	0	1	2	3	4	5	6	7	8	9	10
$p(x H_0)$.001	.010	.044	.117	.205	.246	.205	.117	.044	.010	.001

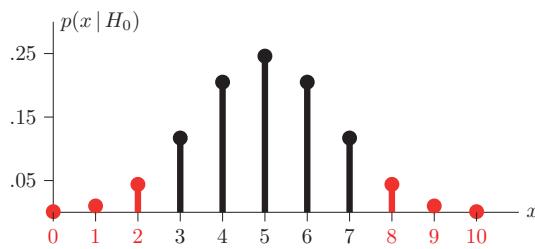
5. Rejection region: under the null hypothesis we expect to get about 5 heads in 10 tosses.

We'll reject H_0 if the number of heads is much fewer or greater than 5. Let's set the rejection region as $\{0, 1, 2, 8, 9, 10\}$. That is, if the number of heads in 10 tosses is in this region we will reject the hypothesis that the coin is fair in favor of the hypothesis that it is not.

We can summarize all this in the graph and probability table below. The rejection region consists of those values of x in red. The probabilities corresponding to it are shaded in red. We also show the null distribution as a stem plot with the rejection values of x in red.

x	0	1	2	3	4	5	6	7	8	9	10
$p(x H_0)$.001	.010	.044	.117	.205	.246	.205	.117	.044	.010	.001

Rejection region and null probabilities as a table for Example 3.



Rejection region and null probabilities as a stem plot for Example 3.

Notes for Example 3:

1. The null hypothesis is the **cautious default**: we won't claim the coin is unfair unless we have compelling evidence.
2. The rejection region consists of data that is **extreme under the null hypothesis**. That is, it consists of the outcomes that are in the tail of the null distribution away from the high probability center. As we'll discuss soon, how far away depends on the significance level α of the test.
3. If we get 3 heads in 10 tosses, then the test statistic is in the non-rejection region. The usual scientific language would be to say that the data 'does not support rejecting the null hypothesis'. Even if we got 5 heads, we would **not** claim that the data proves the null hypothesis is true.

Question: If we have a fair coin what is the probability that we will decide incorrectly it is unfair?

answer: The null hypothesis is that the coin is fair. The question asks for the probability the data from a fair coin will be in the rejection region. That is, the probability that we will get 0, 1, 2, 8, 9 or 10 heads in 10 tosses. This is the sum of the probabilities in red. That is,

$$P(\text{rejecting } H_0 \mid H_0 \text{ is true}) = 0.11$$

Below we will continue with Example 3, define more terms used in NHST and see how to quantify properties of the significance test.

4.1 Simple and composite hypotheses

Definition: simple hypothesis: A **simple hypothesis** is one for which we can specify its distribution completely. A typical simple hypothesis is that a parameter of interest takes a specific value.

Definition: composite hypotheses: If its distribution cannot be fully specified, we say that the hypothesis is **composite**. A typical composite hypothesis is that a parameter of interest lies in a range of values.

In Example 3 the null hypothesis is that $\theta = 0.5$, so the null distribution is $\text{binomial}(10, 0.5)$. Since the null distribution is fully specified, H_0 is simple. The alternative hypothesis is that $\theta \neq 0.5$. This is really many hypotheses in one: θ could be 0.51, 0.7, 0.99, etc. Since the alternative distribution $\text{binomial}(10, \theta)$ is not fully specified, H_A is composite.

Example 4. Suppose we have data x_1, \dots, x_n . Suppose also that our hypotheses are

H_0 : the data is drawn from $N(0, 1)$

H_A : the data is drawn from $N(1, 1)$.

These are both **simple hypotheses** – each hypothesis completely specifies a distribution.

Example 5. (Composite hypotheses.) Now suppose that our hypotheses are

H_0 : the data is drawn from a Poisson distribution of unknown parameter.

H_A : the data is not drawn from a Poisson distribution.

These are both composite hypotheses, as they don't fully specify the distribution.

Example 6. In an ESP experiment a subject is asked to identify the suits of 100 cards drawn (with replacement) from a deck of cards. Let T be the number of successes. The (simple) null hypothesis that the subject does not have ESP is given by

$$H_0: T \sim \text{binomial}(100, 0.25)$$

The (composite) alternative hypothesis that the subject has ESP is given by

$$H_A: T \sim \text{binomial}(100, p) \text{ with } p > 0.25$$

Another (composite) alternative hypothesis that something besides pure chance is going on, i.e. the subject has ESP or anti-ESP. This is given by

$$H_A: T \sim \text{binomial}(100, p), \text{ with } p \neq 0.25$$

Values of $p < 0.25$ represent hypotheses that the subject has a kind of anti-esp.

4.2 Types of error

There are two types of errors we can make. We can incorrectly reject the null hypothesis when it is true or we can incorrectly fail to reject it when it is false. These are unimaginatively labeled **type I** and **type II errors**. We summarize this in the following table.

		True state of nature	
		H_0	H_A
Our decision	Reject H_0	Type I error	correct decision
	'Don't reject' H_0	correct decision	Type II error

Type I: false rejection of H_0

Type II: false non-rejection ('acceptance') of H_0

4.3 Significance level and power

Significance level and power are used to quantify the quality of the significance test. Ideally a significance test would not make errors. That is, it would not reject H_0 when H_0 was true

and would reject H_0 in favor of H_A when H_A was true. Altogether there are 4 important probabilities corresponding to the 2×2 table just above.

$$\begin{array}{ll} P(\text{reject } H_0|H_0) & P(\text{reject } H_0|H_A) \\ P(\text{do not reject } H_0|H_0) & P(\text{do not reject } H_0|H_A) \end{array}$$

The two probabilities we focus on are:

$$\begin{aligned} \text{Significance level} &= P(\text{reject } H_0|H_0) \\ &= \text{probability we incorrectly reject } H_0 \\ &= P(\text{type I error}). \end{aligned}$$

$$\begin{aligned} \text{Power} &= \text{probability we correctly reject } H_0 \\ &= P(\text{reject } H_0|H_A) \\ &= 1 - P(\text{type II error}). \end{aligned}$$

Ideally, a hypothesis test should have a small significance level (near 0) and a large power (near 1). Here are two analogies to help you remember the meanings of significance and power.

Some analogies

1. Think of H_0 as the hypothesis ‘nothing noteworthy is going on’, i.e. ‘the coin is fair’, ‘the treatment is no better than placebo’ etc. And think of H_A as the opposite: ‘something interesting is happening’. Then power is the probability of detecting something interesting when it’s present and significance level is the probability of mistakenly claiming something interesting has occurred.
2. In the U.S. criminal defendants are presumed innocent until proven guilty beyond a reasonable doubt. We can phrase this in NHST terms as

H_0 : the defendant is innocent (the default)

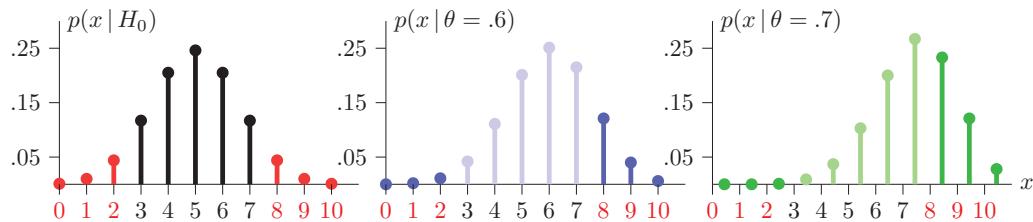
H_A : the defendant is guilty.

Significance level is the probability of finding an innocent person guilty. Power is the probability of correctly finding a guilty party guilty. ‘Beyond a reasonable doubt’ means we should demand the significance level be very small.

Composite hypotheses

H_A is composite in Example 3, so the power is different for different values of θ . We expand the previous probability table to include some alternate values of θ . We do the same with the stem plots. As always in the NHST game, we look at [likelihoods](#): the probability of the data given a hypothesis.

x	0	1	2	3	4	5	6	7	8	9	10
$H_0 : p(x \theta = 0.5)$.001	.010	.044	.117	.205	.246	.205	.117	.044	.010	.001
$H_A : p(x \theta = 0.6)$.000	.002	.011	.042	.111	.201	.251	.215	.121	.040	.006
$H_A : p(x \theta = 0.7)$.000	.0001	.001	.009	.037	.103	.200	.267	.233	.121	.028



Rejection region and null and alternative probabilities for example 3

We use the probability table to compute the significance level and power of this test.

Significance level = probability we reject H_0 when it is true
 = probability the test statistic is in the rejection region when H_0 is true
 = probability the test stat. is in the rejection region of the H_0 row of the table
 = sum of red boxes in the $\theta = 0.5$ row
 = 0.11

Power when $\theta = 0.6$ = probability we reject H_0 when $\theta = 0.6$
 = probability the test statistic is in the rejection region when $\theta = 0.6$
 = probability the test stat. is in the rejection region of the $\theta = 0.6$ row of the table
 = sum of dark blue boxes in the $\theta = 0.6$ row
 = 0.180

Power when $\theta = 0.7$ = probability we reject H_0 when $\theta = 0.7$
 = probability the test statistic is in the rejection region when $\theta = 0.7$
 = probability the test stat. is in the rejection region of the $\theta = 0.7$ row of the table
 = sum of dark green boxes in the $\theta = 0.7$ row
 = 0.384

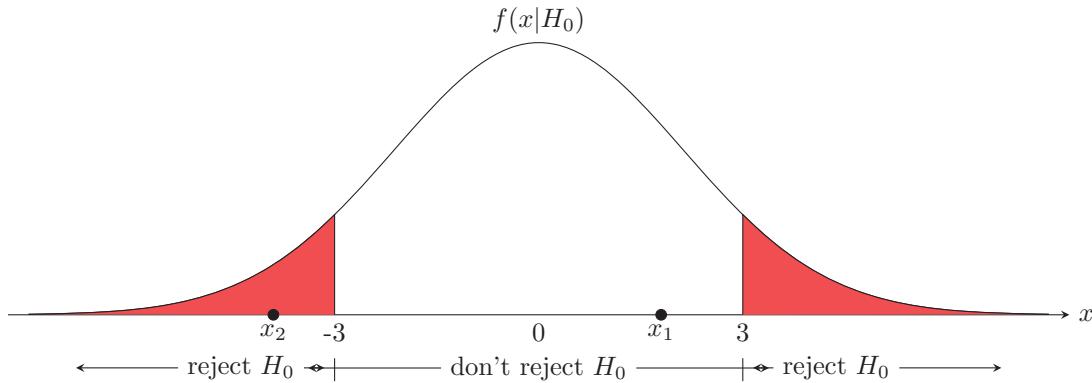
We see that the power is greater for $\theta = 0.7$ than for $\theta = 0.6$. This isn't surprising since we expect it to be easier to recognize that a 0.7 coin is unfair than it is to recognize 0.6 coin is unfair. Typically, we get higher power when the alternate hypothesis is farther from the null hypothesis. In Example 3, it would be quite hard to distinguish a fair coin from one with $\theta = 0.51$.

4.4 Conceptual sketches

We illustrate the notions of null hypothesis, rejection region and power with some sketches of the pdfs for the null and alternative hypotheses.

4.4.1 Null distribution: rejection and non-rejection regions

The first diagram below illustrates a null distribution with rejection and non-rejection regions. Also shown are two possible test statistics: x_1 and x_2 .



The test statistic x_1 is in the non-rejection region. So, if our data produced the test statistic x_1 then we would not reject the null hypothesis H_0 . On the other hand the test statistic x_2 is in the rejection region, so if our data produced x_2 then we would reject the null hypothesis in favor of the alternative hypothesis.

There are several things to note in this picture.

1. The rejection region consists of values far from the center of the null distribution.
2. The rejection region is two-sided. We will also see examples of one-sided rejection regions as well.
3. The alternative hypothesis is not mentioned. We reject or don't reject H_0 based only on the likelihood $f(x|H_0)$, i.e. the probability of the test statistic conditioned on H_0 . As we will see, the alternative hypothesis H_A should be considered when choosing a rejection region, but formally it does not play a role in rejecting or not rejecting H_0 .
4. Sometimes we rather lazily call the non-rejection region the [acceptance region](#). This is technically incorrect because we never truly accept the null hypothesis. We either reject or say the data does not support rejecting H_0 . This is often summarized by the statement: [you can never prove the null hypothesis](#).

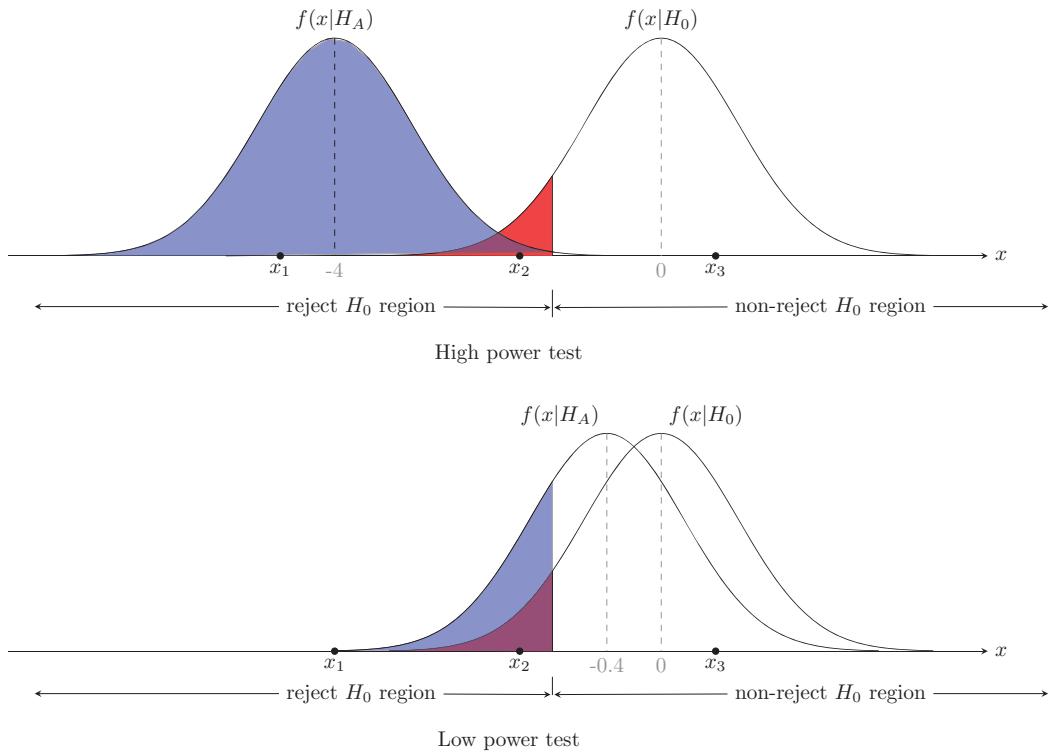
4.4.2 High and low power tests

The next two figures show high and low power tests.

The shaded area under $f(x|H_0)$ represents the significance level. Remember the significance level is

- The probability of falsely rejecting the null hypothesis when it is true.
- The probability the test statistic falls in the rejection region even though H_0 is true.

Likewise, the shaded area under $f(x|H_A)$ represents the power, i.e. the probability that the test statistic is in the rejection (of H_0) region when H_A is true. Both tests have the same significance level, but if $f(x|H_A)$ has considerable overlap with $f(x|H_0)$ the power is much lower. It is well worth your while to thoroughly understand these graphical representations of significance testing.



In both tests both distributions are standard normal. The null distribution, rejection region and significance level are all the same. (The significance level is the red/purple area under $f(x|H_0)$ and above the rejection region.) In the top figure we see the means of the two distributions are 4 standard deviations apart. Since the areas under the densities have very little overlap the test has high power. That is if the data x is drawn from H_A it will almost certainly be in the rejection region. For example x_3 would be a very surprising outcome for the H_A distribution.

In the bottom figure we see the means of the two distributions are just 0.4 standard deviations apart. Since the areas under the densities have a lot of overlap the test has low power. That is if the data x is drawn from H_A it is highly likely to be in the non-rejection region. For example x_3 would be not be a very surprising outcome for the H_A distribution.

Typically we can increase the power of a test by increasing the amount of data and thereby decreasing the variance of the null and alternative distributions. In experimental design it is important to determine ahead of time the number of trials or subjects needed to achieve a desired power.

Example 7. Suppose a drug for a disease is being compared to a placebo. We choose our null and alternative hypotheses as

H_0 = the drug does not work better than the placebo

H_A = the drug works better than the placebo

The power of the hypothesis test is the probability that the test will conclude that the drug is better, if it is indeed truly better. The significance level is the probability that the test will conclude that the drug works better, when in fact it does not.

5 Designing a hypothesis test

Formally all a hypothesis test requires is H_0 , H_A , a test statistic and a rejection region. In practice the design is often done using the following steps.

1. Pick the null hypothesis H_0 .

The choice of H_0 and H_A is not mathematics. It's art and custom. We often choose H_0 to be simple. Or we often choose H_0 to be the simplest or most cautious explanation, i.e. no effect of drug, no ESP, no bias in the coin.

2. Decide if H_A is one-sided or two-sided.

In the example 3 we wanted to know if the coin was unfair. An unfair coin could be biased for or against heads, so $H_A : \theta \neq 0.5$ is a two-sided hypothesis. If we only care whether or not the coin is biased for heads we could use the one-sided hypothesis $H_A : \theta > 0.5$.

3. Pick a test statistic.

For example, the sample mean, sample total, or sample variance. Often the choice is obvious. Some standard statistics that we will encounter are z , t , and χ^2 . We will learn to use these statistics as we work examples over the next few classes. One thing we will say repeatedly is that the distributions that go with these statistics are always conditioned on the null hypothesis. That is, we will compute likelihoods such as $f(z | H_0)$.

4. Pick a significance level and determine the rejection region.

We will usually use α to denote the significance level. The Neyman-Pearson paradigm is to pick α in advance. Typical values are 0.1, 0.05, 0.01. Recall that the significance level is the probability of a type I error, i.e. of incorrectly rejecting the null hypothesis when it is true. The value we choose will depend on the consequences of a type I error.

Once the significance level is chosen we can determine the rejection region in the tail(s) of the null distribution. In Example 3, H_A is two sided so the rejection region is split between the two tails of the null distribution. This distribution is given in the following table:

x	0	1	2	3	4	5	6	7	8	9	10
$p(x H_0)$.001	.010	.044	.117	.205	.246	.205	.117	.044	.010	.001

If we set $\alpha = 0.05$ then the rejection region must contain at most .05 probability. For a two-sided rejection region we get

$$\{0, 1, 9, 10\}.$$

If we set $\alpha = 0.01$ the rejection region is

$$\{0, 10\}.$$

Suppose we change H_A to 'the coin is biased in favor of heads'. We now have a one-sided hypothesis $\theta > 0.5$. Our rejection region will now be in the right-hand tail since we don't want to reject H_0 in favor of H_A if we get a small number of heads. Now if $\alpha = 0.05$ the rejection region is the one-sided range

$$\{9, 10\}.$$

If we set $\alpha = 0.01$ then the rejection region is

$$\{10\}.$$

5. Determine the power(s).

As we saw in Example 3, once the rejection region is set we can determine the power of the test at various values of the alternate hypothesis.

Example 8. (Consequences of significance) If $\alpha = 0.1$ then we'd expect a 10% type I error rate. That is, we expect to reject the null hypothesis in 10% of those experiments where the null hypothesis is true. Whether 0.1 is a reasonable significance level depends on the decisions that will be made using it.

For example, if you were running an experiment to determine if your chocolate is more than 72% cocoa then a 10% error type I error rate is probably okay. That is, falsely believing some 72% chocolate is greater than 72%, is probably acceptable. On the other hand, if your forensic lab is identifying fingerprints for a murder trial then a 10% type I error rate, i.e. mistakenly claiming that fingerprints found at the crime scene belonged to someone who was truly innocent, is definitely not acceptable.

Significance for a composite null hypothesis. If H_0 is composite then $P(\text{type I error})$ depends on which member of H_0 is true. In this case the significance level is defined as the maximum of these probabilities.

6 Critical values

Critical values are like quantiles except they refer to the probability to the right of the value instead of the left.

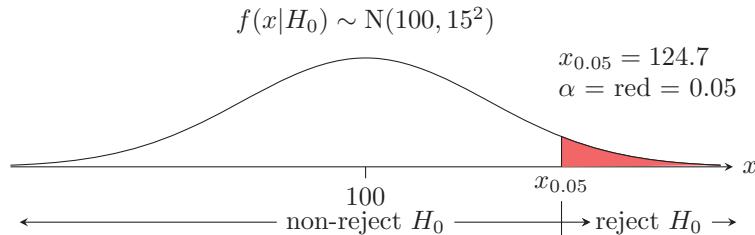
Example 9. Use R to find the 0.05 critical value for the standard normal distribution.

answer: We label this critical value $z_{0.05}$. The critical value $z_{0.05}$ is just the 0.95 quantile, i.e. it has 5% probability to its right and therefore 95% probability to its left. We computed it with the R function qnorm: `qnorm(0.95, 0, 1)`, which returns 1.64.

In a typical significance test the rejection region consists of one or both tails of the null distribution. The value of the test significant that marks the start of the rejection region is a **critical value**. We show this and the notation used in some examples.

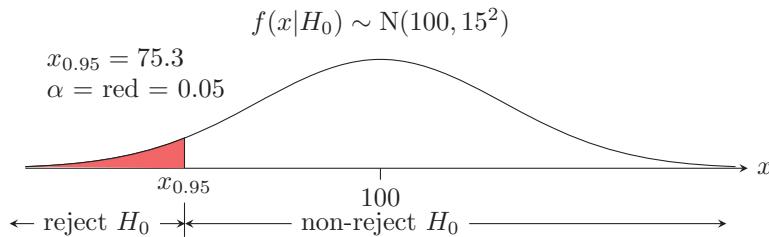
Example 10. Critical values and rejection regions. Suppose our test statistic x has null distribution is $N(100, 15^2)$, i.e. $f(x|H_0) \sim N(100, 15^2)$. Suppose also that our rejection region is right-sided and we have a significance level of 0.05. Find the critical value and sketch the null distribution and rejection region.

answer: The notation used for the critical value with right tail containing probability 0.05 is $x_{0.05}$. The critical value $x_{0.05}$ is just the 0.95 quantile, i.e. it has 5% probability to its right and therefore 95% probability to its left. We computed it with the R function qnorm: `qnorm(0.95, 100, 15)`, which returned 124.7. This is shown in the figure below.



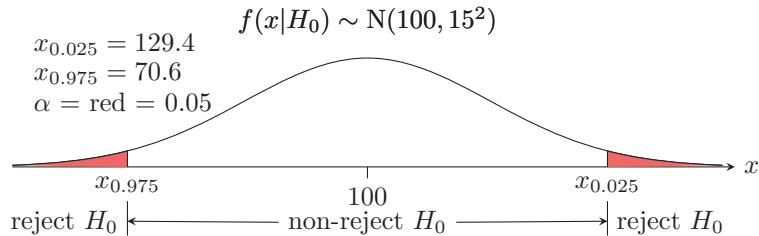
Example 11. Critical values and rejection regions. Repeat the previous example for a left-sided rejection region with significance level 0.05. In this case, the start of the rejection region is at the 0.05 quantile. Since there is 95%

answer: In this case the critical value has 0.05 probability to its left and therefore 0.95 probability to its right. So we label it $x_{0.95}$. Since it is the 0.05 quantile compute it with the R function: `qnorm(0.05, 100, 15)`, which returned 75.3.



Example 12. Critical values. Repeat the previous example for a two-sided rejection region. Put half the significance in each tail.

answer: To have a total significance of 0.05 we put 0.025 in each tail. That is, the left tail starts at $x_{0.975} = q_{0.025}$ and the right tail starts at $x_{0.025} = q_{0.975}$. We compute these values with `qnorm(0.025, 100, 15)` and `qnorm(0.975, 100, 15)`. The values are shown in the figure below.



7 p-values

In practice people often specify the significance level and do the significance test using what are called **p-values**. We will first define *p*-value and see that

If the *p*-value is less than the significance level α then we reject H_0 . Otherwise we do not reject H_0 .

Definition. The *p*-value is the probability, assuming the null hypothesis, of seeing data **at least as extreme as the experimental data**. What ‘at least as extreme’ means depends on

the experimental design.

We illustrate the definition and use of p -values with a simple one-sided example. In later classes we will look at two-sided examples. This example also introduces the [z-test](#). All this means is that our test statistic is standard normal (or approximately standard normal).

Example 13. The z-test for normal hypotheses

IQ is normally distributed in the population according to a $N(100, 15^2)$ distribution. We suspect that most MIT students have above average IQ so we frame the following hypotheses.

- H_0 = MIT student IQs are distributed identically to the general population
= MIT IQ's follow a $N(100, 15^2)$ distribution.
- H_A = MIT student IQs tend to be higher than those of the general population
= the average MIT student IQ is greater than 100.

Notice that H_A is one-sided.

Suppose we test 9 students and find they have an average IQ of $\bar{x} = 112$. Can we reject H_0 at a significance level $\alpha = 0.05$?

answer: To compute p we first standardize the data: Under the null hypothesis $\bar{x} \sim N(100, 15^2/9)$ and therefore

$$z = \frac{\bar{x} - 100}{15/\sqrt{9}} = \frac{36}{15} = 2.4 \sim N(0, 1).$$

That is, the null distribution for z is standard normal. We call z a [z-statistic](#), we will use it as our test statistic.

For a right-sided alternative hypothesis the phrase ‘data at least as extreme’ is a one-sided tail to the right of z . The p -value is then

$$p = P(Z \geq 2.4) = 1 - \text{pnorm}(2.4, 0, 1) = 0.0081975.$$

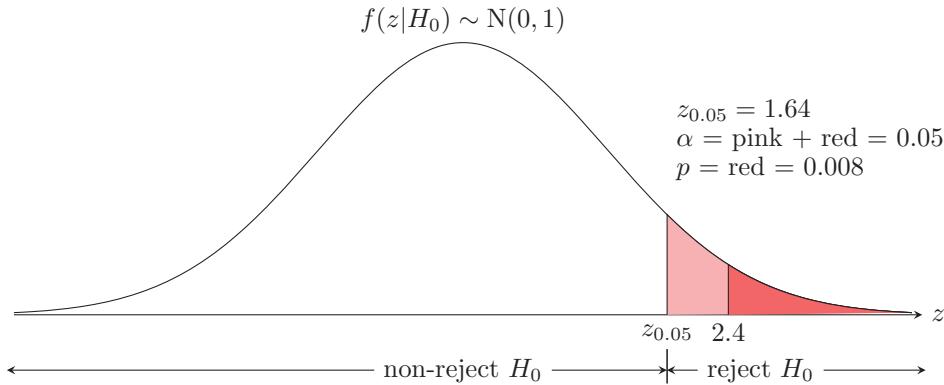
Since $p \leq \alpha$ we reject the null hypothesis. The reason this works is explained below. We phrase our conclusion as

We reject the null hypothesis in favor of the alternative hypothesis that MIT students have higher IQs on average. We have done this at significance level 0.05 with a p -value of 0.008.

Notes: 1. The average $\bar{x} = 112$ is random: if we ran the experiment again we could get a different value for \bar{x} .

2. We could use the statistic \bar{x} directly. Standardizing is fairly standard because, with practice, we will have a good feel for the meaning of different z -values.

The justification for rejecting H_0 when $p \leq \alpha$ is given in the following figure.



In this example $\alpha = 0.05$, $z_{0.05} = 1.64$ and the rejection region is the range to the right of $z_{0.05}$. Also, $z = 2.4$ and the p -value is the probability to the right of z . The picture illustrates that

- $z = 2.64$ is in the rejection region
- is the same as z is to the right of $z_{0.05}$
- is the same as the probability to the right of z is less than 0.05
- which means $p < 0.05$.

8 More examples

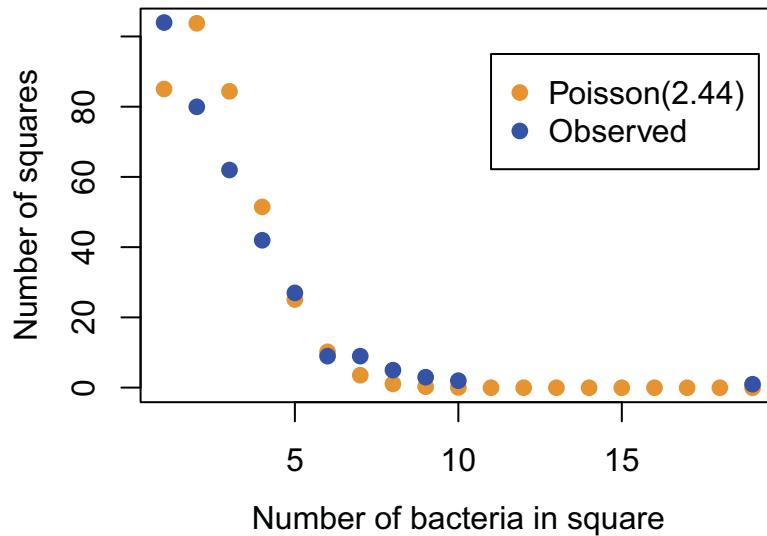
Hypothesis testing is widely used in inferential statistics. We don't expect that the following examples will make perfect sense at this time. Read them quickly just to get a sense of how hypothesis testing is used. We will explore the details of these examples in class.

Example 14. The chi-square statistic and goodness of fit. (Rice, example B, p.313)

To test the level of bacterial contamination, milk was spread over a grid with 400 squares. The amount of bacteria in each square was counted. We summarize in the table below. The bottom row of the table is the number of different squares that had a given amount of bacteria.

Amount of bacteria	0	1	2	3	4	5	6	7	8	9	10	19
Number of squares	56	104	80	62	42	27	9	9	5	3	2	1

We compute that the average amount of bacteria per square is 2.44. Since the Poisson(λ) distribution is used to model counts of relatively rare events and the parameter λ is the expected value of the distribution. we decide to see if these counts could come from a Poisson distribution. To do this we first graphically compare the observed frequencies with those expected from Poisson(2.44).



The picture is suggestive, so we do a hypothesis test with

H_0 : the samples come from a $\text{Poisson}(2.44)$ distribution.

H_A : the samples come from a different distribution.

We use a chi-square statistic, so called because it (approximately) follows a chi-square distribution. To compute X^2 we first combine the last few cells in the table so that the minimum expected count is around 5 (a general rule-of-thumb in this game.)

The expected number of squares with a certain amount of bacteria comes from considering 400 trials from a $\text{Poisson}(2.44)$ distribution, e.g., with $l = 2.44$ the expected number of squares with 3 bacteria is $400 \times e^{-l} \frac{l^3}{3!} = 84.4$.

The chi-square statistic is $\sum \frac{(O_i - E_i)^2}{E_i}$, where O_i is the observed number and E_i is the expected number.

Number per square	0	1	2	3	4	5	6	> 6
Observed	56	104	80	62	42	27	9	20
Expected	34.9	85.1	103.8	84.4	51.5	25.1	10.2	5.0
Component of X^2	12.8	4.2	5.5	6.0	1.7	0.14	0.15	44.5

Summing up we get $X^2 = 74.9$.

Since the mean (2.44) and the total number of trials (400) are fixed, the 8 cells only have 6 degrees of freedom. So, assuming H_0 , our chi-square statistic follows (approximately) a χ_6^2 distribution. Using this distribution, $P(X^2 > 74.59) = 0$ (to at least 6 decimal places). Thus we decisively reject the null hypothesis in favor of the alternate hypothesis that the distribution is not $\text{Poisson}(2.44)$.

To analyze further, look at the individual components of X^2 . There are large contributions in the tail of the distribution, so that is where the fit goes awry.

Example 15. Student's t test.

Suppose we want to compare a medical treatment for increasing life expectancy with a placebo. We give n people the treatment and m people the placebo. Let X_1, \dots, X_n be the number of years people live after receiving the treatment. Likewise, let Y_1, \dots, Y_m be the number of years people live after receiving the placebo. Let \bar{X} and \bar{Y} be the sample means. We want to know if the difference between \bar{X} and \bar{Y} is statistically significant. We frame this as a hypothesis test. Let μ_X and μ_Y be the (unknown) means.

$$H_0 : \mu_X = \mu_Y, \quad H_A : \mu_X \neq \mu_Y.$$

With certain assumptions and a proper formula for the pooled standard error s_p the test statistic $t = \frac{\bar{X} - \bar{Y}}{s_p}$ follow a t distribution with $n + m - 2$ degrees of freedom. So our rejection region is determined by a threshold t_0 with $P(t > t_0) = \alpha$.

MIT OpenCourseWare
<https://ocw.mit.edu>

18.05 Introduction to Probability and Statistics
Spring 2014

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.

Null Hypothesis Significance Testing II

Class 18, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to list the steps common to all null hypothesis significance tests.
2. Be able to define and compute the probability of Type I and Type II errors.
3. Be able to look up and apply one and two sample t -tests.

2 Introduction

We continue our study of significance tests. In these notes we will introduce two new tests: one-sample t -tests and two-sample t -tests. You should pay careful attention to the fact that every test makes some assumptions about the data – often that is drawn from a normal distribution. You should also notice that all the tests follow the same pattern. It is just the computation of the test statistic and the type of the null distribution that changes.

3 Review: setting up and running a significance test

There is a fairly standard set of steps one takes to set up and run a null hypothesis significance test.

1. Design an experiment to collect data and choose a test statistic x to be computed from the data. The key requirement here is to know the null distribution $f(x|H_0)$. To compute power, one must also know the alternative distribution $f(x|H_A)$.
2. Decide if the test is one or two-sided based on H_A and the form of the null distribution.
3. Choose a significance level α for rejecting the null hypothesis. If applicable, compute the corresponding power of the test.
4. Run the experiment to collect data x_1, x_2, \dots, x_n .
5. Compute the test statistic x .
6. Compute the p -value corresponding to x using the null distribution.
7. If $p < \alpha$, reject the null hypothesis in favor of the alternative hypothesis.

Notes.

1. Rather than choosing a significance level, you could instead choose a rejection region and reject H_0 if x falls in this region. The corresponding significance level is then the probability that x falls in the rejection region.

2. The null hypothesis is often the ‘cautious hypothesis’. The lower we set the significance level, the more “evidence” we will require before rejecting our cautious hypothesis in favor of a more sensational alternative. It is standard practice to publish the p value itself so that others may draw their own conclusions.

3. **A key point of confusion:** A significance level of 0.05 does **not** mean the test only makes mistakes 5% of the time. It means that **if the null hypothesis is true**, then the probability the test will mistakenly reject it is 5%. The power of the test measures the accuracy of the test when the alternative hypothesis is true. Namely, the power of the test is the probability of rejecting the null hypothesis **if the alternative hypothesis is true**. Therefore the probability of falsely failing to reject the null hypothesis is 1 minus the power.

Errors. We can summarize these two types of errors and their probabilities as follows:

Type I error = rejecting H_0 when H_0 is true.

Type II error = failing to reject H_0 when H_A is true.

$P(\text{type I error})$ = probability of falsely rejecting H_0
= $P(\text{test statistic is in the rejection region} \mid H_0)$
= significance level of the test

$P(\text{type II error})$ = probability of falsely not rejecting H_0
= $P(\text{test statistic is in the acceptance region} \mid H_A)$
= 1 - power.

Helpful analogies. In terms of medical testing for a disease, a Type I error is a false positive and a Type II error is a false negative. In terms of a jury trial, a Type I error is convicting an innocent defendant and a Type II error is acquitting a guilty defendant.

4 Understanding a significance test

Questions to ask:

1. How did they collect data? What is the experimental setup?

2. What are the null and alternative hypotheses?

3. What type of significance test was used?

Does their data match the criteria needed to use this type of test?

How robust is the test to deviations from this criteria.

4. For example, some tests comparing two groups of data assume that the groups are drawn from distributions that have the same variance. This needs to be verified before applying the test. Often the check is done using another significance test designed to compare the variances of two groups of data.

5. How is the p -value computed?

A significance test comes with a test statistic and a null distribution. In most tests the p -value is

$$p = P(\text{data at least as extreme as what we got} \mid H_0)$$

What does ‘data at least as extreme as the data we saw,’ mean? I.e. is the test one or two-sided.

6. What is the significance level α for this test? If $p < \alpha$ then the experimenter will reject H_0 in favor of H_A .

5 t tests

Many significance tests assume that the data are drawn from a normal distribution, so before using such a test you should examine the data to see if the normality assumption is reasonable. We will describe how to do this in more detail later, but plotting a histogram is a good start. Like the z -test, the one-sample and two-sample t -tests we’ll consider below start from this normality assumption.

We don’t expect you to memorize all the computational details of these tests and those to follow. In real life, you have access to textbooks, google, and wikipedia; on the exam, you’ll have your notecard. Instead, you should be able to identify when a t test is appropriate and apply this test after looking up the details and using a table or software like R.

5.1 z -test

Let’s first review the z -test.

- Data: we assume $x_1, x_2, \dots, x_n \sim N(\mu, \sigma^2)$, where μ is unknown and σ is known.
- Null hypothesis: $\mu = \mu_0$ for some specific value μ_0
- Test statistic: $z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$ = standardized mean
- Null distribution: $f(z | H_0)$ is the pdf of $Z \sim N(0, 1)$
- One-sided p -value (right side): $p = P(Z > z | H_0)$
 One-sided p -value (left side): $p = P(Z < z | H_0)$
 Two-sided p -value: $p = P(|Z| > |z|)$.

Example 1. Suppose that we have data that follows a normal distribution of unknown mean μ and known variance 4. Let the null hypothesis H_0 be that $\mu = 0$. Let the alternative hypothesis H_A be that $\mu > 0$. Suppose we collect the following data:

$$1, 2, 3, 6, -1$$

At a significance level of $\alpha = 0.05$, should we reject the null hypothesis?

answer: There are 5 data points with average $\bar{x} = 2.2$. Because we have normal data with a known variance we should use a z test. Our z statistic is

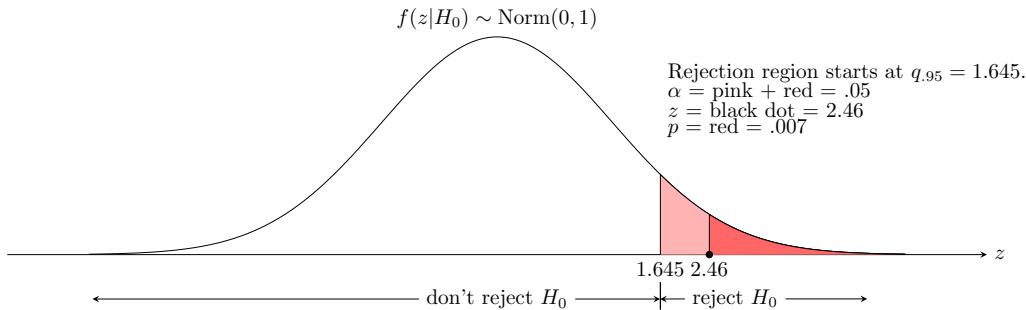
$$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} = \frac{2.2 - 0}{2/\sqrt{5}} = 2.460$$

Our test is one-sided because the alternative hypothesis is one-sided. So (using R) our p -value is

$$p = P(Z > z) = P(Z > 2.460) = 0.007$$

Since $p < .05$, we reject the null hypothesis in favor of the alternative hypothesis $\mu > 0$.

We can visualize the test as follows:



5.2 The Student t distribution

'Student' is the pseudonym used by the William Gosset who first described this test and this test and distribution. See http://en.wikipedia.org/wiki/Student's_t-test

The t -distribution is symmetric and bell-shaped like the normal distribution. It has a parameter df which stands for degrees of freedom. For df small the t -distribution has more probability in its tails than the standard normal distribution. As df increases $t(df)$ becomes more and more like the standard normal distribution.

Here is a simple applet that shows $t(df)$ and compares it to the standard normal distribution: <http://mathlets.org/mathlets/t-distribution/>

As usual in R, the functions `pt`, `dt`, `qt`, `rt` correspond to cdf, pdf, quantiles, and random sampling for a t distribution. Remember that you can type `?dt` in RStudio to view the help file specifying the parameters of `dt`. For example, `pt(1.65, 3)` computes the probability that x is less than or equal 1.65 given that x is sampled from the t distribution with 3 degrees of freedom, i.e. $P(x \leq 1.65)$ given that $x \sim t(3)$.

5.3 One sample t -test

For the z -test, we assumed that the variance of the underlying distribution of the data was known. However, it is often the case that we don't know σ and therefore we must estimate it from the data. In these cases, we use a one sample t -test instead of a z -test and the **studentized** mean in place of the standardized mean

- Data: we assume $x_1, x_2, \dots, x_n \sim N(\mu, \sigma^2)$, where both μ and σ are unknown.
- Null hypothesis: $\mu = \mu_0$ for some specific value μ_0
- Test statistic:

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$$

where

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

Here t is called the **Studentized mean** and s^2 is called the **sample variance**. The latter is an estimate of the true variance σ^2 .

- Null distribution: $f(t|H_0)$ is the pdf of $T \sim t(n-1)$, the t distribution with $n-1$ degrees of freedom.*
- One-sided p -value (right side): $p = P(T > t|H_0)$
One-sided p -value (left side): $p = P(T < t|H_0)$
Two-sided p -value: $p = P(|T| > |t|)$.

*It's a theorem (not an assumption) that if the data is normal with mean μ_0 then the Studentized mean follows a t -distribution. A proof would take us too far afield, but you can look it up if you want: http://en.wikipedia.org/wiki/Student%27s_t-distribution#Derivation

Example 2. Now suppose that in the previous example the variance is unknown. That is, we have data that follows a normal distribution of unknown mean μ and unknown variance σ . Suppose we collect the same data as before:

$$1, 2, 3, 6, -1$$

As above, let the null hypothesis H_0 be that $\mu = 0$ and the alternative hypothesis H_A be that $\mu > 0$. At a significance level of $\alpha = 0.05$, should we reject the null hypothesis?

answer: There are 5 data points with average $\bar{x} = 2.2$. Because we have normal data with unknown mean and unknown variance we should use a one-sample t test. Computing the sample variance we get

$$s^2 = \frac{1}{4} ((1-2.2)^2 + (2-2.2)^2 + (3-2.2)^2 + (6-2.2)^2 + (-1-2.2)^2) = 6.7$$

Our t statistic is

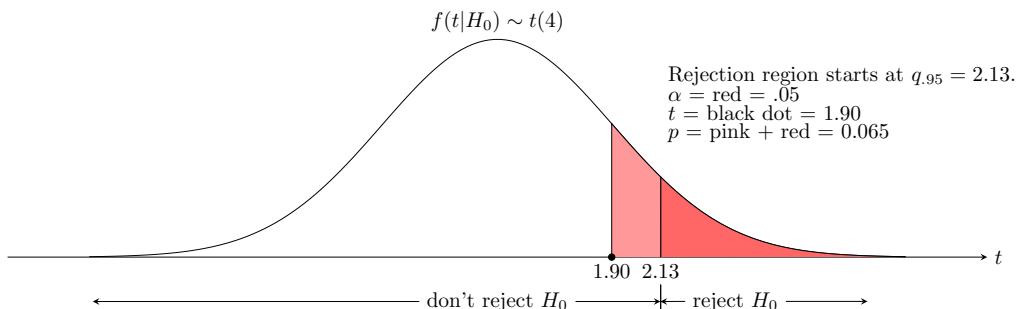
$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} = \frac{2.2 - 0}{\sqrt{6.7}/\sqrt{5}} = 1.901$$

Our test is one-sided because the alternative hypothesis is one-sided. So (using R) the p -value is

$$p = P(T > t) = P(T > 1.901) = 1 - pt(1.901, 4) = 0.065$$

Since $p > .05$, we do not reject the null hypothesis.

We can visualize the test as follows:



5.4 Two-sample t -test with equal variances

We next consider the case of comparing the means of two samples. For example, we might be interested in comparing the mean efficacies of two medical treatments.

- Data: We assume we have two sets of data drawn from normal distributions

$$\begin{aligned}x_1, x_2, \dots, x_n &\sim N(\mu_1, \sigma^2) \\y_1, y_2, \dots, y_m &\sim N(\mu_2, \sigma^2)\end{aligned}$$

where the means μ_1 and μ_2 and the variance σ^2 are all unknown. Notice the assumption that the two distributions have the **same variance**. Also notice that there are n samples in the first group and m samples in the second.

- Null hypothesis: $\mu_1 = \mu_2$ (the values of μ_1 and μ_2 are not specified)

- Test statistic:

$$t = \frac{\bar{x} - \bar{y}}{s_p},$$

where s_p^2 is the **pooled variance**

$$s_p^2 = \frac{(n-1)s_x^2 + (m-1)s_y^2}{n+m-2} \left(\frac{1}{n} + \frac{1}{m} \right)$$

Here s_x^2 and s_y^2 are the sample variances of the x_i and y_j respectively. The expression for t is somewhat complicated, but the basic idea remains the same and it still results in a known null distribution.

- Null distribution: $f(t | H_0)$ is the pdf of $T \sim t(n+m-2)$.
- One-sided p -value (right side): $p = P(T > t | H_0)$
One-sided p -value (left side): $p = P(T < t | H_0)$
Two-sided p -value: $p = P(|T| > |t|)$.

Note 1: Some authors use a different notation. They define the pooled variance as

$$s_{p\text{-other-authors}}^2 = \frac{(n-1)s_x^2 + (m-1)s_y^2}{n+m-2}$$

and what we called the pooled variance they point out is the estimated variance of $\bar{x} - \bar{y}$. That is,

$$s_p^2 = s_{p\text{-other-authors}}^2 \times (1/n + 1/m) \approx s_{\bar{x}-\bar{y}}^2$$

Note 2: There is a version of the two-sample t -test that allows the two groups to have different variances. In this case the test statistic is a little more complicated but R will handle it with equal ease.

Example 3. The following data comes from a real study in which 1408 women were admitted to a maternity hospital for (i) medical reasons or through (ii) unbooked emergency

admission. The duration of pregnancy is measured in complete weeks from the beginning of the last menstrual period. We can summarize the data as follows:

Medical: 775 observations with $\bar{x}_M = 39.08$ and $s_M^2 = 7.77$.

Emergency: 633 observations with $\bar{x}_E = 39.60$ and $s_E^2 = 4.95$

Set up and run a two-sample *t*-test to investigate whether the mean duration differs for the two groups.

What assumptions did you make?

answer: The pooled variance for this data is

$$s_p^2 = \frac{774(7.77) + 632(4.95)}{1406} \left(\frac{1}{775} + \frac{1}{633} \right) = .0187$$

The *t* statistic for the null distribution is

$$\frac{\bar{x}_M - \bar{x}_E}{s_p} = -3.8064$$

We have 1406 degrees of freedom. Using R to compute the two-sided *p*-value we get

$$p = P(|T| > |t|) = 2*dt(-3.8064, 1406) = 0.00015$$

p is very small, much smaller than $\alpha = .05$ or $\alpha = .01$. Therefore we reject the null hypothesis in favor of the alternative that there is a difference in the mean durations.

Rather than compute the two-sided *p*-value exactly using a *t*-distribution we could have noted that with 1406 degrees of freedom the *t* distribution is essentially standard normal and 3.8064 is almost 4 standard deviations. So

$$P(|t| \geq 3.8064) \approx P(|z| \geq 3.8064) < .001$$

We assumed the data was normal and that the two groups had equal variances. Given the large difference between the sample variances this assumption may not be warranted.

In fact, there are other significance tests that test whether the data is approximately normal and whether the two groups have the same variance. In practice one might apply these first to determine whether a *t* test is appropriate in the first place. We don't have time to go into normality tests here, but we will see the *F* distribution used for equality of variances next week.

http://en.wikipedia.org/wiki/Normality_test

http://en.wikipedia.org/wiki/F-test_of_equality_of_variances

MIT OpenCourseWare
<https://ocw.mit.edu>

18.05 Introduction to Probability and Statistics
Spring 2014

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.

Null Hypothesis Significance Testing III

Class 19, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Given hypotheses and data, be able to identify an appropriate significance test from a list of common ones.
2. Given hypotheses, data, and a suggested significance test, know how to look up details and apply the significance test.

2 Introduction

In these notes we will collect together some of the most common significance tests, though by necessity we will leave out many other useful ones. Still, all significance tests follow the same basic pattern in their design and implementation, so by learning the ones we include you should be able to easily apply other ones as needed.

Designing a null hypothesis significance test (NHST):

- Specify null and alternative hypotheses.
- Choose a test statistic whose null distribution and alternative distribution(s) are known.
- Specify a rejection region. Most often this is done implicitly by specifying a significance level α and a method for computing p -values based on the tails of the null distribution.
- Compute power using the alternative distribution(s).

Running a NHST:

- Collect data and compute the test statistic.
- Check if the test statistic is in the rejection region. Most often this is done implicitly by checking if $p < \alpha$. If so, we ‘reject the null hypothesis in favor of the alternative hypothesis’. Otherwise we conclude ‘the data does not support rejecting the null hypothesis’.

Note the careful phrasing: when we fail to reject H_0 , we do not conclude that H_0 is true. The failure to reject may have other causes. For example, we might not have enough data to clearly distinguish H_0 and H_A , whereas more data would indicate that we should reject H_0 .

3 Population parameters and sample statistics

Example 1. If we randomly select 10 men from a population and measure their heights we say we have [sampled the heights](#) from the population. In this case the [sample mean](#), say \bar{x} , is the mean of the sampled heights. It is a statistic and we know its value explicitly. On the other hand, the true average height of the population, say μ , is unknown and we can only estimate its value. We call μ a [population parameter](#).

The main purpose of significance testing is to use sample statistics to draw conclusions about population parameters. For example, we might test if the average height of men in a given population is greater than 70 inches.

4 A gallery of common significance tests related to the normal distribution

We will show a number of tests that all assume [normal data](#). For completeness we will include the z and t tests we've already explored.

You shouldn't try to memorize these tests. It is a hopeless task to memorize the tests given here and even more hopeless to memorize all the tests we've left out. Rather, your goal should be to be able to find the correct test when you need it. Pay attention to the types of hypotheses the tests are designed to distinguish and the assumptions about the data needed for the test to be valid. We will work through the details of these tests in class and on homework.

The null distributions for all of these tests are [all related to the normal distribution](#) by explicit formulas. We will not go into the details of these distributions or the arguments showing how they arise as the null distributions in our significance tests. However, the arguments are accessible to anyone who knows calculus and is interested in understanding them. Given the name of any distribution, you can easily look up the details of its construction and properties online. You can also use R to explore the distribution numerically and graphically.

When analyzing data with any of these tests one thing of key importance is to verify that the assumptions are true or at least approximately true. For example, you shouldn't use a test that assumes the data is normal unless you've checked that the data is approximately normal.

The script [class19.r](#) contains examples of using R to run some of these tests. It is posted in our usual place for R code.

4.1 z -test

- Use: Test if the population mean equals a hypothesized mean.
- Data: x_1, x_2, \dots, x_n .
- Assumptions: The data are independent normal samples:

$$x_i \sim N(\mu, \sigma^2) \text{ where } \mu \text{ is unknown, but } \sigma \text{ is known.}$$
- H_0 : For a specified μ_0 , $\mu = \mu_0$.

- H_A :
 - Two-sided: $\mu \neq \mu_0$
 - one-sided-greater: $\mu > \mu_0$
 - one-sided-less: $\mu < \mu_0$
- Test statistic: $z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$
- Null distribution: $f(z | H_0)$ is the pdf of $Z \sim N(0, 1)$.
- p -value:
 - Two-sided: $p = P(|Z| > z) = 2*(1 - pnorm(abs(z), 0, 1))$
 - one-sided-greater: $p = P(Z > z) = 1 - pnorm(z, 0, 1)$
 - one-sided-less: $p = P(Z < z) = pnorm(z, 0, 1)$
- R code: There does not seem to be a single R function to run a z -test. Of course it is easy enough to get R to compute the z score and p -value.

Example 2. We quickly reprise our example from the class 17 notes.

IQ is normally distributed in the population according to a $N(100, 15^2)$ distribution. We suspect that most MIT students have above average IQ so we frame the following hypotheses.

- H_0 = MIT student IQs are distributed identically to the general population
= MIT IQ's follow a $N(100, 15^2)$ distribution.
- H_A = MIT student IQs tend to be higher than those of the general population
= the average MIT student IQ is greater than 100.

Notice that H_A is one-sided.

Suppose we test 9 students and find they have an average IQ of $\bar{x} = 112$. Can we reject H_0 at a significance level $\alpha = 0.05$?

answer: Our test statistic is

$$z = \frac{\bar{x} - 100}{15/\sqrt{9}} = \frac{36}{15} = 2.4.$$

The right-sided p -value is thereofre

$$p = P(Z \geq 2.4) = 1 - pnorm(2.4, 0, 1) = 0.0081975.$$

Since $p \leq \alpha$ we reject the null hypothesis in favor of the alternative hypothesis that MIT students have higher IQs on average.

4.2 One-sample t -test of the mean

- Use: Test if the population mean equals a hypothesized mean.
- Data: x_1, x_2, \dots, x_n .
- Assumptions: The data are independent normal samples:
 $x_i \sim N(\mu, \sigma^2)$ where both μ and σ are unknown.
- H_0 : For a specified μ_0 , $\mu = \mu_0$

- H_A :

Two-sided: $\mu = \mu_0$

one-sided-greater: $\mu > \mu_0$

one-sided-less: $\mu < \mu_0$

- Test statistic: $t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$,

where s^2 is the sample variance: $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

- Null distribution: $f(t | H_0)$ is the pdf of $T \sim t(n-1)$.

(Student t -distribution with $n-1$ degrees of freedom)

- p -value:

Two-sided: $p = P(|T| > t) = 2 * (1 - pt(abs(t), n-1))$

one-sided-greater: $p = P(T > t) = 1 - pt(t, n-1)$

one-sided-less: $p = P(T < t) = pt(t, n-1)$

- R code example: For data $x = 1, 3, 5, 7, 2$ we can run a one-sample t -test with H_0 :

$\mu = 2.5$ using the R command:

```
t.test(x, mu = mu0, alternative=two.sided=TRUE)
```

This will return a several pieces of information including the mean of the data, t -value and the two-sided p -value. See the help for this function for other argument settings.

Example 3. Look in the class 18 notes or slides for an example of this test. The class 19 example R code also gives an example.

4.3 Two-sample t -test for comparing means

4.3.1 The case of equal variances

We start by describing the test [assuming equal variances](#).

- Use: Test if the population means from two populations differ by a hypothesized amount.
- Data: x_1, x_2, \dots, x_n and y_1, y_2, \dots, y_m .
- Assumptions: Both groups of data are independent normal samples:

$$\begin{aligned} x_i &\sim N(\mu_x, \sigma^2) \\ y_j &\sim N(\mu_y, \sigma^2) \end{aligned}$$

where both μ_x and μ_y are unknown and possibly different. The variance σ^2 is unknown, but the same for both groups.

- H_0 : For a specified μ_0 : $\mu_x - \mu_y = \mu_0$

- H_A :

Two-sided: $\mu_x - \mu_y = \mu_0$

one-sided-greater: $\mu_x - \mu_y > \mu_0$

one-sided-less: $\mu_x - \mu_y < \mu_0$

- Test statistic: $t = \frac{\bar{x} - \bar{y} - \mu_0}{s_p}$,
where s_x^2 and s_y^2 are the sample variances of the x and y data respectively, and s_p^2 is (sometimes called) the pooled sample variance:

$$s_p^2 = \frac{(n-1)s_x^2 + (m-1)s_y^2}{n+m-2} \left(\frac{1}{n} + \frac{1}{m} \right) \quad \text{and} \quad df = n+m-2$$

- Null distribution: $f(t | H_0)$ is the pdf of $T \sim t(df)$, the t -distribution with $df = n+m-2$ degrees of freedom.
- p -value:

Two-sided:	$p = P(T > t)$	=	$2*(1-pt(abs(t), df))$
one-sided-greater:	$p = P(T > t)$	=	$1 - pt(t, df)$
one-sided-less:	$p = P(T < t)$	=	$pt(t, df)$
- R code: The R function `t.test` will run a two-sample t -test. See the example code in `class19.r`

Example 4. Look in the class 18 notes or slides for an example of the two-sample t -test.

Notes: 1. Most often the test is done with $\mu_0 = 0$. That is, the null hypothesis is the [the means are equal](#), i.e. $\mu_x - \mu_y = 0$.

2. If the x and y data have the same length, n , then the formula for s_p^2 becomes simpler:

$$s_p^2 = \frac{s_x^2 + s_y^2}{n}$$

4.3.2 The case of unequal variances

There is a form of the t -test for [when the variances are not assumed equal](#). It is sometimes called [Welch's \$t\$ -test](#).

This looks exactly the same as the case of equal except for a small change in the assumptions and the formula for the pooled variance:

- Use: Test if the population means from two populations differ by a hypothesized amount.
- Data: x_1, x_2, \dots, x_n and y_1, y_2, \dots, y_m .
- Assumptions: Both groups of data are independent normal samples:

$$\begin{aligned} x_i &\sim N(\mu_x, \sigma_x^2) \\ y_j &\sim N(\mu_y, \sigma_y^2) \end{aligned}$$

where both μ_x and μ_y are unknown and possibly different. The variances σ_x^2 and σ_y^2 are unknown and [not assumed to be equal](#).

- H_0, H_A : Exactly the same as the case of equal variances.
- Test statistic: $t = \frac{\bar{x} - \bar{y} - \mu_0}{s_p}$,
where s_x^2 and s_y^2 are the sample variances of the x and y data respectively, and s_p^2

is (sometimes called) the pooled sample variance:

$$s_p^2 = \frac{s_x^2}{n} + \frac{s_y^2}{m} \quad \text{and} \quad df = \frac{(s_x^2/n + s_y^2/m)^2}{(s_x^2/n)^2/(n-1) + (s_y^2/m)^2/(m-1)}$$

- Null distribution: $f(t | H_0)$ is the pdf of $T \sim t(df)$, the t distribution with df degrees of freedom.
- p -value: Exactly the same as the case of equal variances.
- R code: The function `t.test` also handles this case by setting the argument `var.equal=FALSE`.

4.3.3 The paired two-sample t -test

When the data naturally comes in pairs (x_i, y_i) , we can use the [paired two-sample \$t\$ -test](#). (After checking the assumptions are valid!)

Example 5. To measure the effectiveness of a cholesterol lowering medication we might test each subject before and after treatment with the drug. So for each subject we have a pair of measurements: x_i = cholesterol level before treatment and y_i = cholesterol level after treatment.

Example 6. To measure the effectiveness of a cancer treatment we might pair each subject who received the treatment with one who did not. In this case we would want to pair subjects who are similar in terms of stage of the disease, age, sex, etc.

- Use: Test if the average difference between paired values in a population equals a hypothesized value.
 - Data: x_1, x_2, \dots, x_n and y_1, y_2, \dots, y_n must have the same length.
 - Assumptions: The [differences](#) $w_i = x_i - y_i$ between the paired samples are independent draws from a normal distribution $N(\mu, \sigma^2)$, where μ and σ are unknown.
 - **NOTE:** This is just a one-sample t -test using w_i .
 - H_0 : For a specified μ_0 , $\mu = \mu_0$.
 - H_A :
 - Two-sided: $\mu \neq \mu_0$
 - one-sided-greater: $\mu > \mu_0$
 - one-sided-less: $\mu < \mu_0$
 - Test statistic: $t = \frac{\bar{w} - \mu_0}{s/\sqrt{n}}$,
- where s^2 is the sample variance: $s^2 = \frac{1}{n-1} \sum_{i=1}^n (w_i - \bar{w})^2$
- Null distribution: $f(t | H_0)$ is the pdf of $T \sim t(n-1)$.
(Student t -distribution with $n-1$ degrees of freedom)
 - p -value:

Two-sided:	$p = P(T > t) = 2*(1-pt(abs(t), n-1))$
one-sided-greater:	$p = P(T > t) = 1 - pt(t, n-1)$
one-sided-less:	$p = P(T < t) = pt(t, n-1)$

- R code: The R function `t.test` will do a paired two-sample test if you set the argument `paired=TRUE`. You can also run a one-sample *t*-test on $x - y$. There are examples of both of these in `class19.r`

Example 7. The following example is taken from Rice ¹

To study the effect of cigarette smoking on platelet aggregation Levine (1973) drew blood samples from 11 subjects before and after they smoked a cigarette and measured the extent to which platelets aggregated. Here is the data:

Before	25	25	27	44	30	67	53	53	52	60	28
After	27	29	37	56	46	82	57	80	61	59	43
Difference	2	4	10	12	16	15	4	27	9	-1	15

The null hypothesis is that smoking had no effect on platelet aggregation, i.e. that the difference should have mean $\mu_0 = 0$. We ran a paired two-sample *t*-test to test this hypothesis. Here is the R code: (It's also in `class19.r`.)

```
before.cig = c(25,25,27,44,30,67,53,53,52,60,28)
after.cig = c(27,29,37,56,46,82,57,80,61,59,43)
mu0 = 0
result = t.test(after.cig, before.cig, alternative="two.sided", mu=mu0, paired=TRUE)
print(result)
```

Here is the output:

```
Paired t-test
data: after.cig and before.cig
t = 4.2716, df = 10, p-value = 0.001633
alternative hypothesis: true difference in means is not equal to 0
mean of the differences: 10.27273
```

We got the same results with the one-sample *t*-test:

```
t.test(after.cig - before.cig, mu=0)
```

4.4 One-way ANOVA (*F*-test for equal means)

- Use: Test if the population means from n groups are all the same.
- Data: (n groups, m samples from each group)

$$\begin{array}{cccc} x_{1,1}, & x_{1,2}, & \dots, & x_{1,m} \\ x_{2,1}, & x_{2,2}, & \dots, & x_{2,m} \\ \dots \\ x_{n,1}, & x_{n,2}, & \dots, & x_{n,m} \end{array}$$

- Assumptions: Data for each group is an independent normal sample drawn from distributions with (possibly) different means but [the same variance](#):

$$\begin{array}{ll} x_{1,j} & \sim N(\mu_1, \sigma^2) \\ x_{2,j} & \sim N(\mu_2, \sigma^2) \\ \dots \\ x_{n,j} & \sim N(\mu_n, \sigma^2) \end{array}$$

¹John Rice, *Mathematical Statistics and Data Analysis*, 2nd edition, p. 412. This example references P.H Levine (1973) An acute effect of cigarette smoking on platelet function. *Circulation*, 48, 619-623.

The group means μ_i are unknown and possibly different. The variance σ is unknown, but the same for all groups.

- H_0 : All the means are identical $\mu_1 = \mu_2 = \dots = \mu_n$.

- H_A : Not all the means are the same.

- Test statistic: $w = \frac{MS_B}{MS_W}$, where

$$\begin{aligned}\bar{x}_i &= \text{mean of group } i \\ &= \frac{x_{i,1} + x_{i,2} + \dots + x_{i,m}}{m}.\end{aligned}$$

$$\bar{x} = \text{grand mean of all the data.}$$

$$\begin{aligned}s_i^2 &= \text{sample variance of group } i \\ &= \frac{1}{m-1} \sum_{j=1}^m (x_{i,j} - \bar{x}_i)^2.\end{aligned}$$

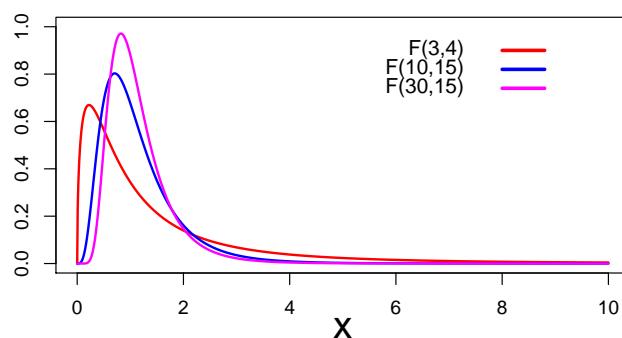
$$MS_B = \text{between group variance}$$

$$\begin{aligned}&= m \times \text{sample variance of group means} \\ &= \frac{m}{n-1} \sum_{i=1}^n (\bar{x}_i - \bar{x})^2.\end{aligned}$$

$$MS_W = \text{average within group variance}$$

$$\begin{aligned}&= \text{sample mean of } s_1^2, \dots, s_n^2 \\ &= \frac{s_1^2 + s_2^2 + \dots + s_n^2}{n}\end{aligned}$$

- Idea: If the μ_i are all equal, this ratio should be near 1. If they are not equal then MS_B should be larger while MS_W should remain about the same, so w should be larger. We won't give a proof of this.
- Null distribution: $f(w | H_0)$ is the pdf of $W \sim F(n-1, n(m-1))$. This is the F -distribution with $(n-1)$ and $n(m-1)$ degrees of freedom. Several F -distributions are plotted below.
- p -value: $p = P(W > w) = 1 - \text{pf}(w, n-1, n*(m-1))$



Notes: 1. ANOVA tests whether all the means are the same. It does not test whether some subset of the means are the same.

2. There is a test where the variances are not assumed equal.

3. There is a test where the groups don't all have the same number of samples.

4. R has a function `aov()` to run ANOVA tests. See:

https://personality-project.org/r/r_guide/r.anova.html#oneway

<http://en.wikipedia.org/wiki/F-test>

Example 8. The table shows patients' perceived level of pain (on a scale of 1 to 6) after 3 different medical procedures.

T_1	T_2	T_3
2	3	2
4	4	1
1	6	3
5	1	3
3	4	5

(1) Set up and run an F-test comparing the means of these 3 treatments.

(2) Based on the test, what might you conclude about the treatments?

answer: Using the code below, the F statistic is 0.325 and the p -value is 0.729 At any reasonable significance level we will fail to reject the null hypothesis that the average pain level is the same for all three treatments..

Note, it is not reasonable to conclude the the null hypothesis is true. With just 5 data points per procedure we might simply lack the power to distinguish different means.

R code to perform the test

```
# DATA ----
T1 = c(2,4,1,5,3)
T2 = c(3,4,6,1,4)
T3 = c(2,1,3,3,5)

procedure = c(rep('T1',length(T1)),rep('T2',length(T2)),rep('T3',length(T3)))
pain = c(T1,T2,T3)
data.pain = data.frame(procedure,pain)
aov.data = aov(pain~procedure,data=data.pain) # do the analysis of variance
print(summary(aov.data)) # show the summary table

# class19.r also show code to compute the ANOVA by hand.
```

The summary shows a p -value (shown as $\text{Pr}(>F)$) of 0.729. Therefore we do not reject the null hypothesis that all three group population means are the same.

4.5 Chi-square test for goodness of fit

This is a test of how well a hypothesized probability distribution fits a set of data. The test statistic is called a **chi-square statistic** and the null distribution associated of the chi-square statistic is the **chi-square distribution**. It is denoted by $\chi^2(df)$ where the parameter df is called the **degrees of freedom**.

Suppose we have an unknown probability mass function given by the following table.

Outcomes	ω_1	ω_2	\dots	ω_n
Probabilities	p_1	p_2	\dots	p_n

In the chi-square test for goodness of fit we hypothesize a set of values for the probabilities. Typically we will hypothesize that the probabilities follow a known distribution with certain parameters, e.g. binomial, Poisson, multinomial. The test then tries to determine if this set of probabilities could have reasonably generated the data we collected.

- Use: Test whether discrete data fits a specific finite probability mass function.
- Data: An observed count O_i for each possible outcome ω_i .
- Assumptions: None
- H_0 : The data was drawn from a specific discrete distribution.
- H_A : The data was drawn from a different distribution.
- Test statistic: The data consists of observed counts O_i for each ω_i . From the null hypothesis probability table we get a set of expected counts E_i . There are two statistics that we can use:

$$\text{Likelihood ratio statistic } G = 2 * \sum O_i \ln \left(\frac{O_i}{E_i} \right)$$

$$\text{Pearson's chi-square statistic } X^2 = \sum \frac{(O_i - E_i)^2}{E_i}.$$

It is a theorem that under the null hypothesis $X^2 \approx G$ and both are approximately chi-square. Before computers, X^2 was used because it was easier to compute. Now, it is better to use G although you will still see X^2 used quite often.

- Degrees of freedom df : For chi-square tests the number of degrees of freedom can be a bit tricky. In this case $df = n - 1$. It is computed as the number of cell counts that can be freely set under H_A consistent with the statistics needed to compute the expected cell counts assuming H_0 .
- Null distribution: Assuming H_0 , both statistics (approximately) follow a chi-square distribution with df degrees of freedom. That is both $f(G | H_0)$ and $f(X^2 | H_0)$ have the same pdf as $Y \sim \chi^2(df)$.
- p -value:

$$p = P(Y > G) = 1 - \text{pchisq}(G, df)$$

$$p = P(Y > X^2) = 1 - \text{pchisq}(X^2, df)$$
- R code: The R function `chisq.test` can be used to do the computations for a chi-square test use X^2 . For G you either have to do it by hand or find a package that has a function. (It will probably be called `likelihood.test` or `G.test`.)

Notes. 1. When the likelihood ratio statistic G is used the test is also called a [G-test](#) or a [likelihood ratio test](#).

Example 9. First chi-square example. Suppose we have an experiment that produces numerical data. For this experiment the possible outcomes are 0, 1, 2, 3, 4, 5 or more. We run 51 trials and count the frequency of each outcome, getting the following data:

Outcomes	0	1	2	3	4	≥ 5
Observed counts	3	10	15	13	7	3

Suppose our null hypothesis H_0 is that the data is drawn from 51 trials of a binomial(8, 0.5) distribution and our alternative hypothesis H_A is that the data is drawn from some other distribution. Do all of the following:

1. Make a table of the observed and expected counts.
2. Compute both the likelihood ratio statistic G and Pearson's chi-square statistic X^2 .

3. Compute the degrees of freedom of the null distribution.

4. Compute the p -values corresponding to G and X^2 .

answer: All of the R code used for this example is in class19.r.

1. Assuming H_0 the data truly comes from a binomial(8, 0.5) distribution. We have 51 total observations, so the expected count for each outcome is just 51 times its probability. We computed the binomial(8, 0.5) probabilities and expected counts in R:

Outcomes	0	1	2	3	4	≥ 5
Observed counts	3	10	15	13	7	3
H_0 probabilities	0.0039	0.0313	0.1094	0.2188	0.2734	0.3633
Expected counts	0.19	1.53	5.36	10.72	13.40	17.80

2. Using the formulas above we compute that $X^2 = 116.41$ and $G = 66.08$

3. The only statistic used in computing the expected counts was the total number of observations 51. So, the degrees of freedom is 5, i.e we can set 5 of the cell counts freely and the last is determined by requiring that the total number is 51.

4. The p -values are $pG = 1 - \text{pchisq}(G, 5)$ and $pX2 = 1 - \text{pchisq}(X^2, 5)$. Both p -values are effectively 0. For almost any significance level we would reject H_0 in favor of H_A .

Example 10. ([Degrees of freedom.](#)) Suppose we have the same data as in the previous example, but our null hypothesis is that the data comes from independent trials of binomial($8, \theta$) distribution, where θ can be anything. (H_A is that the data comes from some other distribution.) In this case we must estimate θ from the data, e.g. using the MLE. In total we have computed two values from the data: the total number of counts and the estimate of θ . So, the degrees of freedom is $6 - 2 = 4$.

Example 11. Mendel's genetic experiments (Adapted from Rice *Mathematical Statistics and Data Analysis*, 2nd ed., example C, p.314)

In one of his experiments on peas Mendel crossed 556 smooth, yellow male peas with wrinkled green female peas. Assuming the smooth and wrinkled genes occur with equal frequency we'd expect 1/4 of the pea population to have two smooth genes (SS), 1/4 to have two wrinkled genes (ss), and the remaining 1/2 would be heterozygous Ss . We also expect these fractions for yellow (Y) and green (y) genes. If the color and smoothness genes are inherited independently and smooth and yellow are both dominant we'd expect the following table of frequencies for phenotypes.

	Yellow	Green	
Smooth	9/16	3/16	3/4
Wrinkled	3/16	1/16	1/4
	3/4	1/4	1

Probability table for the null hypothesis

So from the 556 crosses the expected number of smooth yellow peas is $556 \times 9/16 = 312.75$. Likewise for the other possibilities. Here is a table giving the observed and expected counts from Mendel's experiments.

	Observed count	Expected count
Smooth yellow	315	312.75
Smooth green	108	104.25
Wrinkled yellow	102	104.25
Wrinkled green	31	34.75

The null hypothesis is that the observed counts are random samples distributed according to the frequency table given above. We use the counts to compute our statistics

The likelihood ratio statistic is

$$\begin{aligned}
 G &= 2 * \sum O_i \ln \left(\frac{O_i}{E_i} \right) \\
 &= 2 * \left(315 \ln \left(\frac{315}{412.75} \right) + 108 \ln \left(\frac{108}{104.25} \right) + 102 \ln \left(\frac{102}{104.25} \right) + 31 \ln \left(\frac{31}{34.75} \right) \right) \\
 &= 0.618
 \end{aligned}$$

Pearson's chi-square statistic is

$$X^2 = \sum \frac{(O_i - E_i)^2}{E_i} = \frac{2.75}{312.75} + \frac{3.75}{104.25} + \frac{2.25}{104.25} + \frac{3.75}{34.75} = 0.604$$

You can see that the two statistics are very close. This is usually the case. In general the likelihood ratio statistic is more robust and should be preferred.

The degrees of freedom is 3, because there are 4 observed quantities and one relation between them, i.e. they sum to 556. So, under the null hypothesis G follows a $\chi^2(3)$ distribution. Using R to compute the p -value we get

$$p = 1 - \text{pchisq}(0.618, 3) = 0.892$$

Assuming the null hypothesis we would see data at least this extreme almost 90% of the time. We would not reject the null hypothesis for any reasonable significance level.

The p -value using Pearson's statistic is 0.985 –nearly identical.

The script class19.r shows these calculations and also how to use `chisq.test` to run a chi-square test directly.

4.6 Chi-square test for homogeneity

This is a test to see if several independent sets of random data are all drawn from the same distribution. (The meaning of homogeneity in this case is that all the distributions are the same.)

- Use: Test whether m different independent sets of discrete data are drawn from the same distribution.
- Outcomes: $\omega_1, \omega_2, \dots, \omega_n$ are the possible outcomes. These are the same for each set of data.
- Data: We assume m independent sets of data giving counts for each of the possible outcomes. That is, for data set i we have an observed count $O_{i,j}$ for each possible outcome ω_j .

- Assumptions: None
- H_0 : Each data set is drawn from the same distribution. (We don't specify what this distribution is.)
- H_A : The data sets are not all drawn from the same distribution.
- Test statistic: See the example below. There are mn cells containing counts for each outcome for each data set. Using the null distribution we can estimate expected counts for each of the data sets. The statistics X^2 and G are computed exactly as above.
- Degrees of freedom df : $(m - 1)(n - 1)$. (See the example below.)
- The null distribution $\chi^2(df)$. The p -values are computed just as in the chi-square test for goodness of fit.
- R code: The R function `chisq.test` can be used to do the computations for a chi-square test use X^2 . For G you either have to do it by hand or find a package that has a function. (It will probably be called `likelihood.test` or `G.test`.)

Example 12. Someone claims to have found a long lost work by William Shakespeare. She asks you to test whether or not the play was actually written by Shakespeare .

You go to <http://www.opensourceshakespeare.org> and pick a random 12 pages from *King Lear* and count the use of common words. You do the same thing for the ‘long lost work’. You get the following table of counts.

Word	a	an	this	that
<i>King Lear</i>	150	30	30	90
Long lost work	90	20	10	80

Using this data, set up and evaluate a significance test of the claim that the long lost book is by William Shakespeare. Use a significance level of 0.1.

answer: The null hypothesis H_0 : For the 4 words counted the long lost book has the same relative frequencies as the counts taken from *King Lear*.

The total word count of both books combined is 500, so the the maximum likelihood estimate of the relative frequencies assuming H_0 is simply the total count for each word divided by the total word count.

Word	a	an	this	that	Total count
<i>King Lear</i>	150	30	30	90	300
Long lost work	90	20	10	80	200
totals	240	50	40	170	500
rel. frequencies under H_0	240/500	50/500	40/500	170/500	500/500

Now the expected counts for each book under H_0 are the total count for that book times the relative frequencies in the above table. The following table gives the counts: (observed, expected) for each book.

Word	a	an	this	that	Totals
<i>King Lear</i>	(150, 144)	(30, 30)	(30, 24)	(90, 102)	(300, 300)
Long lost work	(90, 96)	(20, 20)	(10, 16)	(80, 68)	(200, 200)
Totals	(249, 240)	(50, 50)	(40, 40)	(170, 170)	(500, 500)

The chi-square statistic is

$$\begin{aligned} X^2 &= \sum \frac{(O_i - E_i)^2}{E_i} \\ &= \frac{6^2}{144} + \frac{0^2}{30} + \frac{6^2}{24} + \frac{12^2}{102} + \frac{6^2}{96} + \frac{0^2}{20} + \frac{6^2}{16} + \frac{12^2}{68} \\ &\approx 7.9 \end{aligned}$$

There are 8 cells and all the marginal counts are fixed because they were needed to determine the expected counts. To be consistent with these statistics we could freely set the values in 3 cells in the table, e.g. the 3 blue cells, then the rest of the cells are determined in order to make the marginal totals correct. Thus $df = 3$. (Or we could recall that $df = (m - 1)(n - 1) = (3)(1) = 3$, where m is the number of columns and n is the number of rows.)

Using R we find $p = 1 - \text{pchisq}(7.9, 3) = 0.048$. Since this is less than our significance level of 0.1 we reject the null hypothesis that the relative frequencies of the words are the same in both books.

If we make the further assumption that all of Shakespeare's plays have similar word frequencies (which is something we could check) we conclude that the book is probably not by Shakespeare.

4.7 Other tests

There are far too many other tests to even make a dent. We will see some of them in class and on psets. Again, we urge you to master the paradigm of NHST and recognize the importance of choosing a test statistic with a known null distribution.

MIT OpenCourseWare
<https://ocw.mit.edu>

18.05 Introduction to Probability and Statistics
Spring 2014

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.

Comparison of frequentist and Bayesian inference.
Class 20, 18.05
Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to explain the difference between the p -value and a posterior probability to a doctor.

2 Introduction

We have now learned about two schools of statistical inference: Bayesian and frequentist. Both approaches allow one to evaluate evidence about competing hypotheses. In these notes we will review and compare the two approaches, starting from Bayes' formula.

3 Bayes' formula as touchstone

In our first unit (probability) we learned Bayes' formula, a perfectly abstract statement about conditional probabilities of events:

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}.$$

We began our second unit (Bayesian inference) by reinterpreting the events in Bayes' formula:

$$P(\mathcal{H} | \mathcal{D}) = \frac{P(\mathcal{D} | \mathcal{H})P(\mathcal{H})}{P(\mathcal{D})}.$$

Now \mathcal{H} is a hypothesis and \mathcal{D} is data which may give evidence for or against \mathcal{H} . Each term in Bayes' formula has a name and a role.

- The prior $P(\mathcal{H})$ is the probability that \mathcal{H} is true before the data is considered.
- The posterior $P(\mathcal{H} | \mathcal{D})$ is the probability that \mathcal{H} is true after the data is considered.
- The likelihood $P(\mathcal{D} | \mathcal{H})$ is the evidence about \mathcal{H} provided by the data \mathcal{D} .
- $P(\mathcal{D})$ is the total probability of the data taking into account all possible hypotheses.

If the prior and likelihood are known for all hypotheses, then Bayes' formula computes the posterior exactly. Such was the case when we rolled a die randomly selected from a cup whose contents you knew. We call this the deductive logic of probability theory, and it gives a direct way to compare hypotheses, draw conclusions, and make decisions.

In most experiments, the prior probabilities on hypotheses are not known. In this case, our recourse is the art of statistical inference: we either make up a prior (Bayesian) or do our best using only the likelihood (frequentist).

The Bayesian school models uncertainty by a probability distribution over hypotheses. One's ability to make inferences depends on one's degree of confidence in the chosen prior, and the robustness of the findings to alternate prior distributions may be relevant and important.

The frequentist school only uses conditional distributions of data given specific hypotheses. The presumption is that some hypothesis (parameter specifying the conditional distribution of the data) is true and that the observed data is sampled from that distribution. In particular, the frequentist approach does not depend on a subjective prior that may vary from one investigator to another.

These two schools may be further contrasted as follows:

Bayesian inference

- uses probabilities for both hypotheses and data.
- depends on the prior and likelihood of observed data.
- requires one to know or construct a ‘subjective prior’.
- dominated statistical practice before the 20th century.
- may be computationally intensive due to integration over many parameters.

Frequentist inference (NHST)

- never uses or gives the probability of a hypothesis (no prior or posterior).
- depends on the likelihood $P(\mathcal{D} | \mathcal{H})$ for both observed and unobserved data.
- does not require a prior.
- dominated statistical practice during the 20th century.
- tends to be less computationally intensive.

Frequentist measures like p -values and confidence intervals continue to dominate research, especially in the life sciences. However, in the current era of powerful computers and big data, Bayesian methods have undergone an enormous renaissance in fields like machine learning and genetics. There are now a number of large, ongoing clinical trials using Bayesian protocols, something that would have been hard to imagine a generation ago. While professional divisions remain, the consensus forming among top statisticians is that the most effective approaches to complex problems often draw on the best insights from both schools working in concert.

4 Critiques and defenses

4.1 Critique of Bayesian inference

1. The main critique of Bayesian inference is that a subjective prior is, well, subjective. There is no single method for choosing a prior, so different people will produce different priors and may therefore arrive at different posteriors and conclusions.

2. Furthermore, there are philosophical objections to assigning probabilities to hypotheses, as hypotheses do not constitute outcomes of repeatable experiments in which one can measure long-term frequency. Rather, a hypothesis is either true or false, regardless of whether one knows which is the case. A coin is either fair or unfair; treatment 1 is either better or worse than treatment 2; the sun will or will not come up tomorrow.

4.2 Defense of Bayesian inference

1. The probability of hypotheses is exactly what we need to make decisions. When the doctor tells me a screening test came back positive I want to know what is the probability this means I'm sick. That is, I want to know the probability of the hypothesis "I'm sick".
2. Using Bayes' theorem is logically rigorous. Once we have a prior all our calculations have the certainty of deductive logic.
3. By trying different priors we can see how sensitive our results are to the choice of prior.
4. It is easy to communicate a result framed in terms of probabilities of hypotheses.
5. Even though the prior may be subjective, one can specify the assumptions used to arrive at it, which allows other people to challenge it or try other priors.
6. The evidence derived from the data is independent of notions about 'data more extreme' that depend on the exact experimental setup (see the "Stopping rules" section below).
7. Data can be used as it comes in. There is no requirement that every contingency be planned for ahead of time.

4.3 Critique of frequentist inference

1. It is ad-hoc and does not carry the force of deductive logic. Notions like 'data more extreme' are not well defined. The p -value depends on the exact experimental setup (see the "Stopping rules" section below).
2. Experiments must be fully specified ahead of time. This can lead to paradoxical seeming results. See the 'voltmeter story' in:
http://en.wikipedia.org/wiki/Likelihood_principle
3. The p -value and significance level are notoriously prone to misinterpretation. Careful statisticians know that a significance level of 0.05 means the probability of a type I error is 5%. That is, if the null hypothesis is true then 5% of the time it will be rejected due to randomness. Many (most) other people erroneously think a p -value of 0.05 means that the probability of the null hypothesis is 5%.

Strictly speaking you could argue that this is not a critique of frequentist inference but, rather, a critique of popular ignorance. Still, the subtlety of the ideas certainly contributes to the problem. (see "Mind your p 's" below).

4.4 Defense of frequentist inference

1. It is objective: all statisticians will agree on the p -value. Any individual can then decide if the p -value warrants rejecting the null hypothesis.

2. Hypothesis testing using frequentist significance testing is applied in the statistical analysis of scientific investigations, evaluating the strength of evidence against a null hypothesis with data. The interpretation of the results is left to the user of the tests. Different users may apply different significance levels for determining statistical significance. Frequentist statistics does not pretend to provide a way to choose the significance level; rather it explicitly describes the trade-off between type I and type II errors.
3. Frequentist experimental design demands a careful description of the experiment and methods of analysis before starting. This helps control for experimenter bias.
4. The frequentist approach has been used for over 100 years and we have seen tremendous scientific progress. Although the frequentist herself would not put a probability on the belief that frequentist methods are valuable shouldn't this history give the Bayesian a strong prior belief in the utility of frequentist methods?

5 Mind your p 's.

We run a two-sample t -test for equal means, with $\alpha = 0.05$, and obtain a p -value of 0.04. What are the odds that the two samples are drawn from distributions with the same mean?

- (a) 19/1 (b) 1/19 (c) 1/20 (d) 1/24 (e) unknown

answer: (e) unknown. Frequentist methods only give probabilities of statistics conditioned on hypotheses. They do not give probabilities of hypotheses.

6 Stopping rules

When running a series of trials we need a rule on when to stop. Two common rules are:

1. Run exactly n trials and stop.
2. Run trials until you see a certain result and then stop.

In this example we'll consider two coin tossing experiments.

Experiment 1: Toss the coin exactly 6 times and report the number of heads.

Experiment 2: Toss the coin until the first tails and report the number of heads.

Jon is worried that his coin is biased towards heads, so before using it in class he tests it for fairness. He runs an experiment and reports to Jerry that his sequence of tosses was $HHHHHT$. But Jerry is only half-listening, and he forgets which experiment Jon ran to produce the data.

Frequentist approach.

Since he's forgotten which experiment Jon ran, Jerry the frequentist decides to compute the p -values for both experiments given Jon's data.

Let θ be the probability of heads. We have the null and one-sided alternative hypotheses

$$H_0 : \theta = 0.5, \quad H_A : \theta > 0.5.$$

Experiment 1: The null distribution is binomial(6, 0.5) so, the one sided p -value is the probability of 5 or 6 heads in 6 tosses. Using R we get

$$p = 1 - \text{pb}inom(4, 6, 0.5) = 0.1094.$$

Experiment 2: The null distribution is geometric(0.5) so, the one sided p -value is the probability of 5 or more heads before the first tails. Using R we get

$$p = 1 - \text{pgeom}(4, 0.5) = 0.0313.$$

Using the typical significance level of 0.05, the same data leads to opposite conclusions! We would reject H_0 in experiment 2, but not in experiment 1.

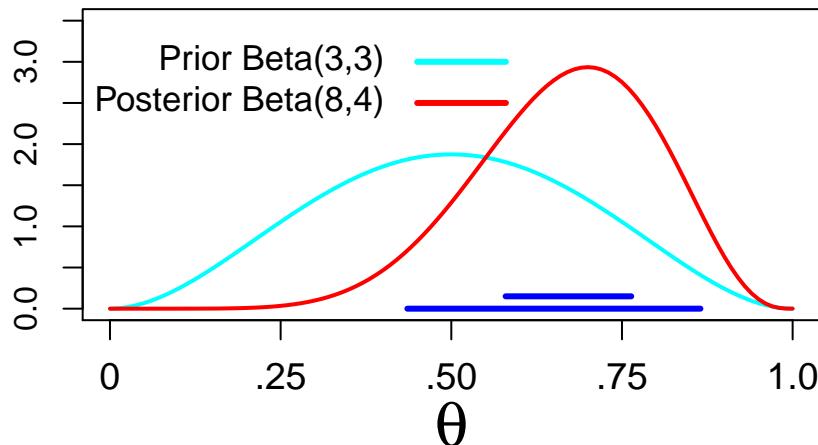
The frequentist is fine with this. The set of possible outcomes is different for the different experiments so the notion of extreme data, and therefore p -value, is different. For example, in experiment 1 we would consider $THHHHHH$ to be as extreme as $HHHHHT$. In experiment 2 we would never see $THHHHHH$ since the experiment would end after the first tails.

Bayesian approach.

Jerry the Bayesian knows it doesn't matter which of the two experiments Jon ran, since the binomial and geometric likelihood functions (columns) for the data $HHHHHT$ are proportional. In either case, he must make up a prior, and he chooses Beta(3,3). This is a relatively flat prior concentrated over the interval $0.25 \leq \theta \leq 0.75$.

See <http://mathlets.org/mathlets/beta-distribution/>

Since the beta and binomial (or geometric) distributions form a conjugate pair the Bayesian update is simple. Data of 5 heads and 1 tails gives a posterior distribution Beta(8,4). Here is a graph of the prior and the posterior. The blue lines at the bottom are 50% and 90% probability intervals for the posterior.



Prior and posterior distributions with 0.5 and 0.9 probability intervals

Here are the relevant computations in R:

Posterior 50% probability interval: `qbeta(c(0.25, 0.75), 8, 4) = [0.58 0.76]`
 Posterior 90% probability interval: `qbeta(c(0.05, 0.95), 8, 4) = [0.44 0.86]`
 $P(\theta > 0.50 | \text{data}) = 1 - \text{pbeta}(0.5, \text{posterior.a}, \text{posterior.b}) = 0.89$

Starting from the prior Beta(3,3), the posterior probability that the coin is biased toward heads is 0.89.

7 Making decisions

Quite often the goal of statistical inference is to help with making a decision, e.g. whether or not to undergo surgery, how much to invest in a stock, whether or not to go to graduate school, etc.

In statistical decision theory, consequences of taking actions are measured by a utility function. The utility function assigns a weight to each possible outcome; in the language of probability, it is simply a random variable.

For example, in my investments I could assign a utility of d to the outcome of a gain of d dollars per share of a stock (if $d < 0$ my utility is negative). On the other hand, if my tolerance for risk is low, I will assign a more negative utility to losses than to gains (say, $-d^2$ if $d < 0$ and d if $d \geq 0$).

A decision rule combines the expected utility with evidence for each hypothesis given by the data (e.g., p -values or posterior distributions) into a formal statistical framework for making decisions.

In this setting, the frequentist will consider the expected utility given a hypothesis

$$E(U | \mathcal{H})$$

where U is the random variable representing utility. There are frequentist methods for combining the expected utility with p -values of hypotheses to guide decisions.

The Bayesian can combine $E(U | \mathcal{H})$ with the posterior (or prior if it's before data is collected) to create a Bayesian decision rule.

In either framework, two people considering the same investment may have different utility functions and make different decisions. For example, a riskier stock (with higher potential upside and downside) will be more appealing with respect to the first utility function above than with respect to the second (loss-averse) one.

A significant theoretical result is that for any decision rule there is a Bayesian decision rule which is, in a precise sense, at least as good a rule.

MIT OpenCourseWare
<https://ocw.mit.edu>

18.05 Introduction to Probability and Statistics
Spring 2014

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.

Confidence intervals based on normal data
Class 22, 18.05
Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to determine whether an expression defines a valid interval statistic.
2. Be able to compute z and t confidence intervals for the mean given normal data.
3. Be able to compute the χ^2 confidence interval for the variance given normal data.
4. Be able to define the confidence level of a confidence interval.
5. Be able to explain the relationship between the z confidence interval (and confidence level) and the z non-rejection region (and significance level) in NHST.

2 Introduction

We continue to survey the tools of frequentist statistics. Suppose we have a model (probability distribution) for observed data with an unknown parameter. We have seen how NHST uses data to test the hypothesis that the unknown parameter has a particular value.

We have also seen how point estimates like the MLE use data to provide an estimate of the unknown parameter. On its own, a point estimate like $\bar{x} = 2.2$ carries no information about its accuracy; it's just a single number, regardless of whether its based on ten data points or one million data points.

For this reason, statisticians augment point estimates with confidence intervals. For example, to estimate an unknown mean μ we might be able to say that our best estimate of the mean is $\bar{x} = 2.2$ with a 95% confidence interval $[1.2, 3.2]$. Another way to describe the interval is: $\bar{x} \pm 1$.

We will leave to later the explanation of exactly what the 95% confidence level means. For now, we'll note that taken together the width of the interval and the confidence level provide a measure on the strength of the evidence supporting the hypothesis that the μ is close to our estimate \bar{x} . You should think of the confidence level of an interval as analogous to the significance level of a NHST. As explained below, it is no accident that we often see significance level $\alpha = 0.05$ and confidence level $0.95 = 1 - \alpha$.

We will first explore confidence intervals in situations where you will easily be able to compute by hand: z and t confidence intervals for the mean and χ^2 confidence intervals for the variance. We will use R to handle all the computations in more complicated cases. Indeed, the challenge with confidence intervals is not their computation, but rather interpreting them correctly and knowing how to use them in practice.

3 Interval statistics

Recall that our working definition of a statistic is anything that can be computed from data. In particular, the formula for a statistic cannot include unknown quantities.

Example 1. Suppose x_1, \dots, x_n is drawn from $N(\mu, \sigma^2)$ where μ and σ are unknown.

- (i) \bar{x} and $\bar{x} - 5$ are statistics.
- (ii) $\bar{x} - \mu$ is not a statistic since μ is unknown.

(iii) If μ_0 a known value, then $\bar{x} - \mu_0$ is a statistic. This case arises when we consider the null hypothesis $\mu = \mu_0$. For example, if the null hypothesis is $\mu = 5$, then the statistic $\bar{x} - \mu_0$ is just $\bar{x} - 5$ from (i).

We can play the same game with intervals to define [interval statistics](#)

Example 2. Suppose x_1, \dots, x_n is drawn from $N(\mu, \sigma^2)$ where μ is unknown.

- (i) The interval $[\bar{x} - 2.2, \bar{x} + 2.2] = \bar{x} \pm 2.2$ is an interval statistic.
- (ii) If σ is [known](#), then $\left[\bar{x} - \frac{2\sigma}{\sqrt{n}}, \bar{x} + \frac{2\sigma}{\sqrt{n}} \right]$ is an interval statistic.
- (iii) On the other hand, if σ is [unknown](#) then $\left[\bar{x} - \frac{2\sigma}{\sqrt{n}}, \bar{x} + \frac{2\sigma}{\sqrt{n}} \right]$ is **not** an interval statistic.
- (iv) If s^2 is the sample variance, then $\left[\bar{x} - \frac{2s}{\sqrt{n}}, \bar{x} + \frac{2s}{\sqrt{n}} \right]$ is an interval statistic because s^2 is computed from the data.

We will return to (ii) and (iv), as these are respectively the z and t confidence intervals for estimating μ .

Technically an interval statistic is nothing more than a pair of point statistics giving the lower and upper bounds of the interval. Our reason for emphasizing that the interval is a statistic is to highlight the following:

1. The interval is random – new random data will produce a new interval.
2. As frequentists we are perfectly happy using it because it doesn't depend on the value of an unknown parameter or hypothesis.
3. As usual with frequentist statistics we have to assume a certain hypothesis, e.g. value of μ , before we can compute probabilities about the interval.

Example 3. Suppose we draw n samples x_1, \dots, x_n from a $N(\mu, 1)$ distribution, where μ is unknown. Suppose we wish to know the probability that 0 is in the interval $[\bar{x} - 2, \bar{x} + 2]$. Without knowing the value of μ this is impossible. However, we can compute this probability for any given (hypothesized) value of μ .

4. [A warning which will be repeated:](#) Be careful in your thinking about these probabilities. Confidence intervals are a frequentist notion. Since frequentists do not compute probabilities of hypotheses, the confidence level is never a probability that the unknown parameter is in the confidence interval.

4 z confidence intervals for the mean

Throughout this section we will assume that we have normally distributed data:

$$x_1, x_2, \dots, x_n \sim N(\mu, \sigma^2).$$

As we often do, we will introduce the main ideas through examples, building on what we know about rejection and non-rejection regions in NHST until we have constructed a confidence interval.

4.1 Definition of z confidence intervals for the mean

We start with z confidence intervals for the mean. First we'll give the formula. Then we'll walk through the derivation in one entirely numerical example. This will give us the basic idea. Then we'll repeat this example, replacing the explicit numbers by symbols. Finally we'll work through a computational example.

Definition: Suppose the data $x_1, \dots, x_n \sim N(\mu, \sigma^2)$, with unknown mean μ and known variance σ^2 . The $(1 - \alpha)$ confidence interval for μ is

$$\left[\bar{x} - \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}}, \bar{x} + \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}} \right], \quad (1)$$

where $z_{\alpha/2}$ is the right critical value $P(Z > z_{\alpha/2}) = \alpha/2$.

For example, if $\alpha = 0.05$ then $z_{\alpha/2} = 1.96$ so the 0.95 (or 95%) confidence interval is

$$\left[\bar{x} - \frac{1.96\sigma}{\sqrt{n}}, \bar{x} + \frac{1.96\sigma}{\sqrt{n}} \right].$$

We've created an applet that generates normal data and displays the corresponding z confidence interval for the mean. It also shows the t -confidence interval, as discussed in the next section. Play around to get a sense for random intervals!

<http://mathlets.org/mathlets/confidence-intervals/>

Example 4. Suppose we collect 100 data points from a $N(\mu, 3^2)$ distribution and the sample mean is $\bar{x} = 12$. Give the 95 % confidence interval for μ .

answer: Using the formula this is trivial to compute: the 95% confidence interval for μ is

$$\left[\bar{x} - \frac{1.96\sigma}{\sqrt{n}}, \bar{x} + \frac{1.96\sigma}{\sqrt{n}} \right] = \left[12 - \frac{1.96 \cdot 3}{\sqrt{100}}, 12 + \frac{1.96 \cdot 3}{\sqrt{100}} \right]$$

4.2 Explaining the definition part 1: rejection regions

Our next goal is to explain the definition (1) starting from our knowledge of rejection/non-rejection regions. The phrase '[non-rejection region](#)' is not pretty, but we will discipline ourselves to use it instead of the inaccurate phrase 'acceptance region'.

Example 5. Suppose that $n = 12$ data points are drawn from $N(\mu, 5^2)$ where μ is unknown. Set up a two-sided significance test of $H_0 : \mu = 2.71$ using the statistic \bar{x} at significance level $\alpha = 0.05$. Describe the rejection and non-rejection regions.

answer: Under the null hypothesis $\mu = 2.71$ we have $x_i \sim N(2.71, 5^2)$ and thus

$$\bar{x} \sim N(2.71, 5^2/12)$$

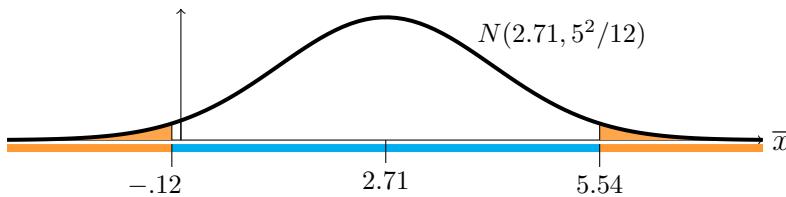
where $5^2/12$ is the variance $(\sigma_{\bar{x}})^2$ of \bar{x} . We know that significance $\alpha = 0.05$ corresponds to a rejection region outside 1.96 standard deviations from the hypothesized mean. That is, the non-rejection and rejection regions are separated by the critical values $\bar{x} \pm 1.96 \sigma_{\bar{x}}$.

Non-rejection region:

$$\left[2.71 - \frac{1.96 \cdot 5}{\sqrt{12}}, \quad 2.71 + \frac{1.96 \cdot 5}{\sqrt{12}} \right] = [-0.12, 5.54].$$

$$\text{Rejection region: } \left(-\infty, 2.71 - \frac{1.96 \cdot 5}{\sqrt{12}} \cup 2.71 + \frac{1.96 \cdot 5}{\sqrt{12}}, \infty \right) = (-\infty, -0.12] \cup [5.54, \infty)$$

The following figure shows the rejection and non-rejection regions for \bar{x} . The regions represent ranges of \bar{x} so they are represented by the colored bars on the \bar{x} axis. The area of the shaded region is the significance level.



The rejection (orange) and non-rejection (blue) regions for \bar{x} .

Let's redo the previous example using symbols for the known quantities as well as for μ .

Example 6. Suppose that n data points are drawn from $N(\mu, \sigma^2)$ where μ is unknown and σ is known. Set up a two-sided significance test of $H_0 : \mu = \mu_0$ using the statistic \bar{x} at significance level $\alpha = 0.05$. Describe the rejection and non-rejection regions.

answer: Under the null hypothesis $\mu = \mu_0$ we have $x_i \sim N(\mu_0, \sigma^2)$ and thus

$$\bar{x} \sim N(\mu_0, \sigma^2/n),$$

where σ^2/n is the variance $(\sigma_{\bar{x}})^2$ of \bar{x} and μ_0, σ and n are all known values.

Let $z_{\alpha/2}$ be the critical value: $P(Z > z_{\alpha/2}) = \alpha/2$. Then the non-rejection and rejection regions are separated by the values of \bar{x} that are $z_{\alpha/2} \cdot \sigma_{\bar{x}}$ from the hypothesized mean.

Since $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$ we have

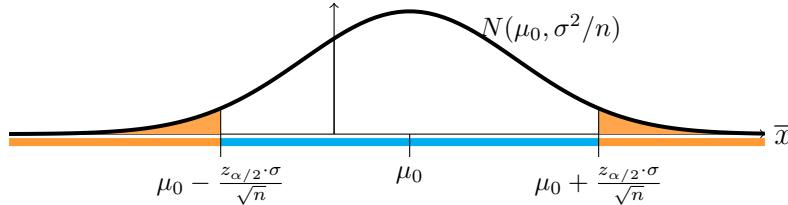
Non-rejection region:

$$\left[\mu_0 - \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}}, \quad \mu_0 + \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}} \right] \tag{2}$$

Rejection region:

$$-\infty, \mu_0 - \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}} \cup \mu_0 + \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}}, \infty .$$

We get the same figure as above, with the explicit numbers replaced by symbolic values.



The rejection (orange) and non-rejection (blue) regions for \bar{x} .

4.3 Manipulating intervals: pivoting

We need to get comfortable manipulating intervals. In general, we will make use of the type of ‘obvious’ statements that are very hard to get across. One key is to be clear about the various items.

Here is a quick summary of intervals around \bar{x} and μ_0 and what is called **pivoting**. Pivoting is the idea the \bar{x} is in $\mu_0 \pm a$ says exactly the same thing as μ_0 is in $\bar{x} \pm a$.

Example 7. Suppose we have the sample mean \bar{x} and hypothesized mean $\mu_0 = 2.71$. Suppose also that the null distribution is $N(\mu_0, 3^2)$. Then with a significance level of 0.05 we have:

- $\mu_0 + 1.96\sigma = 2.71 + 1.96(3) = 2.71 + 5.88$ is the 0.025 critical value
- $\mu_0 - 1.96\sigma = 2.71 - 1.96(3) = 2.71 - 5.88$ is the 0.975 critical value
- The non-rejection region is centered on $\mu_0 = 2.71$. That is, we don’t reject H_0 if \bar{x} is in the interval

$$[\mu_0 - 1.96\sigma, \mu_0 + 1.96\sigma] = [2.71 - 5.88, 2.71 + 5.88]$$

- The confidence interval is centered on \bar{x} . The 0.95 confidence interval uses the same width as the non-rejection region. It is the interval

$$[\bar{x} - 1.96\sigma, \bar{x} + 1.96\sigma] = [\bar{x} - 5.88, \bar{x} + 5.88]$$

There is a symmetry here: \bar{x} is in the interval $[2.71 - 1.96\sigma, 2.71 + 1.96\sigma]$ is equivalent to 2.71 is in the interval $[\bar{x} - 1.96\sigma, \bar{x} + 1.96\sigma]$.

This symmetry is called pivoting. Here are some simple numerical examples of pivoting.

- Example 8.** (i) 1.5 is in the interval $[0 - 2.3, 0 + 2.3]$, so 0 is in the interval $[1.5 - 2.3, 1.5 + 2.3]$
(ii) Likewise 1.5 is not in the interval $[0 - 1, 0 + 1]$, so 0 is not in the interval $[1.5 - 1, 1.5 + 1]$.

The symmetry might be most clear if we talk in terms of distances: the statement

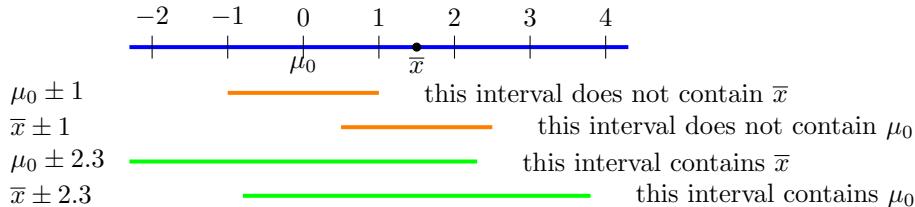
'1.5 is in the interval $[0 - 2.3, 0 + 2.3]$ '

says that the distance from 1.5 to 0 is at most 2.3. Likewise, the statement

'0 is in the interval $[1.5 - 2.3, 1.5 + 2.3]$ '

says exactly the same thing, i.e. the distance from 0 to 1.5 is less than 2.3.

Here is a visualization of *pivoting* from intervals around μ_0 to intervals around \bar{x} .



The distance between \bar{x} and μ is 1.5. Now, since $1 < 1.5$, $\mu \pm 1$, does not stretch far enough to contain \bar{x} . Likewise the interval $\bar{x} \pm 1$ does not stretch far enough to contain μ_0 . In contrast, since $2.3 > 1.5$, we have \bar{x} is in the interval $\mu_0 \pm 2.3$ and μ_0 is in the interval $\bar{x} \pm 2.3$.

4.4 Explaining the definition part 2: translating the non-rejection region to a confidence interval

The previous examples are nice if we happen to have a null hypothesis. But what if [we don't have a null hypothesis?](#) In this case, we have the point estimate \bar{x} but we still want to use the data to estimate an interval range for the unknown mean. That is, we want an interval statistic. This is given by a confidence interval.

Here we will show how to translate the notion of a non-rejection region to that of a confidence interval. The confidence level will control the rate of certain types of errors in much the same way the significance level does for NHST.

The trick is to give a little thought to the non-rejection region. Using the numbers from Example 5 we would say that at significance level 0.05 we don't reject if

$$\bar{x} \text{ is in the interval } 2.71 \pm \frac{1.96 \cdot 5}{\sqrt{12}} = 2.71 \pm 1.96 \cdot 5/\sqrt{12}. \quad (3)$$

The roles of \bar{x} and 2.71 are symmetric. The equation just above can be read as \bar{x} is within $1.96 \cdot 5/\sqrt{12}$ of 2.71. This is exactly equivalent to saying that we don't reject if

$$2.71 \text{ is in the interval } \bar{x} \pm \frac{1.96 \cdot 5}{\sqrt{12}}, \quad (4)$$

i.e. 2.71 is within $1.96 \cdot 5/\sqrt{12}$ of \bar{x} .

Now we have magically arrived at our goal of an interval statistic estimating the unknown mean. We can rewrite equation (4) as: at significance level 0.05 we don't reject if

$$\text{the interval } \left[\bar{x} - \frac{1.96 \cdot 5}{\sqrt{12}}, \bar{x} + \frac{1.96 \cdot 5}{\sqrt{12}} \right] \text{ contains } 2.71. \quad (5)$$

Thus, different values of \bar{x} generate different intervals.

The interval in equation (5) is exactly the [confidence interval](#) defined in Equation (1). We make a few observations about this confidence interval.

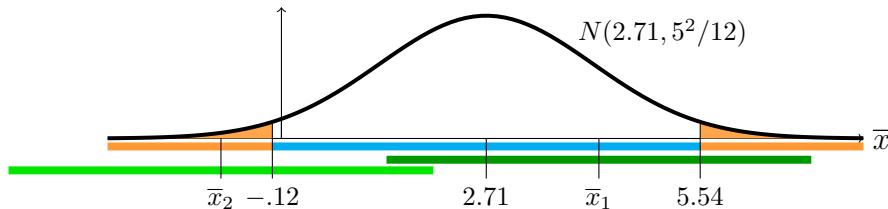
1. It only depends on \bar{x} , so it is a statistic.
2. The significance level $\alpha = 0.05$ means that, [assuming the null hypothesis that \$\mu = 2.71\$ is true](#), random data will lead us to reject the null hypothesis 5% of the time (a Type I error).
3. Again [assuming](#) that $\mu = 2.71$, then 5% of the time the confidence interval will not contain 2.71, and conversely, 95% of the time it will contain 2.71

The following figure illustrates how we don't reject H_0 if the confidence interval around it contains μ_0 and we reject H_0 if the confidence interval doesn't contain μ_0 . There is a lot in the figure so we will list carefully what you are seeing:

1. We started with the figure from Example 5 which shows the null distribution for $\mu_0 = 2.71$ and the rejection and non-rejection regions.
2. We added two possible values of the statistic \bar{x} , i.e. \bar{x}_1 and \bar{x}_2 , and their confidence intervals. Note that the width of each interval is exactly the same as the width of the non-rejection region since both use $\pm \frac{1.96 \cdot 5}{\sqrt{12}}$.

The first value, \bar{x}_1 , is in the non-rejection region and its interval includes the null hypothesis $\mu_0 = 2.71$. This illustrates that [not rejecting \$H_0\$](#) corresponds to the confidence interval [containing \$\mu_0\$](#) .

The second value, \bar{x}_2 , is in the rejection region and its interval does not contain μ_0 . This illustrates that [rejecting \$H_0\$](#) corresponds to the confidence interval [not containing \$\mu_0\$](#) .



The non-rejection region (blue) and two confidence intervals (green).

We can still wring one more essential observation out of this example. Our choice of null hypothesis $\mu = 2.71$ was completely arbitrary. If we replace $\mu = 2.71$ by any other hypothesis $\mu = \mu_0$ then the interval (5) will come out the same.

We call the interval (5) a 95% [confidence interval](#) because, [assuming \$\mu = \mu_0\$](#) , on average it will contain μ_0 in 95% of random trials.

4.5 Explaining the definition part 3: translating a general non-rejection region to a confidence interval

Note that the specific values of σ and n in the preceding example were of no particular consequence, so they can be replaced by their symbols. In this way we can take Example (6) quickly through the same steps as Example (5).

In words, Equation (2) and the corresponding figure say that we don't reject if

$$\bar{x} \text{ is in the interval } \mu_0 \pm \frac{z_{\alpha/2}\sigma}{\sqrt{n}}.$$

This is exactly equivalent to saying that we don't reject if

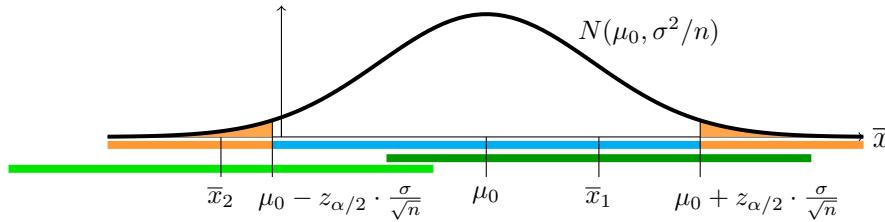
$$\mu_0 \text{ is in the interval } \bar{x} \pm \frac{z_{\alpha/2}\sigma}{\sqrt{n}}. \quad (6)$$

We can rewrite equation (6) as: at significance level α we don't reject if

$$\text{the interval } \left[\bar{x} - \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}}, \bar{x} + \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}} \right] \text{ contains } \mu_0. \quad (7)$$

We call the interval (7) a $(1 - \alpha)$ **confidence interval** because, assuming $\mu = \mu_0$, on average it will contain μ_0 in the fraction $(1 - \alpha)$ of random trials.

The following figure illustrates the point that μ_0 is in the $(1 - \alpha)$ confidence interval around \bar{x} is equivalent to \bar{x} is in the non-rejection region (at significance level α) for $H_0 : \mu_0 = \mu$.



\bar{x}_1 is in non-rejection region for $\mu_0 \Leftrightarrow$ the confidence interval around \bar{x}_1 contains μ_0 .

4.6 Computational example

Example 9. Suppose the data 2.5, 5.5, 8.5, 11.5 was drawn from a $N(\mu, 10^2)$ distribution with unknown mean μ .

(a) Compute the point estimate \bar{x} for μ and the corresponding 50%, 80% and 95% confidence intervals.

(b) Consider the null hypothesis $\mu = 1$. Would you reject H_0 at $\alpha = 0.05$? $\alpha = 0.20$? $\alpha = 0.50$? Do these two ways: first by checking if the hypothesized value of μ is in the relevant confidence interval and second by constructing a rejection region.

answer: (a) We compute that $\bar{x} = 7.0$. The critical points are

$$z_{0.025} = \text{qnorm}(0.975) = 1.96, \quad z_{0.1} = \text{qnorm}(0.9) = 1.28, \quad z_{0.25} = \text{qnorm}(0.75) = 0.67.$$

Since $n = 4$ we have $\bar{x} \sim N(\mu, 10^2/4)$, i.e. $\sigma_{\bar{x}} = 5$. So we have:

$$\begin{aligned} \text{95\% conf. interval} &= [\bar{x} - z_{0.025}\sigma_{\bar{x}}, \bar{x} + z_{0.025}\sigma_{\bar{x}}] = [7 - 1.96 \cdot 5, 7 + 1.96 \cdot 5] = [-2.8, 16.8] \\ \text{80\% conf. interval} &= [\bar{x} - z_{0.1}\sigma_{\bar{x}}, \bar{x} + z_{0.1}\sigma_{\bar{x}}] = [7 - 1.28 \cdot 5, 7 + 1.28 \cdot 5] = [0.6, 13.4] \\ \text{50\% conf. interval} &= [\bar{x} - z_{0.75}\sigma_{\bar{x}}, \bar{x} + z_{0.75}\sigma_{\bar{x}}] = [7 - 0.67 \cdot 5, 7 + 0.67 \cdot 5] = [3.65, 10.35] \end{aligned}$$

Each of these intervals is a range estimate of μ . Notice that the higher the confidence level, the wider the interval needs to be.

(b) Since $\mu = 1$ is in the 95% and 80% confidence intervals, we would not reject the null hypothesis at the $\alpha = 0.05$ or $\alpha = 0.20$ levels. Since $\mu = 1$ is not in the 50% confidence interval, we would reject H_0 at the $\alpha = 0.5$ level.

We construct the rejection regions using the same critical values as in part (a). The difference is that rejection regions are intervals centered on the hypothesized value for μ : $\mu_0 = 1$ and confidence intervals are centered on \bar{x} . Here are the rejection regions.

$$\begin{aligned}\alpha = 0.05 &\Rightarrow (-\infty, \mu_0 - z_{0.025} \sigma_{\bar{x}}] \cup [\mu_0 + z_{0.025} \sigma_{\bar{x}}, \infty) = (-\infty, -8.8] \cup [10.8, \infty) \\ \alpha = 0.20 &\Rightarrow (-\infty, \mu_0 - z_{0.1} \sigma_{\bar{x}}] \cup [\mu_0 + z_{0.1} \sigma_{\bar{x}}, \infty) = (-\infty, -5.4] \cup [7.4, \infty) \\ \alpha = 0.25 &\Rightarrow (-\infty, \mu_0 - z_{0.25} \sigma_{\bar{x}}] \cup [\mu_0 + z_{0.25} \sigma_{\bar{x}}, \infty) = (-\infty, -2.35] \cup [4.35, \infty)\end{aligned}$$

To do the NHST we must check whether or not $\bar{x} = 7$ is in the rejection region.

$\alpha = 0.05$: $7 < 10.8$ is not in the rejection region.

We do not reject the hypothesis that $\mu = 1$ at a significance level of 0.05.

$\alpha = 0.2$: $7 < 7.4$ is not in the rejection region.

We do not reject the hypothesis that $\mu = 1$ at a significance level of 0.2.

$\alpha = 0.5$: $7 > 4.35$ is in the rejection region.

We reject the hypothesis that $\mu = 1$ at a significance level 0.5.

We get the same answers using either method.

5 t -confidence intervals for the mean

This will be nearly identical to normal confidence intervals. In this setting σ is not known, so we have to make the following replacements.

1. Use $s_{\bar{x}} = \frac{s}{\sqrt{n}}$ instead of $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$. Here s is the sample variance we used before in t -tests
2. Use t -critical values instead of z -critical values.

5.1 Definition of t -confidence intervals for the mean

Definition: Suppose that $x_1, \dots, x_n \sim N(\mu, \sigma^2)$, where the values of the mean μ and the standard deviation σ are both unknown. The $(1 - \alpha)$ confidence interval for μ is

$$\left[\bar{x} - \frac{t_{\alpha/2} \cdot s}{\sqrt{n}}, \bar{x} + \frac{t_{\alpha/2} \cdot s}{\sqrt{n}} \right], \quad (8)$$

here $t_{\alpha/2}$ is the right critical value $P(T > t_{\alpha/2}) = \alpha/2$ for $T \sim t(n - 1)$ and s^2 is the sample variance of the data.

5.2 Construction of t confidence intervals

Suppose that n data points are drawn from $N(\mu, \sigma^2)$ where μ and σ are unknown. We'll derive the t confidence interval following the same pattern as for the z confidence interval.

Under the null hypothesis $\mu = \mu_0$, we have $x_i \sim N(\mu_0, \sigma^2)$. So the studentized mean follows a Student t distribution with $n - 1$ degrees of freedom:

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \sim t(n - 1).$$

Let $t_{\alpha/2}$ be the critical value: $P(T > t_{\alpha/2}) = \alpha/2$, where $T \sim t(n - 1)$. We know from running one-sample t -tests that the non-rejection region is given by

$$|t| \leq t_{\alpha/2}$$

Using the definition of the t -statistic to write the rejection region in terms of \bar{x} we get: at significance level α we don't reject if

$$\frac{|\bar{x} - \mu_0|}{s/\sqrt{n}} \leq t_{\alpha/2} \Leftrightarrow |\bar{x} - \mu_0| \leq t_{\alpha/2} \cdot \frac{s}{\sqrt{n}}.$$

Geometrically, the right hand side says that we don't reject if

$$\mu_0 \text{ is within } t_{\alpha/2} \cdot \frac{s}{\sqrt{n}} \text{ of } \bar{x}.$$

This is exactly equivalent to saying that we don't reject if

$$\text{the interval } \left[\bar{x} - \frac{t_{\alpha/2} \cdot s}{\sqrt{n}}, \bar{x} + \frac{t_{\alpha/2} \cdot s}{\sqrt{n}} \right] \text{ contains } \mu_0.$$

This interval is the confidence interval defined in (8).

Example 10. Suppose the data 2.5, 5.5, 8.5, 11.5 was drawn from a $N(\mu, \sigma^2)$ distribution with μ and σ both unknown.

Give interval estimates for μ by finding the 95%, 80% and 50% confidence intervals.

answer: By direct computation we have $\bar{x} = 7$ and $s^2 = 15$. The critical points are $t_{0.025} = qt(0.975) = 3.18$, $t_{0.1} = qt(0.9) = 1.64$, and $t_{0.25} = qt(0.75) = 0.76$.

$$\begin{aligned} 95\% \text{ conf. interval} &= \bar{x} - t_{0.025} \cdot \frac{s}{\sqrt{n}}, \quad \bar{x} + t_{0.025} \cdot \frac{s}{\sqrt{n}} = [0.84, 13.16] \\ 80\% \text{ conf. interval} &= \bar{x} - t_{0.1} \cdot \frac{s}{\sqrt{n}}, \quad \bar{x} + t_{0.1} \cdot \frac{s}{\sqrt{n}} = [3.82, 10.18] \\ 50\% \text{ conf. interval} &= \bar{x} - t_{0.25} \cdot \frac{s}{\sqrt{n}}, \quad \bar{x} + t_{0.25} \cdot \frac{s}{\sqrt{n}} = [5.53, 8.47] \end{aligned}$$

All of these confidence intervals give interval estimates for the value of μ . Again, notice that the higher the confidence level, the wider the corresponding interval.

6 Chi-square confidence intervals for the variance

We now turn to an interval estimate for the unknown variance.

Definition: Suppose the data x_1, \dots, x_n is drawn from $N(\mu, \sigma^2)$ with mean μ and standard deviation σ both unknown. The $(1 - \alpha)$ confidence interval for the variance σ^2 is

$$\frac{(n - 1)s^2}{c_{\alpha/2}}, \quad \frac{(n - 1)s^2}{c_{1-\alpha/2}}. \quad (9)$$

Here $c_{\alpha/2}$ is the **right critical value** $P(X^2 > c_{\alpha/2}) = \alpha/2$ for $X^2 \sim \chi^2(n - 1)$ and s^2 is the sample variance of the data.

The derivation of this interval is nearly identical to that of the previous derivations, now starting from the chi-square test for variance. The basic fact we need is that, for data drawn from $N(\mu, \sigma^2)$ with known σ , the statistic

$$\frac{(n-1)s^2}{\sigma^2}$$

follows a chi-square distribution with $n-1$ degrees of freedom. So given the null hypothesis $H_0 : \sigma = \sigma_0$, the test statistic is $(n-1)s^2/\sigma_0^2$ and the non-rejection region at significance level α is

$$c_{1-\alpha/2} < \frac{(n-1)s^2}{\sigma_0^2} < c_{\alpha/2}.$$

A little algebra converts this to

$$\frac{(n-1)s^2}{c_{1-\alpha/2}} > \sigma_0^2 > \frac{(n-1)s^2}{c_{\alpha/2}}.$$

This says we don't reject if

the interval $\left[\frac{(n-1)s^2}{c_{\alpha/2}}, \frac{(n-1)s^2}{c_{1-\alpha/2}} \right]$ contains σ_0^2

This is our $(1 - \alpha)$ confidence interval.

We will continue our exploration of confidence intervals next class. In the meantime, truly the best way is to internalize the meaning of the confidence level is to experiment with the confidence interval applet:

<http://mathlets.org/mathlets/confidence-intervals/>

MIT OpenCourseWare
<https://ocw.mit.edu>

18.05 Introduction to Probability and Statistics
Spring 2014

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.

Confidence Intervals: Three Views

Class 23, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to produce z , t and χ^2 confidence intervals based on the corresponding standardized statistics.
2. Be able to use a hypothesis test to construct a confidence interval for an unknown parameter.
3. Refuse to answer questions that ask, in essence, ‘given a confidence interval what is the probability or odds that it contains the true value of the unknown parameter?’

2 Introduction

Our approach to confidence intervals in the previous reading was a combination of standardized statistics and hypothesis testing. Today we will consider each of these perspectives separately, as well as introduce a third formal viewpoint. Each provides its own insight.

1. **Standardized statistic.** Most confidence intervals are based on standardized statistics with known distributions like z , t or χ^2 . This provides a straightforward way to construct and interpret confidence intervals as a point estimate plus or minus some error.
2. **Hypothesis testing.** Confidence intervals may also be constructed from hypothesis tests. In cases where we don’t have a standardized statistic this method will still work. It agrees with the standardized statistic approach in cases where they both apply.

This view connects the notions of significance level α for hypothesis testing and confidence level $1 - \alpha$ for confidence intervals; we will see that in both cases α is the probability of making a ‘type 1’ error. This gives some insight into the use of the word confidence. This view also helps to emphasize the frequentist nature of confidence intervals.

3. **Formal.** The formal definition of confidence intervals is perfectly precise and general. In a mathematical sense it gives insight into the inner workings of confidence intervals. However, because it is so general it sometimes leads to confidence intervals without useful properties. We will not dwell on this approach. We offer it mainly for those who are interested.

3 Confidence intervals via standardized statistics

The strategy here is essentially the same as in the previous reading. Assuming normal data we have what we called standardized statistics like the standardized mean, Studentized mean, and standardized variance. These statistics have well known distributions which depend on hypothesized values of μ and σ . We then use algebra to produce confidence intervals for μ or σ .

Don't let the algebraic details distract you from the essentially simple idea underlying confidence intervals: we start with a standardized statistic (e.g., z , t or χ^2) and use some algebra to get an interval that depends only on the data and [known](#) parameters.

3.1 z -confidence intervals for μ : normal data with known σ

z -confidence intervals for the mean of normal data are based on the [standardized mean](#), i.e. the z -statistic. We start with n independent normal samples

$$x_1, x_2, \dots, x_n \sim N(\mu, \sigma^2).$$

We assume that μ is the unknown parameter of interest and σ is known.

We know that the standardized mean is standard normal:

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1).$$

For the standard normal critical value $z_{\alpha/2}$ we have: $P(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1 - \alpha$.

Thus,

$$P\left(-z_{\alpha/2} < \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} < z_{\alpha/2} \mid \mu\right) = 1 - \alpha$$

A little bit of algebra puts this in the form of an interval around μ :

$$P\left(\bar{x} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \mid \mu\right) = 1 - \alpha$$

We can emphasize that the interval depends only on the statistic \bar{x} and the known value σ by writing this as

$$P\left(\left[\bar{x} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right] \text{ contains } \mu \mid \mu\right) = 1 - \alpha.$$

This is the $(1 - \alpha)$ z -confidence interval for μ . We often write it using the shorthand

$$\bar{x} \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

Think of it as $\bar{x} \pm$ error.

Make sure you notice that the [probabilities are conditioned on \$\mu\$](#) . As with all frequentist statistics, we have to fix hypothesized values of the parameters in order to compute probabilities.

3.2 t -confidence intervals for μ : normal data with unknown μ and σ

t -confidence intervals for the mean of normal data are based on the Studentized mean, i.e. the t -statistic.

Again we have $x_1, x_2, \dots, x_n \sim N(\mu, \sigma^2)$, but now we assume both μ and σ are unknown. We know that the Studentized mean follows a Student t distribution with $n - 1$ degrees of freedom. That is,

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} \sim t(n - 1),$$

where s^2 is the sample variance.

Now all we have to do is replace the standardized mean by the [Studentized mean](#) and the same logic we used for z gives us the t -confidence interval: start with

$$P\left(-t_{\alpha/2} < \frac{\bar{x} - \mu}{s/\sqrt{n}} < t_{\alpha/2} \mid \mu\right) = 1 - \alpha.$$

A little bit of algebra isolates μ in the middle of an interval:

$$P\left(\bar{x} - t_{\alpha/2} \cdot \frac{s}{\sqrt{n}} < \mu < \bar{x} + t_{\alpha/2} \cdot \frac{s}{\sqrt{n}} \mid \mu\right) = 1 - \alpha$$

We can emphasize that the interval depends only on the statistics \bar{x} and s by writing this as

$$P\left(\left[\bar{x} - t_{\alpha/2} \cdot \frac{s}{\sqrt{n}}, \bar{x} + t_{\alpha/2} \cdot \frac{s}{\sqrt{n}}\right] \text{ contains } \mu \mid \mu\right) = 1 - \alpha.$$

This is the $(1 - \alpha)$ t -confidence interval for μ . We often write it using the shorthand

$$\bar{x} \pm t_{\alpha/2} \cdot \frac{s}{\sqrt{n}}$$

Think of it as $\bar{x} \pm$ error.

3.3 χ^2 -confidence intervals for σ^2 : normal data with unknown μ and σ

You guessed it: χ^2 -confidence intervals for the variance of normal data are based on the [standardized variance](#), i.e. the χ^2 -statistic.

We follow the same logic as above to get a χ^2 -confidence interval for σ^2 . Because this is the third time through it we'll move a little more quickly.

We assume we have n independent normal samples: $x_1, x_2, \dots, x_n \sim N(\mu, \sigma^2)$. We assume that μ and σ are both unknown. The standardized variance is

$$X^2 = \frac{(n-1)s^2}{\sigma^2} \sim \chi^2(n-1).$$

We know that the X^2 statistic follows a χ^2 distribution with $n-1$ degrees of freedom.

For Z and t we used, without comment, the symmetry of the distributions to replace $z_{1-\alpha/2}$ by $-z_{\alpha/2}$ and $t_{1-\alpha/2}$ by $-t_{\alpha/2}$. Because the χ^2 distribution is not symmetric we need to be explicit about the critical values on both the left and the right. That is,

$$P(c_{1-\alpha/2} < X^2 < c_{\alpha/2}) = 1 - \alpha,$$

where $c_{\alpha/2}$ and $c_{1-\alpha/2}$ are [right tail](#) critical values. Thus,

$$P(c_{1-\alpha/2} < \frac{(n-1)s^2}{\sigma^2} < c_{\alpha/2} \mid \sigma) = 1 - \alpha$$

A little bit of algebra puts this in the form of an interval around σ^2 :

$$P\left(\frac{(n-1)s^2}{c_{\alpha/2}} < \sigma^2 < \frac{(n-1)s^2}{c_{1-\alpha/2}} \mid \sigma\right) = 1 - \alpha$$

We can emphasize that the interval depends only on the statistic s^2 by writing this as

$$P\left(\left[\frac{(n-1)s^2}{c_{\alpha/2}}, \frac{(n-1)s^2}{c_{1-\alpha/2}}\right] \text{ contains } \sigma^2 \mid \sigma^2\right) = 1 - \alpha.$$

This is the $(1 - \alpha)$ χ^2 -confidence interval for σ^2 .

4 Confidence intervals via hypothesis testing

Suppose we have data drawn from a distribution with a parameter θ whose value is unknown. A significance test for the value θ has the following short description.

1. Set the null hypothesis $H_0 : \theta = \theta_0$ for some special value θ_0 , e.g. we often have $H_0 : \theta = 0$.
2. Use the data to compute the value of a test statistic, call it x .
3. If x is far enough into the tail of the null distribution (the distribution [assuming](#) the null hypothesis) then we reject H_0 .

In the case where there is no special value to test we may still want to estimate θ . This is the reverse of significance testing; rather than seeing if we should reject a specific value of θ because it doesn't fit the data we want to find the range of values of θ that do, in some sense, fit the data. This gives us the following definitions.

Definition. Given a value x of the test statistic, the $(1 - \alpha)$ confidence interval contains all values θ_0 which are not rejected (at significance level α) when they are the null hypothesis.

Definition. A type 1 CI error occurs when the confidence interval does not contain the true value of θ .

For a $(1 - \alpha)$ confidence interval the type 1 CI error rate is α .

Example 1. Here is an example relating confidence intervals and hypothesis tests. Suppose data x is drawn from a binomial(12, θ) distribution with θ unknown. Let $\alpha = 0.1$ and create the $(1 - \alpha) = 90\%$ confidence interval for each possible value of x .

Our strategy is to look at one possible value of θ at a time and choose rejection regions for a significance test with $\alpha = 0.1$. Once this is done, we will know, for each value of x , which values of θ are not rejected, i.e. the confidence interval associated with x .

To start we set up a likelihood table for binomial(12, θ) in Table 1. Each row shows the probabilities $p(x|\theta)$ for one value of θ . To keep the size manageable we only show θ in increments of 0.1.

$\theta \setminus x$	0	1	2	3	4	5	6	7	8	9	10	11	12
1.0	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00
0.9	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.02	0.09	0.23	0.38	0.28
0.8	0.00	0.00	0.00	0.00	0.00	0.02	0.05	0.13	0.24	0.28	0.21	0.07	
0.7	0.00	0.00	0.00	0.00	0.01	0.03	0.08	0.16	0.23	0.24	0.17	0.07	0.01
0.6	0.00	0.00	0.00	0.01	0.04	0.10	0.18	0.23	0.21	0.14	0.06	0.02	0.00
0.5	0.00	0.00	0.02	0.05	0.12	0.19	0.23	0.19	0.12	0.05	0.02	0.00	0.00
0.4	0.00	0.02	0.06	0.14	0.21	0.23	0.18	0.10	0.04	0.01	0.00	0.00	0.00
0.3	0.01	0.07	0.17	0.24	0.23	0.16	0.08	0.03	0.01	0.00	0.00	0.00	0.00
0.2	0.07	0.21	0.28	0.24	0.13	0.05	0.02	0.00	0.00	0.00	0.00	0.00	0.00
0.1	0.28	0.38	0.23	0.09	0.02	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.0	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

Table 1. Likelihood table for Binomial(12, θ)

Tables 2-4 below show the rejection region (in orange) and non-rejection region (in blue) for the various values of θ . To emphasize the row-by-row nature of the process the Table 2 just shows these regions for $\theta = 1.0$, then Table 3 adds in regions for $\theta = 0.9$ and Table 4 shows them for all the values of θ .

Immediately following the tables we give a detailed explanation of how the rejection/non-rejection regions were chosen.

$\theta \setminus x$	0	1	2	3	4	5	6	7	8	9	10	11	12	significance
1.0	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	0.000
0.9	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.02	0.09	0.23	0.38	0.28	
0.8	0.00	0.00	0.00	0.00	0.00	0.00	0.02	0.05	0.13	0.24	0.28	0.21	0.07	
0.7	0.00	0.00	0.00	0.00	0.01	0.03	0.08	0.16	0.23	0.24	0.17	0.07	0.01	
0.6	0.00	0.00	0.00	0.01	0.04	0.10	0.18	0.23	0.21	0.14	0.06	0.02	0.00	
0.5	0.00	0.00	0.02	0.05	0.12	0.19	0.23	0.19	0.12	0.05	0.02	0.00	0.00	
0.4	0.00	0.02	0.06	0.14	0.21	0.23	0.18	0.10	0.04	0.01	0.00	0.00	0.00	
0.3	0.01	0.07	0.17	0.24	0.23	0.16	0.08	0.03	0.01	0.00	0.00	0.00	0.00	
0.2	0.07	0.21	0.28	0.24	0.13	0.05	0.02	0.00	0.00	0.00	0.00	0.00	0.00	
0.1	0.28	0.38	0.23	0.09	0.02	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
0.0	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	

Table 2. Likelihood table for binomial(12, θ) with rejection/non-rejection regions for $\theta = 1.0$

$\theta \setminus x$	0	1	2	3	4	5	6	7	8	9	10	11	12	significance
1.0	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	0.000
0.9	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.02	0.09	0.23	0.38	0.28	0.026
0.8	0.00	0.00	0.00	0.00	0.00	0.00	0.02	0.05	0.13	0.24	0.28	0.21	0.07	
0.7	0.00	0.00	0.00	0.00	0.01	0.03	0.08	0.16	0.23	0.24	0.17	0.07	0.01	
0.6	0.00	0.00	0.00	0.01	0.04	0.10	0.18	0.23	0.21	0.14	0.06	0.02	0.00	
0.5	0.00	0.00	0.02	0.05	0.12	0.19	0.23	0.19	0.12	0.05	0.02	0.00	0.00	
0.4	0.00	0.02	0.06	0.14	0.21	0.23	0.18	0.10	0.04	0.01	0.00	0.00	0.00	
0.3	0.01	0.07	0.17	0.24	0.23	0.16	0.08	0.03	0.01	0.00	0.00	0.00	0.00	
0.2	0.07	0.21	0.28	0.24	0.13	0.05	0.02	0.00	0.00	0.00	0.00	0.00	0.00	
0.1	0.28	0.38	0.23	0.09	0.02	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
0.0	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	

Table 3. Likelihood table with rejection/non-rejection regions shown for $\theta = 1.0$ and 0.9

$\theta \setminus x$	0	1	2	3	4	5	6	7	8	9	10	11	12	significance
1.0	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	0.000
0.9	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.02	0.09	0.23	0.38	0.28	0.026
0.8	0.00	0.00	0.00	0.00	0.00	0.00	0.02	0.05	0.13	0.24	0.28	0.21	0.07	0.073
0.7	0.00	0.00	0.00	0.00	0.01	0.03	0.08	0.16	0.23	0.24	0.17	0.07	0.01	0.052
0.6	0.00	0.00	0.00	0.01	0.04	0.10	0.18	0.23	0.21	0.14	0.06	0.02	0.00	0.077
0.5	0.00	0.00	0.02	0.05	0.12	0.19	0.23	0.19	0.12	0.05	0.02	0.00	0.00	0.092
0.4	0.00	0.02	0.06	0.14	0.21	0.23	0.18	0.10	0.04	0.01	0.00	0.00	0.00	0.077
0.3	0.01	0.07	0.17	0.24	0.23	0.16	0.08	0.03	0.01	0.00	0.00	0.00	0.00	0.052
0.2	0.07	0.21	0.28	0.24	0.13	0.05	0.02	0.00	0.00	0.00	0.00	0.00	0.00	0.073
0.1	0.28	0.38	0.23	0.09	0.02	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.026
0.0	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.000

Table 4. Likelihood table with rejection/non-rejection regions for $\theta = 0.0$ to 1.0 **Choosing the rejection and non-rejection regions in the tables**

The first problem we confront is how exactly to choose the rejection region. We used two rules:

1. The total probability of the rejection region, i.e. the significance, should be less than or equal to 0.1. (Since we have a discrete distribution it is impossible to make the significance exactly 0.1.)
2. We build the rejection region by choosing values of x one at a time, always picking the unused value with the smallest probability. We stop when the next value would make the significance more than 0.1.

There are other ways to choose the rejection region which would result in slight differences. Our method is one reasonable way.

Table 2 shows the rejection (orange) and non-rejection (blue) regions for $\theta = 1.0$. This is a special case because most of the probabilities in this row are 0.0. We'll move right on to the next table and step through the process for that.

In Table 3, let's walk through the steps used to find these regions for $\theta = 0.9$.

- The smallest probability is when $x = 0$, so $x = 0$ is in the rejection region.
- The next smallest is when $x = 1$, so $x = 1$ is in the rejection region.
- We continue with $x = 2, \dots, 8$. At this point the total probability in the rejection region is 0.026.
- The next smallest probability is when $x = 9$. Adding this probability (0.09) to 0.026 would put the total probability over 0.1. So we leave $x = 9$ out of the rejection region and stop the process.

Note three things for the $\theta = 0.9$ row:

1. None of the probabilities in this row are truly zero, though some are small enough that they equal 0 to 2 decimal places.
2. We show the significance for this value of θ in the right hand margin. More precisely, we show the significance level of the NHST with null hypothesis $\theta = 0.9$ and the given rejection region.
3. The rejection region consists of values of x . When we say the rejection region is shown in orange we really mean the rejection region contains the values of x corresponding to the probabilities highlighted in orange.

Think: Look back at the $\theta = 1.0$ row and make sure you understand why the rejection region is $x = 0, \dots, 11$ and the significance is 0.000.

Example 2. Using Table 4 determine the 0.90 confidence interval when $x = 8$.

answer: The 90% confidence interval consists of all those θ that would not be rejected by an $\alpha = 0.1$ hypothesis test when $x = 8$. Looking at the table, the blue (non-rejected) entries in the column $x = 8$ correspond to $0.5 \leq \theta \leq 0.8$: the confidence interval is $[0.5, 0.8]$.

Remark: The point of this example is to show how confidence intervals and hypothesis tests are related. Since Table 4 has only finitely many values of θ , our answer is close but not exact. Using a computer we could look at many more values of θ . For this problem we used R to find that, correct to 2 decimal places, the confidence interval is $[0.42, 0.85]$.

Example 3. Explain why the expected type one CI error rate will be at most 0.092, provided that the true value of θ is in the table.

answer: The short answer is that this is the maximum significance for any θ in Table 4. Expanding on that slightly: we make a type one CI error if the confidence interval does not contain the true value of θ , call it θ_{true} . This happens exactly when the data x is in the rejection region for θ_{true} . The probability of this happening is the significance for θ_{true} and this is at most 0.092.

Remark: The point of this example is to show how confidence level, type one CI error rate and significance for each hypothesis are related. As in the previous example, we can use R to compute the significance for many more values of θ . When we do this we find that the maximum significance for any θ is 0.1 occurring when $\theta \approx 0.0452$.

Summary notes:

1. We start with a test statistic x . The confidence interval is random because it depends on x .
2. For each hypothesized value of θ we make a significance test with significance level α by choosing rejection regions.
3. For a specific value of x the associated confidence interval for θ consists of all θ that aren't rejected for that value, i.e. all θ that have x in their non-rejection regions.
4. Because the distribution is discrete we can't always achieve the exact significance level, so our confidence interval is really an 'at least 90% confidence interval'.

Example 4. Open the applet <http://mathlets.org/mathlets/confidence-intervals/>. We want you to play with the applet to understand the random nature of confidence intervals and the meaning of confidence as (1 - type I CI error rate).

- (a) Read the help. It is short and will help orient you in the applet. Play with different settings of the parameters to see how they affect the size of the confidence intervals.
- (b) Set the number of trials to $N = 1$. Click the 'Run N trials' button repeatedly and see that each time data is generated the confidence intervals jump around.
- (c) Now set the confidence level to $c = .5$. As you click the 'Run N trials' button you should see that about 50% of the confidence intervals include the true value of μ . The 'Z correct' and 't correct' values should change accordingly.
- (d) Now set the number of trials to $N = 100$. With $c = .8$. The 'Run N trials' button will now run 100 trials at a time. Only the last confidence interval will be shown in the graph, but the trials all run and the 'percent correct' statistics will be updated based on all 100 trials.

Click the run trials button repeatedly. Watch the correct rates start to converge to the confidence level. To converge even faster, set $N = 1000$.

5 Formal view of confidence intervals

Recall: An interval statistic is an interval I_x computed from data x . An interval is determined by its lower and upper bounds, and these are random because x is random.

We suppose that x is drawn from a distribution with pdf $f(x|\theta)$ where the parameter θ is unknown.

Definition: A $(1 - \alpha)$ confidence interval for θ is an interval statistic I_x such that

$$P(I_x \text{ contains } \theta_0 \mid \theta = \theta_0) = 1 - \alpha$$

for all possible values of θ_0 .

We wish this was simpler, but a definition is a definition and this definition is one way to weigh the evidence provided by the data x . Let's unpack it a bit.

The confidence level of an interval statistic is a probability concerning a random interval and a hypothesized value θ_0 for the unknown parameter. Precisely, it is the probability that the random interval (computed from random data) contains the value θ_0 , given that the model parameter truly is θ_0 . Since the true value of θ is unknown, the frequentist statistician defines a, say, 95% confidence intervals so that the 0.95 probability is valid no matter which hypothesized value of the parameter is actually true.

6 Comparison with Bayesian probability intervals

Confidence intervals are a frequentist notion, and as we've repeated many times, frequentists don't assign probabilities to hypotheses, e.g. the value of an unknown parameter. Rather they compute likelihoods; that is, probabilities about data or associated statistics given a hypothesis (note the condition $\theta = \theta_0$ in the formal view of confidence intervals). Note that the construction of confidence intervals proceeds entirely from the full likelihood table.

In contrast Bayesian posterior probability intervals are truly the probability that the value of the unknown parameter lies in the reported range. We add the usual caveat that this depends on the specific choice of a (possibly subjective) Bayesian prior.

This distinction between the two is subtle because Bayesian posterior probability intervals and frequentist confidence intervals share the following properties:

1. They start from a model $f(x|\theta)$ for observed data x with unknown parameter θ .
2. Given data x , they give an interval $I(x)$ specifying a range of values for θ .
3. They come with a number (say .95) that is the probability of something.

In practice, many people misinterpret confidence intervals as Bayesian probability intervals, forgetting that frequentists never place probabilities on hypotheses (this is analogous to mistaking the p -value in NHST for the probability that H_0 is false). The harm of this misinterpretation is somewhat mitigated by that fact that, given enough data and a reasonable prior, Bayesian and frequentist intervals often work out to be quite similar.

For an amusing example illustrating how they can be quite different, see the first answer here (involving chocolate chip cookies!):

<http://stats.stackexchange.com/questions/2272/whats-the-difference-between-a-confidence-interval-and-a-credible-interval>

This example uses the formal definitions and is really about confidence sets instead of confidence intervals.

MIT OpenCourseWare
<https://ocw.mit.edu>

18.05 Introduction to Probability and Statistics
Spring 2014

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.

Confidence Intervals for the Mean of Non-normal Data

Class 23, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to derive the formula for conservative normal confidence intervals for the proportion θ in Bernoulli data.
2. Be able to compute rule-of-thumb 95% confidence intervals for the proportion θ of a Bernoulli distribution.
3. Be able to compute large sample confidence intervals for the mean of a general distribution.

2 Introduction

So far, we have focused on constructing confidence intervals for data drawn from a normal distribution. We'll now will switch gears and learn about confidence intervals for the mean when the data is not necessarily normal.

We will first look carefully at estimating the probability θ of success when the data is drawn from a $\text{Bernoulli}(\theta)$ distribution –recall that θ is also the mean of the Bernoulli distribution.

Then we will consider the case of a a large sample from an unknown distribution; in this case we can appeal to the central limit theorem to justify the use z -confidence intervals.

3 Bernoulli data and polling

One common use of confidence intervals is for estimating the proportion θ in a $\text{Bernoulli}(\theta)$ distribution. For example, suppose we want to use a political poll to estimate the proportion of the population that supports candidate A, or equivalent the probability θ that a random person supports candidate A. In this case we have a simple rule-of-thumb that allows us to quickly compute a confidence interval.

3.1 Conservative normal confidence intervals

Suppose we have i.i.d. data x_1, x_2, \dots, x_n all drawn from a $\text{Bernoulli}(\theta)$ distribution. then a **conservative normal** $(1 - \alpha)$ confidence interval for θ is given by

$$\bar{x} \pm z_{\alpha/2} \cdot \frac{1}{2\sqrt{n}}. \quad (1)$$

The proof given below uses the central limit theorem and the observation that $\sigma = \sqrt{\theta(1 - \theta)} \leq 1/2$.

You will also see in the derivation below that this formula is conservative, providing an ‘at least $(1 - \alpha)$ ’ confidence interval.

Example 1. A pollster asks 196 people if they prefer candidate A to candidate B and finds that 120 prefer A and 76 prefer B. Find the 95% conservative normal confidence interval for θ , the proportion of the population that prefers A.

answer: We have $\bar{x} = 120/196 = 0.612$, $\alpha = 0.05$ and $z_{.025} = 1.96$. The formula says a 95% confidence interval is

$$I \approx 0.612 \pm \frac{1.96}{2 \cdot 14} = 0.612 \pm 0.007.$$

3.2 Proof of Formula 1

The proof of Formula 1 will rely on the following fact.

Fact. The standard deviation of a $Bernoulli(\theta)$ distribution is at most 0.5.

Proof of fact: Let’s denote this standard deviation by σ_θ to emphasize its dependence on θ . The variance is then $\sigma_\theta^2 = \theta(1 - \theta)$. It’s easy to see using calculus or by graphing this parabola that the maximum occurs when $\theta = 1/2$. Therefore the maximum variance is $1/4$, which implies that the standard deviation σ_θ is less than $\sqrt{1/4} = 1/2$.

Proof of formula (1). The proof relies on the central limit theorem which says that (for large n) the distribution of \bar{x} is approximately normal with mean θ and standard deviation σ_θ/\sqrt{n} . For normal data we have the $(1 - \alpha)$ z -confidence interval

$$\bar{x} \pm z_{\alpha/2} \cdot \frac{\sigma_\theta}{\sqrt{n}}$$

The trick now is to replace σ_θ by $\frac{1}{2}$: since $\sigma_\theta \leq \frac{1}{2}$ the resulting interval around \bar{x}

$$\bar{x} \pm z_{\alpha/2} \cdot \frac{1}{2\sqrt{n}}$$

is always at least as wide as the interval using $\pm \sigma_\theta/\sqrt{n}$. A wider interval is more likely to contain the true value of θ so we have a ‘conservative’ $(1 - \alpha)$ confidence interval for θ .

Again, we call this **conservative** because $\frac{1}{2\sqrt{n}}$ overestimates the standard deviation of \bar{x} , resulting in a wider interval than is necessary to achieve a $(1 - \alpha)$ confidence level.

3.3 How political polls are reported

Political polls are often reported as a value with a margin-of-error. For example you might hear

52% favor candidate A with a margin-of-error of $\pm 5\%$.

The actual precise meaning of this is

if θ is the proportion of the population that supports A then the point estimate for θ is 52% and the 95% confidence interval is $52\% \pm 5\%$.

Notice that reporters of polls in the news do not mention the 95% confidence. You just have to know that that’s what pollsters do.

The 95% rule-of-thumb confidence interval.

Recall that the $(1 - \alpha)$ conservative normal confidence interval is

$$\bar{x} \pm z_{\alpha/2} \cdot \frac{1}{2\sqrt{n}}.$$

If we use the standard approximation $z_{.025} = 2$ (instead of 1.96) we get the **rule-of thumb 95% confidence interval** for θ :

$$\bar{x} \pm \frac{1}{\sqrt{n}}.$$

Example 2. Polling. Suppose there will soon be a local election between candidate A and candidate B . Suppose that the fraction of the voting population that supports A is θ .

Two polling organizations ask voters who they prefer.

1. The firm of *Fast and First* polls 40 random voters and finds 22 support A .
2. The firm of *Quick but Cautious* polls 400 random voters and finds 190 support A .

Find the point estimates and 95% rule-of-thumb confidence intervals for each poll. Explain how the statistics reflect the intuition that the poll of 400 voters is more accurate.

answer: For poll 1 we have

$$\text{Point estimate: } \bar{x} = 22/40 = 0.55$$

$$\text{Confidence interval: } \bar{x} \pm \frac{1}{\sqrt{n}} = 0.55 \pm \frac{1}{\sqrt{40}} = 0.55 \pm 0.16 = 55\% \pm 16\%.$$

For poll 2 we have

$$\text{Point estimate: } \bar{x} = 190/400 = 0.475$$

$$\text{Confidence interval: } \bar{x} \pm \frac{1}{\sqrt{n}} = 0.475 \pm \frac{1}{\sqrt{400}} = 0.475 \pm 0.05 = 47.5\% \pm 5\%.$$

The greater accuracy of the poll of 400 voters is reflected in the smaller margin of error, i.e. 5% for the poll of 400 voters vs. 16% for the poll of 40 voters.

Other binomial proportion confidence intervals

There are many methods of producing confidence intervals for the proportion p of a binomial(n, p) distribution. For a number of other common approaches, see:

http://en.wikipedia.org/wiki/Binomial_proportion_confidence_interval

4 Large sample confidence intervals

One typical goal in statistics is to estimate the mean of a distribution. When the data follows a normal distribution we could use confidence intervals based on standardized statistics to estimate the mean.

But suppose the data x_1, x_2, \dots, x_n is drawn from a distribution with pmf or pdf $f(x)$ that may not be normal or even parametric. If the distribution has finite mean and variance and if n is sufficiently large, then the following version of the central limit theorem shows we can still use a standardized statistic.

Central Limit Theorem: For large n , the sampling distribution of the studentized mean is approximately standard normal: $\frac{\bar{x} - \mu}{s/\sqrt{n}} \approx N(0, 1)$.

So for large n the $(1 - \alpha)$ confidence interval for μ is approximately

$$\left[\bar{x} - \frac{s}{\sqrt{n}} \cdot z_{\alpha/2}, \bar{x} + \frac{s}{\sqrt{n}} \cdot z_{\alpha/2} \right]$$

where $z_{\alpha/2}$ is the $\alpha/2$ critical value for $N(0, 1)$. This is called the [large sample confidence interval](#).

Example 3. How large must n be?

Recall that a type 1 CI error occurs when the confidence interval does not contain the true value of the parameter, in this case the mean. Let's call the value $(1 - \alpha)$ the *nominal* confidence level. We say nominal because unless n is large we shouldn't expect the true type 1 CI error rate to be α .

We can run numerical simulations to approximate of the true confidence level. We expect that as n gets larger the true confidence level of the large sample confidence interval will converge to the nominal value.

We ran such simulations for x drawn from the exponential distribution $\exp(1)$ (which is far from normal). For several values of n and nominal confidence level c we ran 100,000 trials. Each trial consisted of the following steps:

1. draw n samples from $\exp(1)$.
2. compute the sample mean \bar{x} and sample standard deviation s .
3. construct the large sample c confidence interval: $\bar{x} \pm z_{\alpha/2} \cdot \frac{s}{\sqrt{n}}$.
4. check for a type 1 CI error, i.e. see if the true mean $\mu = 1$ is not in the interval.

With 100,000 trials, the empirical confidence level should closely approximate the true level. For comparison we ran the same tests on data drawn from a standard normal distribution. Here are the results.

n	nominal conf. $1 - \alpha$	simulated conf.
20	0.95	0.905
20	0.90	0.856
20	0.80	0.762
50	0.95	0.930
50	0.90	0.879
50	0.80	0.784
100	0.95	0.938
100	0.90	0.889
100	0.80	0.792
400	0.95	0.947
400	0.90	0.897
400	0.80	0.798

Simulations for $\exp(1)$

n	nominal conf. $1 - \alpha$	simulated conf.
20	0.95	0.936
20	0.90	0.885
20	0.80	0.785
50	0.95	0.944
50	0.90	0.894
50	0.80	0.796
100	0.95	0.947
100	0.90	0.896
100	0.80	0.797
400	0.95	0.949
400	0.90	0.898
400	0.80	0.798

Simulations for $N(0, 1)$.

For the $\exp(1)$ distribution we see that for $n = 20$ the simulated confidence of the large sample confidence interval is less than the nominal confidence $1 - \alpha$. But for $n = 100$ the

simulated confidence and nominal confidence are quite close. So for $\exp(1)$, n somewhere between 50 and 100 is large enough for most purposes.

Think: For $n = 20$ why is the simulated confidence for the $N(0, 1)$ distribution is smaller than the nominal confidence?

This is because we used $z_{\alpha/2}$ instead of $t_{\alpha/2}$. For large n these are quite close, but for $n = 20$ there is a noticeable difference, e.g. $z_{.025} = 1.96$ and $t_{.025} = 2.09$.

MIT OpenCourseWare
<https://ocw.mit.edu>

18.05 Introduction to Probability and Statistics
Spring 2014

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.

Bootstrap confidence intervals

Class 24, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to construct and sample from the empirical distribution of data.
2. Be able to explain the bootstrap principle.
3. Be able to design and run an empirical bootstrap to compute confidence intervals.
4. Be able to design and run a parametric bootstrap to compute confidence intervals.

2 Introduction

The [empirical bootstrap](#) is a statistical technique popularized by Bradley Efron in 1979. Though remarkably simple to implement, the bootstrap would not be feasible without modern computing power. The key idea is to perform computations on the data itself to estimate the variation of statistics that are themselves computed from the same data. That is, the data is ‘pulling itself up by its own bootstrap.’ (A google search of ‘by ones own bootstraps’ will give you the etymology of this metaphor.) Such techniques existed before 1979, but Efron widened their applicability and demonstrated how to implement the bootstrap effectively using computers. He also coined the term ‘bootstrap’ ¹.

Our main application of the bootstrap will be to estimate the variation of point estimates; that is, to estimate confidence intervals. An example will make our goal clear.

Example 1. Suppose we have data

$$x_1, x_2, \dots, x_n$$

If we knew the data was drawn from $N(\mu, \sigma^2)$ with the unknown mean μ and known variance σ^2 then we have seen that

$$\left[\bar{x} - 1.96 \frac{\sigma}{\sqrt{n}}, \bar{x} + 1.96 \frac{\sigma}{\sqrt{n}} \right]$$

is a 95% confidence interval for μ .

Now suppose the data is drawn from some completely unknown distribution. To have a name we’ll call this distribution F and its (unknown) mean μ . We can still use the sample mean \bar{x} as a [point estimate](#) of μ . But how can we find a confidence interval for μ around \bar{x} ? Our answer will be to use the bootstrap!

In fact, we’ll see that the bootstrap handles other statistics as easily as it handles the mean. For example: the median, other percentiles or the trimmed mean. These are statistics where, even for normal distributions, it can be difficult to compute a confidence interval from theory alone.

¹Paraphrased from Dekking et al. *A Modern Introduction to Probability and Statistics*, Springer, 2005, page 275.

3 Sampling

In statistics to **sample** from a set is to choose elements from that set. In a random sample the elements are chosen randomly. There are two common methods for random sampling.

Sampling without replacement

Suppose we draw 10 cards at random from a deck of 52 cards without putting any of the cards back into the deck between draws. This is called **sampling without replacement** or **simple random sampling**. With this method of sampling our 10 card sample will have no duplicate cards.

Sampling with replacement

Now suppose we draw 10 cards at random from the deck, but after each draw we put the card back in the deck and shuffle the cards. This is called **sampling with replacement**. With this method, the 10 card sample might have duplicates. It's even possible that we would draw the 6 of hearts all 10 times.

Think: What's the probability of drawing the 6 of hearts 10 times in a row?

Example 2. We can view rolling an 8-sided die repeatedly as sampling with replacement from the set $\{1,2,3,4,5,6,7,8\}$. Since each number is equally likely, we say we are sampling uniformly from the data. There is a subtlety here: each data point is equally probable, but if there are repeated values within the data those values will have a higher probability of being chosen. The next example illustrates this.

Note. In practice if we take a small number from a very large set then it doesn't matter whether we sample with or without replacement. For example, if we randomly sample 400 out of 300 million people in the U.S. then it is so unlikely that the same person will be picked twice that there is no real difference between sampling with or without replacement.

4 The empirical distribution of data

The empirical distribution of data is simply the distribution that you see in the data. Let's illustrate this with an example.

Example 3. Suppose we roll an 8-sided die 10 times and get the following data, written in increasing order:

$$1, 1, 2, 3, 3, 3, 3, 4, 7, 7.$$

Imagine writing these values on 10 slips of paper, putting them in a hat and drawing one at random. Then, for example, the probability of drawing a 3 is $4/10$ and the probability of drawing a 4 is $1/10$. The full empirical distribution can be put in a probability table

value x	1	2	3	4	7
$p(x)$	2/10	1/10	4/10	1/10	2/10

Notation. If we label the true distribution the data is drawn from as F , then we'll label the empirical distribution of the data as F^* . If we have enough data then the law of large numbers tells us that F^* should be a good approximation of F .

Example 4. In the dice example just above, the true and empirical distributions are:

value x	1	2	3	4	5	5	7	8
true $p(x)$	1/8	1/8	1/8	1/8	1/8	1/8	1/8	1/8
empirical $p(x)$	2/10	1/10	4/10	1/10	0	0	2/10	0

The true distribution F and the empirical distribution F^* of the 8-sided die.

Because F^* is derived strictly from data we call it the **empirical distribution** of the data. We will also call it the **resampling distribution**. Notice that we always know F^* explicitly. In particular the expected value of F^* is just the sample mean \bar{x} .

5 Resampling

The empirical bootstrap proceeds by resampling from the data. We continue the dice example above.

Example 5. Suppose we have 10 data points, given in increasing order:

$$1, 1, 2, 3, 3, 3, 3, 4, 7, 7$$

We view this as a **sample** taken from some underlying distribution. To **resample** is to sample with replacement from the empirical distribution, e.g. put these 10 numbers in a hat and draw one at random. Then put the number back in the hat and draw again. You draw as many numbers as the desired size of the resample.

To get us a little closer to implementing this on a computer we rephrase this in the following way. Label the 10 data points x_1, x_2, \dots, x_{10} . To resample is to draw a number j from the uniform distribution on $\{1, 2, \dots, 10\}$ and take x_j as our resampled value. In this case we could do so by rolling a 10-sided die. For example, if we roll a 6 then our resampled value is 3, the 6th element in our list.

If we want a resampled data set of size 5, then we roll the 10-sided die 5 times and choose the corresponding elements from the list of data. If the 5 rolls are

$$5, 3, 6, 6, 1$$

then the resample is

$$3, 2, 3, 3, 1.$$

Notes: 1. Because we are sampling with replacement, the same data point can appear multiple times when we resample.

2. Also because we are sampling with replacement, we can have a resample data set of any size we want, e.g. we could resample 1000 times.

Of course, in practice one uses a software package like R to do the resampling.

5.1 Star notation

If we have sample data of size n

$$x_1, x_2, \dots, x_n$$

then we denote a [resample of size \$m\$](#) by adding a star to the symbols

$$x_1^*, x_2^*, \dots, x_m^*$$

Similarly, just as \bar{x} is the mean of the original data, we write \bar{x}^* for the mean of the [resampled data](#).

6 The empirical bootstrap

Suppose we have n data points

$$x_1, x_2, \dots, x_n$$

drawn from a distribution F . An [empirical bootstrap sample](#) is a resample of the [same size \$n\$](#) :

$$x_1^*, x_2^*, \dots, x_n^*.$$

You should think of the latter as a sample of size n drawn from the empirical distribution F^* . For any statistic v computed from the original sample data, we can define a statistic v^* by the same formula but computed instead using the resampled data. With this notation we can state the bootstrap principle.

6.1 The bootstrap principle

The bootstrap setup is as follows:

1. x_1, x_2, \dots, x_n is a data sample drawn from a distribution F .
2. u is a statistic computed from the sample.
3. F^* is the empirical distribution of the data (the resampling distribution).
4. $x_1^*, x_2^*, \dots, x_n^*$ is a resample of the data [of the same size](#) as the original sample
5. u^* is the statistic computed from the resample.

Then the [bootstrap principle](#) says that

1. $F^* \approx F$.
2. The variation of u is well-approximated by the variation of u^* .

Our real interest is in point 2: we can approximate the variation of u by that of u^* . We will exploit this to estimate the size of confidence intervals.

6.2 Why the resample is the same size as the original sample

This is straightforward: the variation of the statistic u will depend on the size of the sample. If we want to approximate this variation we need to use resamples of the same size.

6.3 Toy example of an empirical bootstrap confidence interval

Example 6. [Toy example](#). We start with a made-up set of data that is small enough to show each step explicitly. The sample data is

$$30, 37, 36, 43, 42, 43, 43, 46, 41, 42$$

Problem: Estimate the mean μ of the underlying distribution and give an 80% bootstrap confidence interval.

Note: R code for this example is shown in the section ‘R annotated transcripts’ below. The code is also implemented in the R script `class24-empiricalbootstrap.r` which is posted with our other R code.

answer: The sample mean is $\bar{x} = 40.3$. We use this as an estimate of the true mean μ of the underlying distribution. As in Example 1, to make the confidence interval we need to know how much the distribution of \bar{x} varies around μ . That is, we’d like to know the distribution of

$$\delta = \bar{x} - \mu.$$

If we knew this distribution we could find $\delta_{.1}$ and $\delta_{.9}$, the 0.1 and 0.9 critical values of δ . Then we’d have

$$P(\delta_{.9} \leq \bar{x} - \mu \leq \delta_{.1} | \mu) = 0.8 \Leftrightarrow P(\bar{x} - \delta_{.9} \geq \mu \geq \bar{x} - \delta_{.1} | \mu) = 0.8$$

which gives an 80% confidence interval of

$$[\bar{x} - \delta_{.1}, \bar{x} - \delta_{.9}].$$

As always with confidence intervals, we hasten to point out that the probabilities computed above are probabilities concerning the statistic \bar{x} given that the true mean is μ .

The bootstrap principle offers a practical approach to estimating the distribution of $\delta = \bar{x} - \mu$. It says that we can approximate it by the distribution of

$$\delta^* = \bar{x}^* - \bar{x}$$

where \bar{x}^* is the mean of an empirical bootstrap sample.

Here’s the beautiful key to this: since δ^* is computed by resampling the original data, we can have a computer simulate δ^* as many times as we’d like. Hence, by the law of large numbers, we can estimate the distribution of δ^* with high precision.

Now let’s return to the sample data with 10 points. We used R to generate 20 bootstrap samples, each of size 10. Each of the 20 columns in the following array is one bootstrap sample.

43	36	46	30	43	43	43	37	42	42	43	37	36	42	43	43	42	43	42	43
43	41	37	37	43	43	46	36	41	43	43	42	41	43	46	36	43	43	43	42
42	43	37	43	46	37	36	41	36	43	41	36	37	30	46	46	42	36	36	43
37	42	43	41	41	42	36	42	42	43	42	43	41	43	36	43	43	41	42	46
42	36	43	43	42	37	42	42	42	46	30	43	36	43	43	42	37	36	42	30
36	36	42	42	36	36	43	41	30	42	37	43	41	41	43	43	42	46	43	37
43	37	41	43	41	42	43	46	46	36	43	42	43	30	41	46	43	46	30	43
41	42	30	42	37	43	43	42	43	43	46	43	30	42	30	42	30	43	43	42
46	42	42	43	41	42	30	37	30	42	43	42	37	37	37	42	43	43	46	
42	43	43	41	42	36	43	30	37	43	42	43	41	36	37	41	43	42	43	43

Next we compute $\delta^* = \bar{x}^* - \bar{x}$ for each bootstrap sample (i.e. each column) and sort them from smallest to biggest:

-1.6, -1.4, -1.4, -0.9, -0.5, -0.2, -0.1, 0.1, 0.2, 0.2, 0.2, 0.4, 0.4, 0.4, 0.7, 0.9, 1.1, 1.2, 1.2, 1.6, 1.6, 2.0

We will approximate the critical values $\delta_{.1}$ and $\delta_{.9}$ by $\delta_{.1}^*$ and $\delta_{.9}^*$. Since $\delta_{.1}^*$ is at the 90th percentile we choose the 18th element in the list, i.e. 1.6. Likewise, since $\delta_{.9}^*$ is at the 10th percentile we choose the 2nd element in the list, i.e. -1.4.

Therefore our bootstrap 80% confidence interval for μ is

$$[\bar{x} - \delta_{.1}^*, \bar{x} - \delta_{.9}^*] = [40.3 - 1.6, 40.3 + 1.4] = [38.7, 41.7]$$

In this example we only generated 20 bootstrap samples so they would fit on the page. Using R, we would generate 10000 or more bootstrap samples in order to obtain a very accurate estimate of $\delta_{.1}^*$ and $\delta_{.9}^*$.

6.4 Justification for the bootstrap principle

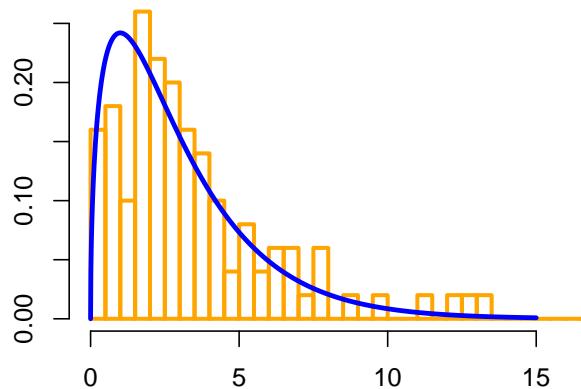
The bootstrap is remarkable because resampling gives us a decent estimate on how the point estimate might vary. We can only give you a ‘hand-waving’ explanation of this, but it’s worth a try. The bootstrap is based roughly on the law of large numbers, which says, in short, that with enough data the empirical distribution will be a good approximation of the true distribution. Visually it says that the histogram of the data should approximate the density of the true distribution.

First let’s note what resampling can’t do for us: it can’t improve our point estimate. For example, if we estimate the mean μ by \bar{x} then in the bootstrap we would compute \bar{x}^* for many resamples of the data. If we took the average of all the \bar{x}^* we would expect it to be very close to \bar{x} . This wouldn’t tell us anything new about the true value of μ .

Even with a fair amount of data the match between the true and empirical distributions is not perfect, so there will be error in estimating the mean (or any other value). But the amount of variation in the estimates is much less sensitive to differences between the density and the histogram. As long as they are reasonably close both the empirical and true distributions will exhibit the similar amounts of variation. So, in general the bootstrap principle is more robust when approximating the distribution of relative variation than when approximating absolute distributions.

What we have in mind is the scenario of our examples. The distribution (over different sets of experimental data) of \bar{x} is ‘centered’ at μ and the distribution of \bar{x}^* is centered at \bar{x} . If there is a significant separation between \bar{x} and μ then these two distributions will also differ significantly. On the other hand the distribution of $\delta = \bar{x} - \mu$ describes the variation of \bar{x} about its center. Likewise the distribution of $\delta^* = \bar{x}^* - \bar{x}$ describes the variation of \bar{x}^* about its center. So even if the centers are quite different the two variations about the centers can be approximately equal.

The figure below illustrates how the empirical distribution approximates the true distribution. To make the figure we generate 100 random values from a chi-square distribution with 3 degrees of freedom. The figure shows the pdf of the true distribution as a blue line and a histogram of the empirical distribution in orange.



The true and empirical distributions are approximately equal.

7 Other statistics

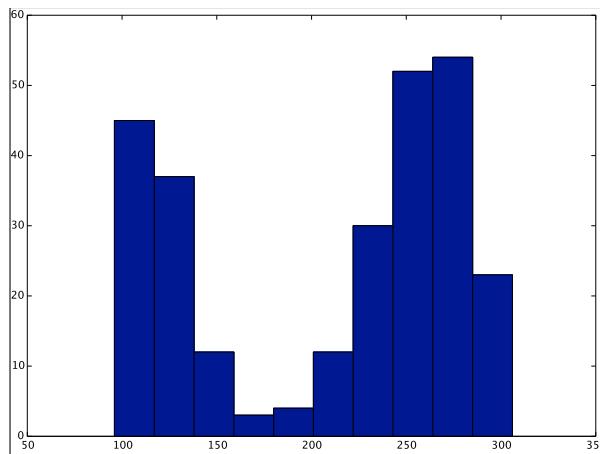
So far in this class we've avoided confidence intervals for the median and other statistics because their sample distributions are hard to describe theoretically. The bootstrap has no such problem. In fact, to handle the median all we have to do is change 'mean' to 'median' in the R code from Example 6.

Example 7. Old Faithful: confidence intervals for the median

Old Faithful is a geyser in Yellowstone National Park in Wyoming:

http://en.wikipedia.org/wiki/Old_Faithful

There is a publicly available data set which gives the durations of 272 consecutive eruptions. Here is a histogram of the data.



Question: Estimate the median length of an eruption and give a 90% confidence interval for the median.

answer: The full answer to this question is in the R file `oldfaithful_simple.r` and the Old Faithful data set. Both are posted on the class R code page. (Look under 'Other R code' for the old faithful script and data.)

Note: the code in `oldfaithful_simple.r` assumes that the data `oldfaithful.txt` is in the current working directory.

Let's walk through a summary of the steps needed to answer the question.

1. Data: x_1, \dots, x_{272}
2. Data median: $x_{\text{median}} = 240$
3. Find the median x_{median}^* of a bootstrap sample x_1^*, \dots, x_{272}^* . Repeat 1000 times.
4. Compute the bootstrap differences

$$\delta^* = x_{\text{median}}^* - x_{\text{median}}$$

Put these 1000 values in order and pick out the .95 and .05 critical values, i.e. the 50th and 950th biggest values. Call these $\delta_{.95}^*$ and $\delta_{.05}^*$.

5. The bootstrap principle says that we can use $\delta_{.95}^*$ and $\delta_{.05}^*$ as estimates of $\delta_{.95}$ and $\delta_{.05}$. So our estimated 90% bootstrap confidence interval for the median is

$$[x_{\text{median}} - \delta_{.05}^*, x_{\text{median}} - \delta_{.95}^*]$$

The bootstrap 90% CI we found for the Old Faithful data was [235, 250]. Since we used 1000 bootstrap samples a new simulation starting from the same sample data should produce a similar interval. If in Step 3 we increase the number of bootstrap samples to 10000, then the intervals produced by simulation would vary even less. One common strategy is to increase the number of bootstrap samples until the resulting simulations produce intervals that vary less than some acceptable level.

Example 8. Using the Old Faithful data, estimate $P(|\bar{x} - \mu| > 5 | \mu)$.

answer: We proceed exactly as in the previous example except using the mean instead of the median.

1. Data: x_1, \dots, x_{272}
2. Data mean: $\bar{x} = 209.27$
3. Find the mean \bar{x}^* of 1000 empirical bootstrap samples: x_1^*, \dots, x_{272}^* .
4. Compute the bootstrap differences

$$\delta^* = \bar{x}^* - \bar{x}$$

5. The bootstrap principle says that we can use the distribution of δ^* as an approximation for the distribution $\delta = \bar{x} - \mu$. Thus,

$$P(|\bar{x} - \mu| > 5 | \mu) = P(|\delta| > 5 | \mu) \approx P(|\delta^*| > 5)$$

Our bootstrap simulation for the Old Faithful data gave 0.225 for this probability.

8 Parametric bootstrap

The examples in the previous sections all used the empirical bootstrap, which makes no assumptions at all about the underlying distribution and draws bootstrap samples by resampling the data. In this section we will look at the [parametric bootstrap](#). The only difference

between the parametric and empirical bootstrap is the source of the bootstrap sample. For the parametric bootstrap, we generate the bootstrap sample from a parametrized distribution.

Here are the elements of using the parametric bootstrap to estimate a confidence interval for a parameter.

0. Data: x_1, \dots, x_n drawn from a distribution $F(\theta)$ with unknown parameter θ .
1. A statistic $\hat{\theta}$ that estimates θ .
2. Our bootstrap samples are drawn from $F(\hat{\theta})$.
3. For each bootstrap sample

$$x_1^*, \dots, x_n^*$$

we compute $\hat{\theta}^*$ and the bootstrap difference $\delta^* = \hat{\theta}^* - \hat{\theta}$.

4. The bootstrap principle says that the distribution of δ^* approximates the distribution of $\delta = \hat{\theta} - \theta$.

5. Use the bootstrap differences to make a bootstrap confidence interval for θ .

Example 9. Suppose the data x_1, \dots, x_{300} is drawn from an $\exp(\lambda)$ distribution. Assume also that the data mean $\bar{x} = 2$. Estimate λ and give a 95% parametric bootstrap confidence interval for λ .

answer: This is implemented in the R script `class24-parametricbootstrap.r` which is posted with our other R code.

It's will be easiest to explain the solution using commented code.

```
# Parametric bootstrap

# Given 300 data points with mean 2.
# Assume the data is exp(lambda)

# PROBLEM: Compute a 95% parametric bootstrap confidence interval for lambda

# We are given the number of data points and mean
n = 300
xbar = 2

# The MLE for lambda is 1/xbar
lambdahat = 1.0/xbar

# Generate the bootstrap samples
# Each column is one bootstrap sample (of 300 resampled values)
nboot = 1000

# Here's the key difference with the empirical bootstrap:
# We draw the bootstrap sample from Exponential(lambdahat)
x = rexp(n*nboot, lambdahat)
bootstrapsample = matrix(x, nrow=n, ncol=nboot)

# Compute the bootstrap lambdastar
lambdastar = 1.0/colMeans(bootstrapsample)

# Compute the differences
deltastar = lambdastar - lambdahat
```

```

# Find the 0.05 and 0.95 quantile for deltastar
d = quantile(deltastar, c(0.05,0.95))

# Calculate the 95% confidence interval for lambda.
ci = lambdahat - c(d[2], d[1])

# This line of code is just one way to format the output text.
# sprintf is an old C function for doing this. R has many other
# ways to do the same thing.
s = sprintf("Confidence interval for lambda: [% .3f, % .3f]", ci[1], ci[2])
cat(s)

```

9 The bootstrap percentile method (should not be used)

Instead of computing the differences δ^* , the bootstrap percentile method uses the distribution of the bootstrap sample statistic as a direct approximation of the data sample statistic.

Example 10. Let's redo Example 6 using the bootstrap percentile method.

We first compute \bar{x}^* from the bootstrap samples given in Example 6. After sorting we get

```
35.7 37.4 38.0 39.5 39.7 39.8 39.8 40.1 40.1 40.6 40.7 40.8 41.1 41.1 41.7 42.0
42.1 42.4 42.4 42.4
```

The percentile method says to use the distribution of \bar{x}^* as an approximation to the distribution of \bar{x} . The 0.9 and 0.1 critical values are given by the 2nd and 18th elements. Therefore the 80% confidence interval is [37.4, 42.4]. This is a bit wider than our answer to Example 6.

The bootstrap percentile method is appealing due to its simplicity. However it depends on the bootstrap distribution of \bar{x}^* based on a [particular](#) sample being a good approximation to the true distribution of \bar{x} . Rice says of the percentile method, “Although this direct equation of quantiles of the bootstrap sampling distribution with confidence limits may seem initially appealing, its rationale is somewhat obscure.”² In short, [don't use the bootstrap percentile method](#). Use the empirical bootstrap instead (we have explained both in the hopes that you won't confuse the empirical bootstrap for the percentile bootstrap).

10 R annotated transcripts

10.1 Using R to generate an empirical bootstrap confidence interval

This code only generates 20 bootstrap samples. In real practice we would generate many more bootstrap samples. It is making a bootstrap confidence interval for the mean. This code is implemented in the R script `class24-empiricalbootstrap.r` which is posted with our other R code.

```

# Data for the example 6
x = c(30,37,36,43,42,43,43,46,41,42)
n = length(x)

```

²John Rice, *Mathematical Statistics and Data Analysis*, 2nd edition, p. 272.

```
# sample mean
xbar = mean(x)

nboot = 20
# Generate 20 bootstrap samples, i.e. an n x 20 array of
# random resamples from x
tmpdata = sample(x,n*nboot, replace=TRUE)
bootstrapsample = matrix(tmpdata, nrow=n, ncol=nboot)

# Compute the means  $\bar{x}^*$ 
bsmeans = colMeans(bootstrapsample)

# Compute  $\delta^*$  for each bootstrap sample
deltastar = bsmeans - xbar

# Find the 0.1 and 0.9 quantile for deltastar
d = quantile(deltastar, c(0.1, 0.9))

# Calculate the 80% confidence interval for the mean.
ci = xbar - c(d[2], d[1])
cat('Confidence interval: ', ci, '\n')

# ALTERNATIVE: the quantile() function is sophisticated about
# choosing a quantile between two data points. A less sophisticated
# approach is to pick the quantiles by sorting deltastar and
# choosing the index that corresponds to the desired quantiles.
# We do this below.

# Sort the results
sorteddeltastar = sort(deltastar)

# Look at the sorted results
hist(sorteddeltastar, nclass=6)
print(sorteddeltastar)

# Find the .1 and .9 critical values of deltastar
d9alt = sorteddeltastar[2]
d1alt = sorteddeltastar[18]

# Find and print the 80% confidence interval for the mean
ciAlt = xbar - c(d1alt,d9alt)
cat('Alternative confidence interval: ', ciAlt, '\n')
```

MIT OpenCourseWare
<https://ocw.mit.edu>

18.05 Introduction to Probability and Statistics
Spring 2014

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.

Linear regression

Class 25, 18.05

Jeremy Orloff and Jonathan Bloom

1 Learning Goals

1. Be able to use the method of least squares to fit a line to bivariate data.
2. Be able to give a formula for the total squared error when fitting any type of curve to data.
3. Be able to say the words homoscedasticity and heteroscedasticity.

2 Introduction

Suppose we have collected bivariate data (x_i, y_i) , $i = 1, \dots, n$. The goal of linear regression is to model the relationship between x and y by finding a function $y = f(x)$ that is a close fit to the data. The modeling assumptions we will use are that x_i is **not** random and that y_i is a function of x_i plus some random noise. With these assumptions x is called the **independent or predictor variable** and y is called the **dependent or response variable**.

Example 1. The cost of a first class stamp in cents over time is given in the following list.

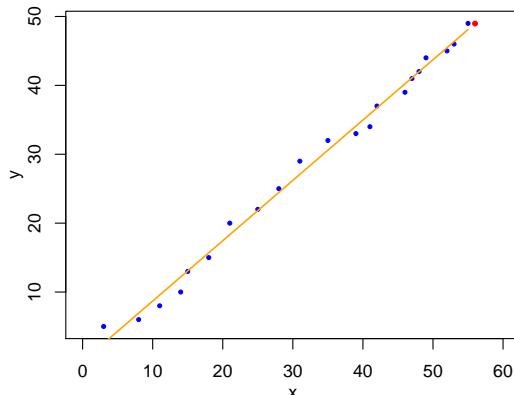
.05 (1963)	.06 (1968)	.08 (1971)	.10 (1974)	.13 (1975)	.15 (1978)	.20 (1981)	.22 (1985)
.25 (1988)	.29 (1991)	.32 (1995)	.33 (1999)	.34 (2001)	.37 (2002)	.39 (2006)	.41 (2007)
.42 (2008)	.44 (2009)	.45 (2012)	.46 (2013)	.49 (2014)			

Using the R function `lm` we found the ‘least squares fit’ for a line to this data is

$$y = -0.06558 + 0.87574x,$$

where x is the number of years since 1960 and y is in cents.

Using this result we ‘predict’ that in 2016 ($x = 56$) the cost of a stamp will be 49 cents (since $-0.06558 + 0.87574 \cdot 56 = 48.98$).

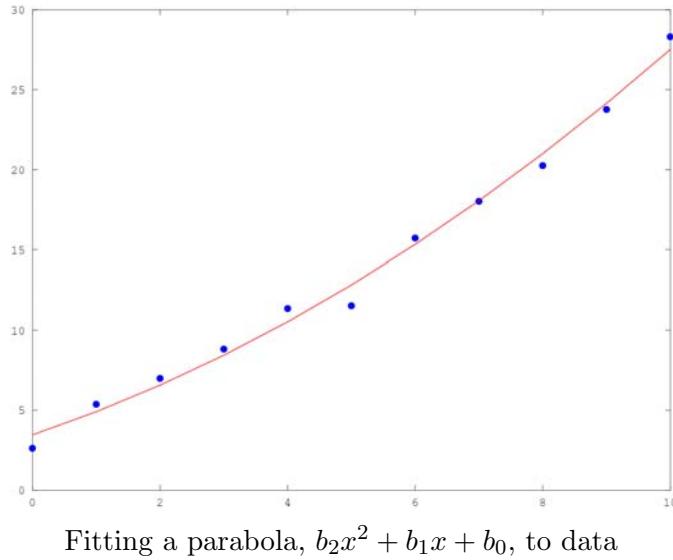


Stamp cost (cents) vs. time (years since 1960). Red dot is predicted cost in 2016.

Note that none of the data points actually lie on the line. Rather this line has the ‘best fit’ with respect to [all the data](#), with a small error for each data point.

Example 2. Suppose we have n pairs of fathers and adult sons. Let x_i and y_i be the heights of the i^{th} father and son, respectively. The least squares line for this data could be used to predict the adult height of a young boy from that of his father.

Example 3. We are not limited to best fit lines. For all positive d , the method of least squares may be used to find a polynomial of degree d with the ‘best fit’ to the data. Here’s a figure showing the least squares fit of a parabola ($d = 2$).



3 Fitting a line using least squares

Suppose we have data (x_i, y_i) as above. The goal is to find the line

$$y = \beta_1 x + \beta_0$$

that ‘best fits’ the data. Our model says that each y_i is predicted by x_i up to some error ϵ_i :

$$y_i = \beta_1 x_i + \beta_0 + \epsilon_i.$$

So

$$\epsilon_i = y_i - \beta_1 x_i - \beta_0.$$

The method of least squares finds the values $\hat{\beta}_0$ and $\hat{\beta}_1$ of β_0 and β_1 that minimize the sum of the squared errors:

$$S(\beta_0, \beta_1) = \sum_i \epsilon_i^2 = \sum_i (y_i - \beta_1 x_i - \beta_0)^2.$$

Using calculus or linear algebra (details in the appendix), we find

$$\hat{\beta}_1 = \frac{s_{xy}}{s_{xx}} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \tag{1}$$

where

$$\bar{x} = \frac{1}{n} \sum x_i, \quad \bar{y} = \frac{1}{n} \sum y_i, \quad s_{xx} = \frac{1}{(n-1)} \sum (x_i - \bar{x})^2, \quad s_{xy} = \frac{1}{(n-1)} \sum (x_i - \bar{x})(y_i - \bar{y}).$$

Here \bar{x} is the sample mean of x , \bar{y} is the sample mean of y , s_{xx} is the sample variance of x , and s_{xy} is the sample covariance of x and y .

Example 4. Use least squares to fit a line to the following data: (0,1), (2,1), (3,4).

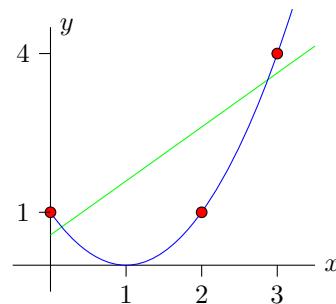
answer: In our case, $(x_1, y_1) = (0, 1)$, $(x_2, y_2) = (2, 1)$ and $(x_3, y_3) = (3, 4)$. So

$$\bar{x} = \frac{5}{3}, \quad \bar{y} = 2, \quad s_{xx} = \frac{14}{9}, \quad s_{xy} = \frac{4}{3}$$

Using the above formulas we get

$$\hat{\beta}_1 = \frac{6}{7}, \quad \hat{\beta}_0 = \frac{4}{7}.$$

So the least squares line has equation $y = \frac{6}{7} + \frac{6}{7}x$. This is shown as the green line in the following figure. We will discuss the blue parabola soon.



Least squares fit of a line (green) and a parabola (blue)

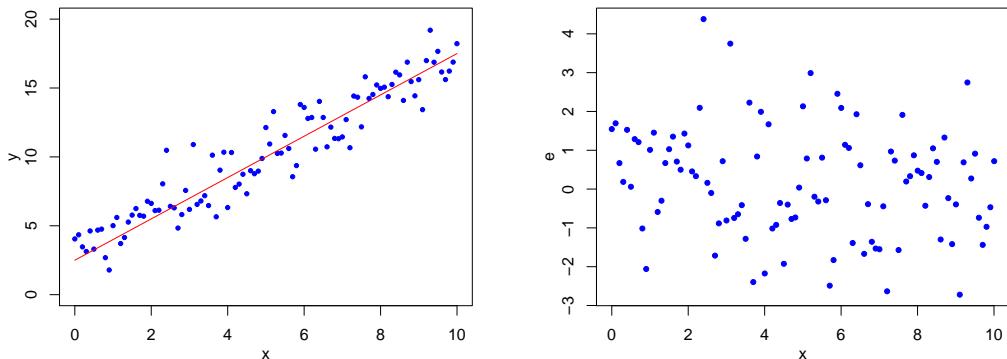
Simple linear regression: It's a little confusing, but the word linear in 'linear regression' does not refer to fitting a line. We will explain its meaning below. However, the most common curve to fit is a line. When we fit a line to bivariate data it is called [simple linear regression](#).

3.1 Residuals

For a line the model is

$$y_i = \hat{\beta}_1 x_i + \hat{\beta}_0 + \epsilon_i.$$

We think of $\hat{\beta}_1 x_i + \hat{\beta}_0$ as predicting or explaining y_i . The left-over term ϵ_i is called the [residual](#), which we think of as random noise or measurement error. A useful visual check of the linear regression model is to plot the residuals. The data points should hover near the regression line. The residuals should look about the same across the range of x .

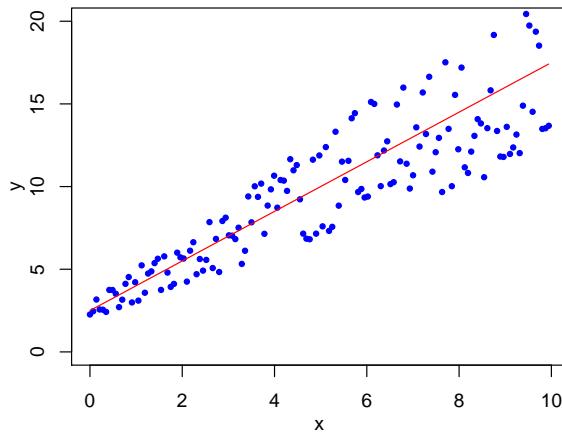


Data with regression line (left) and residuals (right). Note the homoscedasticity.

3.2 Homoscedasticity

An important assumption of the linear regression model is that the residuals ϵ_i have the same variance for all i . This is called [homoscedasticity](#). You can see this is the case for both figures above. The data hovers in the band of fixed width around the regression line and at every x the residuals have about the same vertical spread.

Below is a figure showing [heteroscedastic](#) data. The vertical spread of the data increases as x increases. Before using least squares on this data we would need to transform the data to be homoscedastic.



Heteroscedastic Data

4 Linear regression for fitting polynomials

When we fit a line to data it is called [simple linear regression](#). We can also use linear regression to fit polynomials to data. The use of the word linear in both cases may seem confusing. This is because the word ‘linear’ in linear regression does not refer to fitting a line. Rather it refers to the linear algebraic equations for the unknown parameters β_i , i.e. each β_i has exponent 1.

Example 5. Take the same data as in the Example 4 and use least squares to find the

best fitting parabola to the data.

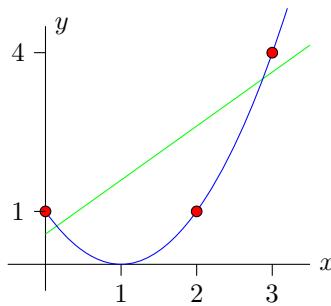
answer: A parabola has the formula $y = \beta_0 + \beta_1 x + \beta_2 x^2$. The squared error is

$$S(\beta_0, \beta_1, \beta_2) = \sum (y_i - (\beta_0 + \beta_1 x_i + \beta_2 x_i^2))^2.$$

After substituting the given values for each x_i and y_i , we can use calculus to find the triple $(\beta_0, \beta_1, \beta_2)$ that minimizes S . With this data, we find that the least squares parabola has equation

$$y = 1 - 2x + x^2.$$

Note that for 3 points the quadratic fit is perfect.



Least squares fit of a line (green) and a parabola (blue)

Example 6. The pairs (x_i, y_i) may give the age and vocabulary size of n children. Since we expect that young children acquire new words at an accelerating pace, we might guess that a higher order polynomial might best fit the data.

Example 7. (Transforming the data) Sometimes it is necessary to transform the data before using linear regression. For example, let's suppose the relationship is exponential, i.e. $y = ce^{ax}$. Then

$$\ln(y) = ax + \ln(c).$$

So we can use simple linear regression to obtain a model

$$\ln(y_i) = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

and then exponentiate to obtain the exponential model

$$y_i = e^{\hat{\beta}_0} e^{\hat{\beta}_1 x_i}.$$

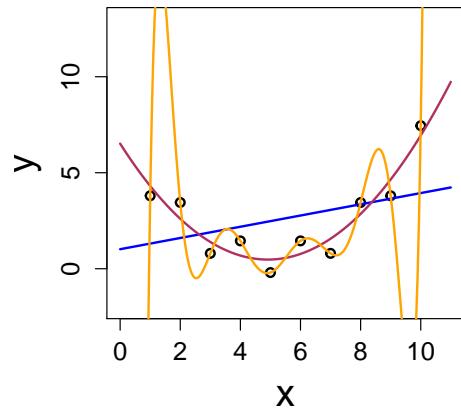
4.1 Overfitting

You can always achieve a better fit by using a higher order polynomial. For instance, given 6 data points (with distinct x_i) one can always find a fifth order polynomial that goes through all of them. This can result in what's called **overfitting**. That is, fitting the noise as well as the true relationship between x and y . An overfit model will fit the original data better but perform less well on predicting y for new values of x . Indeed, a primary challenge of statistical modeling is balancing model fit against model complexity.

Example 8. In the plot below, we fit polynomials of degree 1, 3, and 9 to bivariate data consisting of 10 data points. The degree 2 model (maroon) gives a significantly better fit

than the degree 1 model (blue). The degree 10 model (orange) gives fits the data exactly, but at a glance we would guess it is overfit. That is, we don't expect it to do a good job fitting the next data point we see.

In fact, we generated this data using a quadratic model, so the degree 2 model will tend to perform best fitting new data points.



4.2 R function lm

As you would expect we don't actually do linear regression by hand. Computationally, linear regression reduces to solving simultaneous equations, i.e. to matrix calculations. The R function `lm` can be used to fit any order polynomial to data. ([lm stands for linear model](#)). We will explore this in the next studio class. In fact `lm` can fit many types of functions besides polynomials, as you can explore using R help or google.

5 Multiple linear regression

Data is not always bivariate. It can be trivariate or even of some higher dimension. Suppose we have data in the form of tuples

$$(y_i, x_{1,i}, x_{2,i}, \dots x_{m,i})$$

We can analyze this in a manner very similar to linear regression on bivariate data. That is, we can use least squares to fit the model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_m x_m.$$

Here each x_j is a predictor variable and y is the response variable. For example, we might be interested in how a fish population varies with measured levels of several pollutants, or we might want to predict the adult height of a son based on the height of the mother and the height of the father.

We don't have time in 18.05 to study multiple linear regression, but we wanted you to see the name.

6 Least squares as a statistical model

The linear regression model for fitting a line says that the value y_i in the pair (x_i, y_i) is drawn from a random variable

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

where the ‘error’ terms ε_i are independent random variables with mean 0 and standard deviation σ . The standard assumption is that the ε_i are i.i.d. with distribution $N(0, \sigma^2)$. In any case, the mean of Y_i is given by:

$$E(Y_i) = \beta_0 + \beta_1 x_i + E(\varepsilon_i) = \beta_0 + \beta_1 x_i.$$

From this perspective, the least squares method chooses the values of β_0 and β_1 which minimize the sample variance about the line.

In fact, the least square estimate $(\hat{\beta}_0, \hat{\beta}_1)$ coincides with the maximum likelihood estimate for the parameters (β_0, β_1) ; that is, among all possible coefficients, $(\hat{\beta}_0, \hat{\beta}_1)$ are the ones that make the observed data most probable.

7 Regression to the mean

The reason for the term ‘regression’ is that the predicted response variable y will tend to be ‘closer’ to (i.e., regress to) its mean than the predictor variable x is to its mean. Here closer is in quotes because we have to control for the scale (i.e. standard deviation) of each variable. The way we control for scale is to first standardize each variable.

$$u_i = \frac{x_i - \bar{x}}{\sqrt{s_{xx}}}, \quad v_i = \frac{y_i - \bar{y}}{\sqrt{s_{yy}}}.$$

Standardization changes the mean to 0 and variance to 1:

$$\bar{u} = \bar{v} = 0, \quad s_{uu} = s_{vv} = 1.$$

The algebraic properties of covariance show

$$s_{uv} = \frac{s_{xy}}{\sqrt{s_{xx}s_{yy}}} = \rho,$$

the correlation coefficient. Thus the least squares fit to $v = \beta_0 + \beta_1 u$ has

$$\hat{\beta}_1 = \frac{s_{uv}}{s_{uu}} = \rho \quad \text{and} \quad \hat{\beta}_0 = \bar{v} - \hat{\beta}_1 \bar{u} = 0.$$

So the least squares line is $v = \rho u$. Since ρ is the correlation coefficient, it is between -1 and 1. Let’s assume it is positive and less than 1 (i.e., x and y are positively but not perfectly correlated). Then the formula $v = \rho u$ means that if u is positive then the predicted value of v is less than u . That is, v is closer to 0 than u . Equivalently,

$$\frac{y - \bar{y}}{\sqrt{s_{yy}}} < \frac{x - \bar{x}}{\sqrt{s_{xx}}}$$

i.e., y regresses to \bar{y} . Notice how the standardization takes care of controlling the scale.

Consider the extreme case of 0 correlation between x and y . Then, no matter what the x value, the predicted value of y is always \bar{y} . That is, y has regressed all the way to its mean.

Note also that the regression line always goes through the point (\bar{x}, \bar{y}) .

Example 9. Regression to the mean is important in longitudinal studies. Rice (*Mathematical Statistics and Data Analysis*) gives the following example. Suppose children are given an IQ test at age 4 and another at age 5 we expect the results will be positively correlated. The above analysis says that, on average, those kids who do poorly on the first test will tend to show improvement (i.e. regress to the mean) on the second test. Thus, a useless intervention might be misinterpreted as useful since it seems to improve scores.

Example 10. Another example with practical consequences is reward and punishment. Imagine a school where high performance on an exam is rewarded and low performance is punished. Regression to the mean tells us that (on average) the high performing students will do slightly worse on the next exam and the low performing students will do slightly better. An unsophisticated view of the data will make it seem that punishment improved performance and reward actually hurt performance. There are real consequences if those in authority act on this idea.

8 Appendix

We collect in this appendix a few things you might find interesting. **You will not be asked to know these things for exams.**

8.1 Proof of the formula for least square fit of a line

The most straightforward proof is to use calculus. The sum of the squared errors is

$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_1 x_i - \beta_0)^2.$$

Taking partial derivatives (and remembering that x_i and y_i are the data, hence constant)

$$\begin{aligned}\frac{\partial S}{\partial \beta_0} &= \sum_{i=1}^n -2(y_i - \beta_1 x_i - \beta_0) = 0 \\ \frac{\partial S}{\partial \beta_1} &= \sum_{i=1}^n -2x_i(y_i - \beta_1 x_i - \beta_0) = 0\end{aligned}$$

Summing this up we get two linear equations in the unknowns β_0 and β_1 :

$$\begin{aligned}\left(\sum x_i\right) \beta_1 + n\beta_0 &= \sum y_i \\ \left(\sum x_i^2\right) \beta_1 + \left(\sum x_i\right) \beta_0 &= \sum x_i y_i\end{aligned}$$

Solving for β_1 and β_0 gives the formulas in Equation (1).

A sneakier approach which avoids calculus is to standardize the data, find the best fit line, and then unstandardize. We omit the details.

For a deluge of applications across disciplines see:

http://en.wikipedia.org/wiki/Linear_regression#Applications_of_linear_regression

8.2 Measuring the fit

Once one computes the regression coefficients, it is important to check how well the regression model fits the data (i.e., how closely the best fit line tracks the data). A common but crude ‘goodness of fit’ measure is the coefficient of determination, denoted R^2 . We’ll need some notation to define it. The total sum of squares is given by:

$$\text{TSS} = \sum (y_i - \bar{y})^2.$$

The residual sum of squares is given by the sum of the squares of the residuals. When fitting a line, this is:

$$\text{RSS} = \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2.$$

The RSS is the “unexplained” portion of the total sum of squares, i.e. unexplained by the regression equation. The difference $\text{TSS} - \text{RSS}$ is the “explained” portion of the total sum of squares. The coefficient of determination R^2 is the ratio of the “explained” portion to the total sum of squares:

$$R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}}.$$

In other words, R^2 measures the proportion of the variability of the data that is accounted for by the regression model. A value close to 1 indicates a good fit, while a value close to 0 indicates a poor fit. In the case of simple linear regression, R^2 is simply the square of the correlation coefficient between the observed values y_i and the predicted values $\beta_0 + \beta_1 x_i$.

Example 11. In the overfitting example (8), the values of R^2 are:

degree	R^2
1	0.3968
2	0.9455
9	1.0000

Notice the goodness of fit measure increases as n increases. The fit is better, but the model also becomes more complex, since it takes more coefficients to describe higher order polynomials.

MIT OpenCourseWare
<https://ocw.mit.edu>

18.05 Introduction to Probability and Statistics
Spring 2014

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.