# Single Maths B notes

Clare Wallace

2025-12-19

# Table of contents

# Welcome to Single Maths B

**Welcome** to Single Maths B! These lecture notes contain all the mathematical content you'll need to know to succeed in Maths this year.

If you have questions about any of the content here, try one of the following:

- ask a friend!

- ask me! I like to answer emails, and I am often in my office (MCS3060): you can (and should) pop by to see if I'm around. You can do this during my official office hours (Mondays, 1-2pm) for a guaranteed speedy response, but you definitely shouldn't wait until then, especially if it's a short or quick question.

- Google it, or try a textbook. There are some good ones on the reading list (see below).

These notes have been developed over the years by several members of the Maths department, including (most recently) Martin Kerin, Nabil Iqbal, and Steve Abel.

> **Warning**
>
> There could still be typos. If you find one, let me know about it and you can have a free bag of Skittles.

## How to use these notes

The notes contain all the mathematical content for the course. In lectures, we will start at the beginning and work our way through the whole document, until we reach the end (hopefully, this will happen exactly at the end of term).

Throughout the notes, there are boxes like this one:

> **Try it out**
>
> You might be able to do some of the questions on the problem sheet already.

These contain examples you can work through to check your understanding. Wherever possible, I've also worked examples into the text, but there are some places where I want to give you an extra example. These come in purple boxes.

# 1 Probability

## 1.1 Introduction to Probability

### 1.1.1 What is probability?

Probability is how we quantify uncertainty; it is the *extent to which an event is likely to occur*. We use it to study events whose outcomes we do not (yet) know, whether this is because they have not happened yet, or because we have not yet observed them.

We quantify this uncertainty by assigning each event a number between 0 and 1. The higher the probability of an event, the more likely it is to occur.

Historically, the early theory of probability was developed in the context of gambling. In the seventeenth century, Blaise Pascal, Pierre de Fermat, and the Chevalier de Méré were interested in questions like "If I roll a six-sided die four times, how likely am I to get at least one six?" and "if I roll a *pair* of dice twenty-four times, how likely am I to get at least one pair of sixes?" Many of the examples we'll see in this course still use situations like rolling dice, drawing cards, or sticking your hand into a bag filled with differently-coloured tokens.

Nowadays, probability theory helps us to understand how the world around us works, such as in the study of genetics and quantum mechanics; to model complex systems, such as population growth and financial markets, and to analyse data, via the theory of *statistics*.

We'll see a bit of statistical theory at the end of this chapter, but will mostly stay on the probabilistic side of that line.

### 1.1.2 Events

> **Definition**
>
> We use probability theory to describe scenarios in which we don't know what the outcome will be. We call these scenarios *experiments* or *trials*.
> The set of all possible outcomes of an experiment is its *sample space*, $S$. Subsets of $S$ are called *events*, and may contain several different outcomes.

> **Examples**
>
> In the experiment in which we roll a single six-sided die, we have:
>
> - The sample space is $S = \{1, 2, 3, 4, 5, 6\}$
>
> - An example of a possible outcome is 5 (or "we roll a five")
>
> - An example of an event is $A = \{2, 4, 6\}$ (or "we roll an even number").

Because events are subsets of the sample space, we can treat them as sets.

### 1.1.2.1 Set operations

There are three basic operations we can use to combine and manipulate sets.

> **Definition**
>
> If $A$ and $B$ are events, then
>
> - The event *not A*, which we write $A^c$ (the $c$ is for *complement*), is the set of all outcomes in $S$ which are not in $A$.
>
> - The event *A or B*, which we write $A \cup B$ and call the *union* of $A$ and $B$, is the set of all outcomes which are in at least one of $A$ and $B$.
>
> - The event *A and B*, which we write $A \cap B$ and call the *intersection* of $A$ and $B$, is the set of all outcomes which are in both $A$ and $B$.



(a) (a)



(b) (b)



(c) (c)



(d) (d)

Figure 1.1: Pictures illustrating: (a) $A$; (b) $A \cup B$; (c) $A \cap B$; and (d) $A$ $B$.

### 1.1.2.2 Working with events

When we want to consider all the outcomes in an event $A$ which are *not* in $B$, we write $A \cap B^c = A \backslash B$.

We say that two events are *disjoint* (or incompatible, or mutually exclusive) if they cannot occur at the same time; in other words, if $A$ and $B$ are disjoint, then $A \cap B$ contains no outcomes.

We write $A \cap B = \emptyset$, and we call $\emptyset$ the empty set.

If every outcome in an event $A$ is also in an event $B$, we say that $A$ is a *subset* of $B$, and we write $A \subseteq B$.

---

**Examples**

For example, since all Single Maths students are fans of probability,

$$\{\text{Single Maths students}\} \subseteq \{\text{Fans of probability}\}.$$

We can depict this in a diagram: see Figure 1.2 below.

---



Figure 1.2: Notice that the circle of "probability fans" takes up quite a lot of the sample space.

The following set of basic rules will be helpful when working with events.

**Commutativity:**
$$A \cup B = B \cup A, \quad A \cap B = B \cap A$$

**Associativity:**
$$(A \cup B) \cup C = A \cup (B \cup C), \quad (A \cap B) \cap C = A \cap (B \cap C)$$

**Distributivity:**
$$(A \cap B) \cup C = (A \cup C) \cap (B \cup C), \quad (A \cup B) \cap C = (A \cap C) \cup (B \cap C)$$

**De Morgan's laws:**
$$(A \cup B)^c = A^c \cap B^c, \quad (A \cap B)^c = A^c \cup B^c$$

**Example**

For example, if $A = \{$Dinner is on time$\}$ and $B = \{$Dinner is delicious$\}$, then

$$(A \cap B)^c = \{\text{Dinner is either late or disappointing}\},$$

and

$$(A \cup B)^c = \{\text{Dinner is } \textit{both} \text{ late } \textit{and} \text{ disappointing}\}.$$

### 1.1.3 Axioms of Probability

Once we have decided what our experiment (and hence our sample space) should be, we assign a probability to each event $A \subseteq S$. This probability is a number, which we write $\mathbb{P}(A)$.

**Remember** that $A$ is an event, which is a set, and that $\mathbb{P}(A)$ is a probability, which is a number. It makes sense to take the union of sets, or to add numbers together - but not the other way around!

We need a system of rules (the *axioms*) for how the probabilities are assigned, to make sure everything stays consistent. There are lots of such systems, but we will use Kolmogorov's axioms, from 1933. There's no particular reason to choose one system over another, but these are a popular choice.

**Definition**

The axioms are:

1. The probability of any event is a real number in the interval $[0, 1]$: $0 \leq \mathbb{P}(A) \leq 1$.

2. The probability that *something* in $S$ happens is 1: $\mathbb{P}(S) = 1$.

3. If $A$ and $B$ are disjoint events, then $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B)$.

We can use set operations to see some immediate consequences of the axioms:

- Since $A$ and $A^c$ are disjoint, we have $\mathbb{P}(A^c) = \mathbb{P}(S) - \mathbb{P}(A) = 1 - \mathbb{P}(A)$.

- Impossible events have probability zero: $\mathbb{P}(\emptyset) = 0$.

- For (not necessarily disjoint) events $A$ and $B$, we have $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)$.

- If $A \subseteq B$, then $\mathbb{P}(A) \leq \mathbb{P}(B)$.

**Suggested exercises:** Q1 – Q10.

## 1.2 Counting principles

**Key idea**

When our experiment has $m$ outcomes, each of which is *equally likely*, then each outcome $s$ in the sample space $S$ has probability

$$\mathbb{P}(\{s\}) = \frac{1}{m}$$

and each event $A \subseteq S$ has probability

$$\mathbb{P}(A) = \frac{|A|}{m} = \frac{\text{number of ways A can occur}}{\text{total no. of outcomes}}.$$

In this section, we look at some different ways to count the number of outcomes in an event, when the events are more complex than, say, a roll of a die.

### 1.2.1 The multiplication principle

If our experiment can be broken down into $r$ smaller experiments, in which

- the first experiment has $m_1$ equally likely outcomes
- the second experiment has $m_2$ equally likely outcomes
- ...
- the $r$th experiment has $m_r$ equally likely outcomes,

then there are

$$m_1 \times m_2 \times \cdots \times m_r = \prod_{j=1}^{r} m_j$$

possible, equally likely, outcomes for the whole experiment.

> **Try it out**
>
> - If there are four different routes from Newcastle to Durham, and three different routes from Durham to York, how many different routes are there from Newcastle to York?
>
> - If I toss six coins (1p, 2p, 5p, 10p, 15p, and 20p), how many different ways are there to get one 'heads' and five 'tails'?

In general, sampling $r$ times with replacement from $m$ options gives $m^r$ different possiblities.

### 1.2.2 Permutations

When we select $r$ items from a group of size $n$, *in order* and *without replacement*, we call the result a *permutation of size r from n*.

> **Key idea**
>
> The number of permutations of size $r$ from $n$ is
>
> $$n \times (n-1) \times \cdots \times (n-r+1) = \frac{n!}{(n-r)!}.$$

A special case is when we want to arrange the whole list. Then, there are

$$r \times (r-1) \times \cdots \times 1 = \frac{r!}{0!} = r!$$

different permutations.

**Try it out**

- How many different ways are there to arrange six books on a shelf?

- In a society with twenty members, which must choose one president and one secretary, how many different ways can these roles be filled?

- If six (six-sided) dice are rolled, what is the probability that each of the numbers 1-6 appears exactly once?

### 1.2.3 Combinations

When we select $r$ items from a group of size $n$, *without replacement*, but not in any particular order, then we have a *combination of size $r$ from $n$*.

**Key idea**

There are
$$\binom{n}{r} = \frac{(n!)}{(n-r)!r!}$$
different ways to choose a combination of size $r$ from $n$ objects.

Two useful ways of thinking about combinations:

- You might notice that $\binom{n}{r} = \binom{n}{n-r}$. This is because we can also look at the combination of items we *don't* pick. It's much easier (psychologically, at least) to list the different ways to leave 3 cards in the deck than it is to list the different ways to draw 49 cards!

- There is a relationship between combinations and permutations:

$$\text{the number of combinations} = \frac{1}{r!} \times \text{ the number of permutations.}$$

  This is because each combination counted when the order *doesn't* matter comes up $r!$ different times when the order *does* matter.

**Try it out**

- How many different ways are there to form a subcommittee of eight people, from a group of twenty?

- If I have $n$ points on the circumference of a circle, how many different triangles can I form with vertices among these points?

**Remember:** If we're allowed repeated values, the only tool we need is the multiplication principle.

If there can be no repeats (sampling without replacement), then we use permutations if the objects are all distinct, and combinations if they are not. Usually if we're dealt a hand of cards, or draw a bunch of things out of a bag, then they're indistinguishable. But if we're rolling several dice, or assigning objects to people, then we can (hopefully) tell the dice or people apart.

You might find the flowchart in Figure 1.3 helpful.



Figure 1.3: A decision-making flowchart for permutations, combinations, and the multiplication principle.

### 1.2.4 Multinomial coefficients

When we want to separate a group of size $n$ into $k \geq 2$ groups of possibly different sizes, we use *multinomial coefficients*.

> **Key idea**
>
> If the group sizes are $n_1, n_2, \ldots, n_k$, with $n_1 + n_2 + \cdots + n_k = n$, then the number of different ways to arrange the groups is given by the multinomial coefficient
>
> $$\binom{n}{n_1, n_2, \ldots, n_k} = \frac{n!}{n_1! n_2! \ldots n_k!}.$$

To see how this works, think about choosing the groups in order. There are $\binom{n}{n_1}$ ways to choose the first group; then, there are $\binom{n-n_1}{n_2}$ ways to choose the second group from the remaining objects. Continuing like this until all the groups are selected, by the multiplication principle there are

$$\binom{n}{n_1} \times \binom{n-n_1}{n_2} \times \binom{n-n_1-n_2}{n_3} \times \cdots \times \binom{n_{k-1}+n_k}{n_{k-1}} \times \binom{n_k}{n_k}$$

ways to choose all the groups. Writing each binomial coefficient in terms of factorials, and doing (lots of nice) cancelling, we end up with our expression for the multinomial coefficient.

As it turns out, the multinomial coefficient $\binom{n}{n_1, n_2, \ldots, n_k}$ is also the number of different (i.e. distinguishable) permutations of $n$ objects of which $n_1$ are identical and of type 1, $n_2$ are identical and of type 2, ..., $n_k$ are identical and of type $k$ (where $n = n_1 + n_2 + \cdots + n_k$).

> **Example**
>
> The number of different ways to distribute 10 toys among 3 children, ensuring that the youngest gets exactly one more than its older siblings, is
>
> $$\binom{10}{4, 3, 3} = \frac{10!}{4!3!3!} = 4,200.$$

> **Try it out**
>
> - In how many different (i.e. distinguishable) ways can you arrange the letters in STATISTICS?
>
> - If you arrange the letters S,S,S,T,T,T,I,I,A,C in a random order, what is the probability that they spell 'Statistics'?

**Suggested exercises:** Q11 – Q17.

## 1.3 Conditional Probability and Bayes' Theorem

Sometimes, knowing whether or not one event has occurred can change the probability of another event. For example, if we know that the score on a die was even, there is a one in three chance that we rolled a two (rather than one in six). Gaining the knowledge that our score is even affects how likely it is that we got each possible score.

> **Definition**
>
> We write $\mathbb{P}(A \mid B)$ for the *conditional probability of A, given B*; it is defined by
>
> $$\mathbb{P}(A \mid B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}.$$

> **Key idea**
>
> We can rearrange the definition of conditional probability to get
>
> $$\mathbb{P}(A \cap B) = \mathbb{P}(A \mid B)\,\mathbb{P}(B) = \mathbb{P}(B \mid A)\,\mathbb{P}(A),$$
>
> which leads to **Bayes' theorem**:
>
> $$\mathbb{P}(A \mid B) = \frac{\mathbb{P}(B \mid A)\mathbb{P}(A)}{\mathbb{P}(B)}.$$

Writing conditional probabilities in this way allows us to "invert" them; quite often, one of $\mathbb{P}(A \mid B)$ and $\mathbb{P}(B \mid A)$ is easier to spot than the other.

## 1.4 Independence

> **Definition**
>
> We say that two events are *independent* if the occurrence of one has no bearing on the occurrence of the other, that is,
> $$\mathbb{P}(A \mid B) = \mathbb{P}(A).$$

> **Examples**
>
> - The scores obtained from rolling two separate dice are independent.
>
> - Height and shoe size of people are usually not independent.
>
> - Lecture attendance and exam grades are not independent!

When events $A$ and $B$ are independent, we have

$$\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B).$$

## 1.5 Partitions

Suppose we can separate our sample space into $n$ *mutually disjoint events* $E_1, E_2, \ldots, E_n$: we know that exactly one of these events must happen. We call the collection $\{E_1, E_2, \ldots, E_n\}$ a *partition*, and we can use it to break down the probabilities of different events $A \subseteq S$.

First, we can write

$$A = (A \cap E_1) \cup (A \cap E_2) \cup \cdots \cup (A \cap E_n),$$

so that

$$\mathbb{P}(A) = \mathbb{P}(A \cap E_1) + \mathbb{P}(A \cap E_2) + \cdots + \mathbb{P}(A \cap E_n).$$

We can also introduce conditional probability, to get **the partition theorem**:

$$\mathbb{P}(A) = \mathbb{P}(A \mid E_1)\,\mathbb{P}(E_1) + \mathbb{P}(A \mid E_2)\,\mathbb{P}(E_2) + \cdots + \mathbb{P}(A \mid E_n)\,\mathbb{P}(E_n).$$

The partition theorem is useful whenever we can break an event down into cases, each of which is straightforward.

> **Try it out**
>
> One of the most well-known (especially recently!) examples of the partition theorem is in testing for diseases.
> Suppose that a disease affects one in 10,000 people. We have a test for this disease which correctly identifies 90% of people who *do* have the disease (so gives false negatives to 10% of people with the disease), and gives false positives to 1% of people who *do not* have the disease.
> If a randomly chosen person is tested, what is the probability that their test result is positive?
> Given that the test result is positive, what is the probability that they have the disease?

**Suggested exercises:** Q18 – Q26.

## 1.6 Random variables

**Definition**

A *random variable X* is a function $X : S \to \mathbb{R}$ which assigns a numerical value to each possible outcome of an experiment (i.e. to each element of the sample space $S$), such that the probability $\mathbb{P}(X \leq b)$ is well defined for all $b \in \mathbb{R}$ (i.e. the event $\{X \leq b\} = \mathbb{P}(\{s \in S \mid X(s) \leq b\}$ can always be assigned a probability).
We say that a random variable is *discrete* if we can list its possible values, or *continuous* if it can take any value in a range.

We won't ever need to worry about the "well-defined" part of the definition in this module, but, strictly speaking, there do exist complicated real-valued functions on certain sample spaces which are not random variables.

**Example**

If the *experiment* is "toss four coins", then some of the elements of the sample space are HHHH, HHHT, HHTH, HHTT,... . One random variable we can define is

$$X = \text{ Number of heads.}$$

Then if our *outcome* is HHTT, we have $X(\text{HHTT}) = 2$.

### 1.6.1 Discrete random variables

To describe a discrete random variable $X : S \to \mathbb{R}$, we can use its *probability distribution*, which is sometimes called a probability mass function.

**Key idea**

The probability distribution is often displayed in a table, which shows the different values $X$ can take, along with the associated probabilities:

| values | $x_1$ | $x_2$ | ... | $x_n$ |
|---|---|---|---|---|
| probabilities | $\mathbb{P}(X = x_1)$ | $\mathbb{P}(X = x_2)$ | ... | $\mathbb{P}(X = x_n)$ |

Recall here that the event $\{X = x\}$ is given by $\{X = x\} = \{s \in S \mid X(s) = x\} \subseteq S$. In a probability distribution, the probabilities must be *non-negative* and must *sum to 1*. To find the probability that $X$ takes values in an interval $[a, b]$, we have

$$\mathbb{P}(a \leq X \leq b) = \sum_{a \leq x_i \leq b} \mathbb{P}(X = x_i).$$

### 1.6.1.1 Joint and marginal distributions

> **Definition**
>
> When we have two (or more) discrete random variables, $X$ and $Y$ (and $Z$ and...), the **joint probability distribution** is the table of probabilities $\mathbb{P}(X = x, Y = y)$ of every possible combination of values $x$ for $X$ and $y$ for $Y$:
>
> | | $x_1$ | ... | $x_n$ |
> |---|---|---|---|
> | $y_1$ | $\mathbb{P}(X = x_1, Y = y_1)$ | ... | $\mathbb{P}(X = x_n, Y = y_1)$ |
> | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ |
> | $y_m$ | $\mathbb{P}(X = x_1, Y = y_m)$ | ... | $\mathbb{P}(X = x_n, Y = y_m)$ |

Recall here that the event $\{X = x, Y = y\}$ is given by $\{X = x, Y = y\} = \{X = x\} \cap \{Y = y\} \subseteq S$. Moreover, as in the case of the probability distribution of a single random variable, the probabilities in a joint probability distribution must be *non-negative* and must *sum to 1*.

We can find the **marginal probability distributions** of $X$ and $Y$ from the joint distribution, by summing across the rows or columns:

$$\mathbb{P}(X = x_k) = \sum_j \mathbb{P}(X = x_k, Y = y_j),$$

$$\mathbb{P}(Y = y_j) = \sum_k \mathbb{P}(X = x_k, Y = y_j).$$

Two discrete random variables $X$ and $Y$ are said to be **independent** if

$$\mathbb{P}(X = x_k, Y = y_j) = \mathbb{P}(X = x_k)\,\mathbb{P}(Y = y_j)$$

for all possible pairs $(x_k, y_j)$ of values of $X$ and $Y$.

> **Try it out**
>
> Let $X$ be the random variable which takes value 3 when a fair coin lands heads up, and takes value 0 otherwise. Let $Y$ be the value shown after rolling a fair die. Write down the distributions of $X$ and $Y$, and the joint distribution of $X$ and $Y$. You may assume that $X$ and $Y$ are independent. Use your table to find the probability that $X > Y$.

> **Example**
>
> Let $S = \{(a, b) \mid a, b \in \{1, \dots, 6\}\}$ be the sample space when a pair of fair dice is tossed. Let $X : S \to \mathbb{R}$ and $Y : S \to \mathbb{R}$ be the (discrete) random variables defined by
>
> $$X(a, b) = a + b \quad \text{and} \quad Y(a, b) = \max\{a, b\}$$
>
> respectively. Then the joint distribution of $X$ and $Y$ is
>
> | | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
> |---|---|---|---|---|---|---|---|---|---|---|---|
> | 1 | $\frac{1}{36}$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
> | 2 | 0 | $\frac{2}{36}$ | $\frac{1}{36}$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
> | 3 | 0 | 0 | $\frac{2}{36}$ | $\frac{2}{36}$ | $\frac{1}{36}$ | 0 | 0 | 0 | 0 | 0 | 0 |
> | 4 | 0 | 0 | 0 | $\frac{2}{36}$ | $\frac{2}{36}$ | $\frac{2}{36}$ | $\frac{1}{36}$ | 0 | 0 | 0 | 0 |

|   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 5 | 0 | 0 | 0 | 0 | $\frac{2}{36}$ | $\frac{2}{36}$ | $\frac{2}{36}$ | $\frac{2}{36}$ | $\frac{1}{36}$ | 0 | 0 |
| 6 | 0 | 0 | 0 | 0 | 0 | $\frac{2}{36}$ | $\frac{2}{36}$ | $\frac{2}{36}$ | $\frac{2}{36}$ | $\frac{2}{36}$ | $\frac{1}{36}$ |

For example, the event that both $X = 5$ and $Y = 3$ occurs only for the outcomes $(2,3)$ and $(3,2)$, yielding a probability $\mathbb{P}(X = 5, Y = 3) = \frac{2}{36}$. The (marginal) probability that $X = 5$ is

$$\mathbb{P}(X = 5) = \sum_{k=1}^{6} \mathbb{P}(X = 5, Y = k)$$
$$= \mathbb{P}(X = 5, Y = 3) + \mathbb{P}(X = 5, Y = 4)$$
$$= \frac{2}{36} + \frac{2}{36} = \frac{4}{36}$$

as expected, since $X = 5$ occurs only for the outcomes $(1,4), (2,3), (3,2)$ and $(4,1)$. Similarly, the (marginal) probability that $Y = 3$ is

$$\mathbb{P}(Y = 3) = \sum_{m=2}^{12} \mathbb{P}(X = m, Y = 3)$$
$$= \mathbb{P}(X = 4, Y = 3) + \mathbb{P}(X = 5, Y = 3) + \mathbb{P}(X = 6, Y = 3)$$
$$= \frac{2}{36} + \frac{2}{36} + \frac{1}{36} = \frac{5}{36}$$

since $Y = 3$ occurs only for the outcomes $(1,3), (2,3), (3,3), (3,2)$ and $(3,1)$.

Finally, observe that $X$ and $Y$ are not independent random variables since, for example, $\mathbb{P}(X = 2, Y = 3) = 0$, whereas $\mathbb{P}(X = 2) = \frac{1}{36}$ and $\mathbb{P}(Y = 3) = \frac{5}{36}$, so that $\mathbb{P}(X = 2)\,\mathbb{P}(Y = 3) \neq 0$.

### 1.6.2 Continuous random variables

When our random variable is continuous, we cannot describe its probability distribution using a list of probabilities. Instead, we use a **probability density function** (pdf), $f_X(x)$.

> **Key idea**
>
> The density function $f_X(x)$ describes a curve over the possible values taken by the random variable $X$. In a density function, the values must be *non-negative* and *integrate to 1*.

To find the probability that $X$ lies in an interval $[a, b]$, we have

$$\mathbb{P}(a \leq X \leq b) = \int_a^b f_X(x)\,dx.$$

**Remember** that the density $f_X(x)$ is not the same thing as $\mathbb{P}(X = x)$. In fact, for every $x$, we have $\mathbb{P}(X = x) = 0$.

Another way of specifying the distribution of a continuous random variable is through its **cumulative distribution function** (cdf) $F_X : \mathbb{R} \to [0, 1]$, given by

$$F_X(x) = \mathbb{P}(X \leq x) = \int_{-\infty}^{x} f_X(t)\,dt.$$

**Example**

A random variable $X$ is said to have a *uniform distribution* (often denoted $\mathrm{Unif}(a,b)$) on an interval $[a,b]$ if its probability density function $f_X$ satisfies

$$f_X(x) = \begin{cases} \frac{1}{b-a}, & \text{for } x \in [a,b], \\ 0, & \text{otherwise.} \end{cases}$$

If $[c,d] \subset [a,b]$ is another interval (that is, if $a \leq c < d \leq b$), then

$$\mathbb{P}(c \leq X \leq d) = \int_c^d f_X(x)\,dx = \frac{1}{b-a}\int_c^d dx = \frac{d-c}{b-a}.$$

Similarly, if the interval $[c',d']$ has $c' < a \leq d' \leq b$, then, since $f_X(x) = 0$ for all $x < a$,

$$\mathbb{P}(c' \leq X \leq d') = \int_{c'}^{d'} f_X(x)\,dx = \int_a^{d'} f_X(x)\,dx = \frac{1}{b-a}\int_a^{d'} dx = \frac{d'-a}{b-a}.$$

Via similar calculations, we see that the cumulative distribution function $F_X : \mathbb{R} \to [0,1]$ is given by

$$F_X(x) = \begin{cases} 0, & \text{for } x < a, \\ \frac{x-a}{b-a}, & \text{for } a \leq x \leq b, \\ 1, & \text{for } b < x. \end{cases}$$

**Try it out**

Let $X$ be a continuous random variable with probability density function:

$$f_X(x) = \begin{cases} \beta e^{-\beta x}, & \text{for } x > 0, \\ 0, & \text{for } x \leq 0. \end{cases}$$

Check that $f_X(x)$ is a valid probability density function when $\beta > 0$. Find the cumulative distribution function of $X$ and, hence, find $\mathbb{P}(X > 3)$.

### 1.6.2.1 Joint and marginal distributions

**Definition**

The joint probability distribution of two (or more) continuous random variables $X$ and $Y$ (and $Z$ and...) can be described using their **joint probability density function $f_{X,Y}(x,y)$**. This is a function of two variables describing how the pair of random variables $X$ and $Y$ are "spread out".

As it is a density, the function $f_{X,Y}$ is non-negative and must integrate to 1. The probability that $X$ and $Y$ take values in a region $A$ of the $xy$-plane is given by the double integral (to be discussed in **?@sec-integration**

$$\mathbb{P}((X,Y) \in A) = \iint_A f_{X,Y}(x,y)\,dxdy.$$

We can find the **marginal probability distributions** of $X$ and $Y$ from the joint distribution, by

integrating out one of the variables:

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x,y)\,dy$$

$$f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x,y)\,dx.$$

**Suggested exercises:** Q27 – Q32.

## 1.7 Expectation and Variance

While the probability distribution or probability density function tells us everything about the distribution of a random variable, this can often be too much information. Quantities which instead summarise the distribution can be useful to convey information about our random variable without trying to describe it in its entirety.

Summaries of a distribution include the expectation, the variance, the skewness and the kurtosis. In this course, we're only interested in the expectation, which tells us about the *location* of the distribution, and the variance, which tells us about its *spread*. The skewness tells us about the *symmetry* of the distribution about its expectation, while the kurtosis tells us about the likelihood of the random variable taking values far away from the mean.

### 1.7.1 Expectation

---
**Definition**

The **expectation** of a random variable $X$ is given by

$$\mathbb{E}[X] = \begin{cases} \sum_x x\,\mathbb{P}(X=x)\,, & \text{if } X \text{ is discrete,} \\ \int_{-\infty}^{\infty} x f_X(x)\,dx\,, & \text{if } X \text{ is continuous.} \end{cases}$$

The expectation is sometimes called the *mean* or the *average* of the random variable $X$.

---

#### 1.7.1.1 Properties of Expectation

**Linearity:** If $X$ is a random variable and $a$ and $b$ are (real) constants, then

$$\mathbb{E}[aX+b] = a\,\mathbb{E}[X] + b.$$

**Additivity:** If $X_1, X_2, \ldots, X_n$ are random variables, then

$$\mathbb{E}[X_1 + X_2 + \cdots + X_n] = \mathbb{E}[X_1] + \mathbb{E}[X_2] + \cdots + \mathbb{E}[X_n].$$

**Positivity:** If $X$ is a positive random variable (that is, if $\mathbb{P}(X \geq 0) = 1$), then $\mathbb{E}[X] \geq 0$.

**Independence:** If $X$ and $Y$ are independent random variables, then

$$\mathbb{E}[XY] = \mathbb{E}[X]\,\mathbb{E}[Y].$$

**Expectation of a function:** If $X$ is a random variable and $r$ is a (nice[1]) function, then $r(X) = r \circ X$ is a random variable with expectation

$$\mathbb{E}[r(X)] = \begin{cases} \sum_x r(x)\, \mathbb{P}(X = x), & \text{if } X \text{ is discrete,} \\ \int_{-\infty}^{\infty} r(x) f_X(x)\, dx, & \text{if } X \text{ is continuous.} \end{cases}$$

### 1.7.2 Variance

> **Definition**
>
> For a random variable $X$ with expectation $\mathbb{E}[X] = \mu$, the **variance** of $X$ is given by
>
> $$\text{Var}(X) = \mathbb{E}[(X - \mu)^2].$$

By expanding out the brackets and using the linearity of the expectation, we can rewrite the variance as

$$\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2.$$

The variance is always positive, because it is the expectation of a positive random variable. The **standard deviation** is the square root of the variance:

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

### 1.7.2.1 Properties of Variance

**Affine transformations:** If $X$ is a random variable and $a$ and $b$ are (real) constants, then

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

**Independence:** If $X$ and $Y$ are independent random variables, then

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y).$$

> **Try it out**
>
> - Let $X$ be a continuous random variable with probability density function:
>
>   $$f_X(x) = \begin{cases} \beta e^{-\beta x} & \text{for } x > 0, \\ 0 & \text{for } x \leq 0. \end{cases}$$
>
>   What are the expectation and variance of $X$?
>
> - Let $Y$ be a random variable with the following probability distribution:

---

[1]Actually this is generally only true *locally*, i.e. the function $f(x, y)$ that we build might have some problems if we try to define it in all space. In fact this works only when the spaces we are considering are "simple"; if they have holes in them etc. then we cannot globally define $f$. This is thus a connection between "topology" (i.e. global properties of spaces) and calculus. It is also a surprisingly important subject to physicists: google "de Rham cohomology" to find out more.

| $y$ | 1 | 2 | 3 |
|---|---|---|---|
| $\mathbb{P}(Y = y)$ | $\frac{1}{6}$ | $\frac{2}{6}$ | $\frac{3}{6}$ |

Find $\mathbb{E}[X]$, $\mathrm{Var}(X)$, and $\mathbb{E}\left[\frac{1}{X}\right]$.

**Suggested exercises:** Revisit Q30; Q33 – Q37.

## 1.8 The Binomial Distribution

> **Definition**
>
> If $X$ is the number of successes (i.e. 0 or 1) from a single experiment which succeeds with probability $p$ and fails with probability $1 - p$, then the random variable $X$ has probability distribution
>
> | $x$ | $0$ | $1$ |
> |---|---|---|
> | $\mathbb{P}(X=x)$ | $1-p$ | $p$ |
>
> In such a case, we say that $X$ has a *Bernoulli distribution with parameter $p$* and write $X \sim \mathrm{Bern}(p)$.

The expectation and variance of $X \sim \mathrm{Bern}(p)$ are:

$$\mathbb{E}[X] = p$$
$$\mathrm{Var}(X) = p(1-p).$$

Suppose we have $n$ Bernoulli-style trials, which succeed or fail *independently* of each other, and such that all trials have the same probability $p$ of succeeding. We count the *total number of successes* across all the trials.

> **Definition**
>
> If $Y$ is the total number of successes from $n$ independent Bernoulli trials (each with parameter $p$), we say that $Y$ has a binomial distribution with parameters $n$ and $p$, and we write $Y \sim \mathrm{Bin}(n,p)$.

If $0 \le k \le n$, we have

$$\mathbb{P}(Y = k) = \binom{n}{k} p^k (1-p)^{n-k}.$$

This is because each configuration of $k$ successes and $n - k$ successes has probability $p^k(1-p)^{n-k}$, by the multiplication principle; and there are $\binom{n}{k}$ different ways of arranging the $k$ successes and $n - k$ failures among the trials.

> **Try it out**
>
> Check that the probabilities in the binomial distribution are all non-negative and sum to 1.

The expectation and variance of $Y \sim \text{Bin}(n, p)$ are:

$$\mathbb{E}[Y] = np$$
$$\text{Var}(Y) = np(1 - p).$$

---

**Examples**

- If I toss six coins, the total number of heads has a $\text{Bin}(6, \frac{1}{2})$ distribution.

- If each SMB student decides to skip a lecture with probability 0.2, then the number of students who turn up has a $\text{Bin}(217, 0.8)$ distribution (assuming you all decide independently of each other!). In particular, the expected number of students at each lecture is $217 \times 0.8 \approx 174$.

---

## 1.9 The Poisson Distribution

While the binomial distribution is about counting successes in a *fixed* number of trials, the Poisson distribution lets us count how many times something happens without a fixed upper limit. This is useful in a lot of real-world contexts, for example:

- the number of people who visit a website

- the number of yeast cells in a sample (such as in experiments by Gossett at Guinness in the 1920s)

- the number of particles emitted from a radioactive sample.

---

**Definition**

The Poisson distribution is used to model scenarios in which events happen randomly, independently, and at a *constant rate* $r$. If $X$ is the total number of these events that happen in a time period of length $s$, then $X$ has a Poisson distribution with parameter $\lambda = rs$, and we write $X \sim \text{Po}(\lambda)$.
If $k \in \mathbb{N}$, we have

$$\mathbb{P}(X = k) = e^{-\lambda} \frac{\lambda^k}{k!}.$$

---

**Try it out**

Check that the probabilities in the Poisson distribution are all non-negative and sum to 1.

---

The expectation and variance of $X$ are

$$\mathbb{E}[X] = \text{Var}(X) = \lambda.$$

### 1.9.1 Using the Poisson distribution to approximate the binomial distribution

Instead of thinking about our time period $[0, s]$ as one long interval, we can split it up into $n$ smaller ones (each one will have length $\frac{s}{n}$).

Suppose we count the number of sub-intervals in which events occur. If the sub-intervals are small enough, it is very unlikely that there will be multiple events in any of them, and the probability that there is one event will be $p \approx \frac{rs}{n} = \frac{\lambda}{n}$.

We can view the sub-intervals as $n$ independent trials, and the total number of successes becomes binomially distributed.

This is a good approximation because the probabilities $\mathbb{P}(X = k)$ in the binomial distribution $\text{Bin}(n, \frac{\lambda}{n})$ and the Poisson distribution $\text{Po}(\lambda)$ are similar as long as $n$ is big enough. That is, for large $n$ we have

$$\binom{n}{k}\left(\frac{\lambda}{n}\right)^k\left(1 - \frac{\lambda}{n}\right)^{n-k} = \frac{n(n-1)\dots(n-k+1)}{k!}\frac{\lambda^k}{n^k}\left(1 - \frac{\lambda}{n}\right)^{n-k}$$

$$= \frac{n(n-1)\dots(n-k+1)}{n^k} \times \left(1 - \frac{\lambda}{n}\right)^{n-k} \times \frac{\lambda^k}{k!}$$

$$\approx 1 \times e^{-\lambda} \times \frac{\lambda^k}{k!},$$

This approximation is good if $n \geq 20$ and $p \leq 0.05$, and excellent if $n \geq 100$ and $np \leq 10$.

**Suggested exercises:** Revisit Q38–Q41.

## 1.10 The Normal Distribution

Unlike the binomial and Poisson distributions, the normal (or Gaussian) distribution is continuous. It is one of the most used (and most useful) distributions. A random variable whose "large-scale" randomness comes from many small-scale contributions is usually normally distributed: for example, people's heights are determined by many different genetic and environmental factors. All of these different factors have tiny impacts on your final height; overall, the distribution of the height of a random person is roughly normal.

### 1.10.1 The standard normal distribution

The first version of the normal distribution we will meet is the standard normal.

> **Definition**
>
> We say that a continuous random variable $Z$ has a *standard normal distribution*, and we write $Z \sim \mathcal{N}(0, 1)$, if its probability density function is given by
>
> $$f_Z(x) = \frac{1}{\sqrt{2\pi}}e^{-\frac{x^2}{2}}$$
>
> for all $x \in \mathbb{R}$.

**Properties of the standard normal distribution**

- The probability density function $f_Z$ of a random variable with standard normal distribution is symmetric about 0. Then
$$\mathbb{P}(Z \leq z) = \mathbb{P}(Z \geq -z) = \mathbb{P}(-Z \leq z),$$
which implies, in particular, that the random variable $-Z$ has the same (normal) distribution as $Z$.

- This symmetry also means that $xf_Z(x)$ is an odd function; so the expectation of $Z$ is zero.

- The variance of $Z$ is
$$\text{Var}(Z) = \mathbb{E}[Z^2] - 0$$
$$= \int_{-\infty}^{\infty} x^2 f_Z(x) \, dx$$
$$= \frac{1}{\sqrt{2\pi}} \int_{\infty}^{\infty} x^2 e^{-\frac{x^2}{2}} \, dx = 1.$$

(You can find this via integration by parts.)

**The cumulative distribution function for $Z$**

> **Definition**
>
> The cumulative distribution function for $Z$ is denoted $\Phi(z)$ and is given by
> $$\Phi(z) = \mathbb{P}(Z \leq z) = \int_{-\infty}^{z} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \, dx.$$

There is no neat ("algebraic") expression for $\Phi(z)$: in practice, when we need to evaluate it we use numerical methods to get (usually very good) approximations. These values are traditionally recorded in tables but usually, they're built into computer software and some calculators.

Some useful properties of $\Phi(z)$, which reduce the number of values we need in the tables, are:

- Because $f_Z(x)$ is symmetric, we have
$$\Phi(z) = \mathbb{P}(Z \leq z) = \mathbb{P}(Z \geq -z) = 1 - \Phi(-z).$$

- We have $\Phi(0) = \frac{1}{2}$.
- $\mathbb{P}(a \leq Z \leq b) = \Phi(b) - \Phi(a)$.

**Interpolation:** When the value we need to find isn't in a table we have access to, we can interpolate. If $a < b < c$ and we know $\Phi(a)$ and $\Phi(b)$, we approximate:
$$\Phi(b) \approx \Phi(a) + \frac{b-a}{c-a} \left( \Phi(c) - \Phi(a) \right).$$

For example, most normal tables only go to two decimal places, but $\Phi(0.553)$ will be approximately $(0.553 - 0.55)/(0.56 - 0.55) = 0.3$ of the way between $\Phi(0.55)$ and $\Phi(0.56)$.

### 1.10.2 General normal distributions

> **Definition**
>
> We say that a continuous random variable $X$ has a *normal distribution with parameters $\mu$ and $\sigma^2$*, and we write $X \sim \mathcal{N}(\mu, \sigma^2)$, if the random variable $Z = \frac{X-\mu}{\sigma}$ has a standard normal distribution.

We can also write this in the other direction: $X \sim \mathcal{N}(\mu, \sigma^2)$ if $X = \sigma Z + \mu$. Since the distribution of $Z$ is symmetric, we use the convention $\sigma > 0$.

**Properties of general normal distributions**

- The expectation of $X$ is
$$\mathbb{E}[X] = \mathbb{E}[\sigma Z + \mu]$$
$$= \mu + \sigma \, \mathbb{E}[Z]$$
$$= \mu + 0 = \mu.$$

- The variance of $X$ is
$$\mathrm{Var}(X) = \mathrm{Var}(\sigma Z + \mu)$$
$$= \sigma^2 \, \mathrm{Var}(Z)$$
$$= \sigma^2.$$

- The probability density function of $X$ is
$$f_X(x) = \frac{1}{\sigma} f_Z \left( \frac{x - \mu}{\sigma} \right) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \left( \frac{x - \mu}{\sigma} \right)^2 \right\}.$$

- The cumulative distribution function of $X$ is given by
$$\mathbb{P}(X \le x) = \mathbb{P}\left( \sigma Z + \mu \le \frac{x - \mu}{\sigma} \right)$$
$$= \mathbb{P}\left( Z \le \frac{x - \mu}{\sigma} \right)$$
$$= \Phi\left( \frac{x - \mu}{\sigma} \right).$$

We can use the table for the standard normal distribution to evaluate the cumulative distribution function of *any* normal distribution, by using this transformation.

---

**Try it out**

1. If $X \sim \mathcal{N}(12, 25)$, what is $\mathbb{P}(X \le 3)$?

2. If $Y \sim \mathcal{N}(1, 4)$, what is $\mathbb{P}(-1 < Y < 2)$?

---

### 1.10.3 Using the normal distribution to approximate the binomial and Poisson distributions

Just as we can use the Poisson distribution to approximate specific probabilities in the binomial distribution, we can use the normal distribution to approximate cumulative probabilities. If $n$ is large and $X \sim \mathrm{Bin}(n, p)$, then approximately we have $X \sim \mathcal{N}(np, np(1 - p))$.

In particular,
$$\mathbb{P}(X \le k) \approx \Phi\left( \frac{k - np}{\sqrt{np(1 - p)}} \right).$$

This is a useful approximation when both $np$ and $np(1 - p)$ are at least 10; as the two values increase, the approximation gets better.

> **Try it out**
>
> A machine produces $n = 1500$ gadgets every day. Each individual gadget is defective with probability $p = 0.02$. Find (approximately) the probability that more than 40 of the items produced in one day are defective.

Similarly, we can use the normal distribution to approximate the cumulative probabilities in the Poisson distribution: if $X \sim \text{Po}(\lambda)$, then approximately we have $X \sim \mathcal{N}(\lambda, \lambda)$ and

$$\mathbb{P}(X \le k) \approx \Phi\left(\frac{k - \lambda}{\sqrt{\lambda}}\right).$$

This is a useful approximation when $\lambda$ is at least 5, and gets better as $\lambda$ increases.

**Suggested exercises:** Q42 – Q45.


## 1.11 The Central Limit Theorem

### 1.11.1 Experimental errors

> **Definition**
>
> When we are measuring a quantity whose "true value" is $\mu$, our measurement takes the form $X = \mu + \varepsilon$, where $\varepsilon$ is the *experimental error*. Before we do the experiment, we can think of both $\varepsilon$ and $X$ as random quantities. Afterwards, $X$ is a fixed and known quantity, and $\mu$ and $\varepsilon$ are fixed but unknown quantities (to us). Our goal is to use $X$ to infer something about $\mu$.

**Assumption:** We will assume that there are no systematic errors or bias in the experiment; in other words, $\mathbb{E}[\varepsilon] = 0$.

If the variance of $\varepsilon$ is $\text{Var}(\varepsilon) = \sigma^2$, then

$$\mathbb{E}[X] = \mu + \mathbb{E}[\varepsilon] = \mu + 0 = \mu$$
$$\text{Var}(X) = 0 + \text{Var}(\varepsilon) = \sigma^2.$$

This means that, on average, the value of our measurement is a good estimate of the value of $\mu$; however, if the variance of $\varepsilon$ is large, our measurement will have quite a high probability of being far from the true value.

To improve our estimate, we can do one of two things:

- try to improve our measurement technique, to reduce the variance

- take more measurements!

### 1.11.2 The sample mean

> **Definition**
>
> When we take $n$ independent random variables $X_1, X_2, \ldots, X_n$ which all have the same distribution, we say that $X_1, X_2, \ldots, X_n$ are *independent and identically distributed (i.i.d.)*.

> **Example**
>
> We might obtain i.i.d. samples by repeating our measurement, or experiment, $n$ times, or by sampling $n$ people from a large population.

> **Definition**
>
> If $X_1, X_2, \ldots, X_n$ are random variables, then the *sample mean* is the average
>
> $$\overline{X} = \frac{1}{n} \sum_{j=1}^{n} X_j.$$

Before we take our measurements, this is also a random variable; afterwards, it is just a number. To distinguish between the two situations, we use $\overline{X}$ for the random variable, and $\overline{x}$ for the number.

It is perhaps worth remarking that the definition of the sample mean is a little ambiguous, as $\overline{X}$ is not defined on the same sample space as $X_1, X_2, \ldots, X_n$. Indeed, if $S$ is the sample space on which each $X_j$ is defined (i.e. $X_j : S \to \mathbb{R}$), then the sample mean is defined on the $n$-fold product sample space $S \times S \times \cdots \times S$ via $\overline{X}(s_1, s_2, \ldots, s_n) = \frac{1}{n} \sum_{j=1}^{n} X_j(s_j)$. That is, we take the list $(s_1, s_2, \ldots, s_n) \in S \times S \times \cdots \times S$ of $n$ outcomes (e.g. of an experiment repeated $n$ times), then evaluate the $j^{\text{th}}$ random variable $X_j$ on the $j^{\text{th}}$ outcome $s_j$ and, finally, compute the average of the values obtained. Let's look at an example.

> **Try it out**
>
> We toss a pair of fair dice eight times. For each toss, the sample space is given by $S = \{(a, b) \mid a, b \in \{1, \ldots, 6\}\}$. Let $X_1, X_2, \ldots, X_8$ be random variables, where $X_j : S \to \mathbb{R}$ is defined as the sum $X_j(a_j, b_j) = a_j + b_j$ of the outcome $(a_j, b_j) \in S$ of the $j^{\text{th}}$ toss. If the outcomes of the eight tosses are
>
> $$(1, 3), (5, 2), (3, 3), (5, 6), (1, 1), (4, 3), (2, 3) \text{ and } (1, 3)$$
>
> , respectively, find the sample and population means.
> **Answer:** The sample mean (i.e. its value after all tosses have been completed) is given by
>
> $$\overline{x} = \frac{1}{8} \left( X_1(1, 3) + X_2(5, 2) + X_3(3, 3) + X_4(5, 6) + X_5(1, 1) + X_6(4, 3) + X_7(2, 3) + X_8(1, 3) \right)$$
> $$= \frac{1}{8} (4 + 7 + 6 + 11 + 2 + 7 + 5 + 4)$$
> $$= \frac{46}{8} = 5.75.$$
>
> On the other hand, the *population mean*/expectation in this case would be $\frac{252}{36} = 7$, since the total of the sums $a + b$ of all 36 possible outcomes $(a, b) \in S$ is 252. (It's not hard to check directly that $\mathbb{E}[X_j] = 7$ for each $j$.) With a much larger number of tosses, we could expect that the sample mean would be close to the population mean/expectation.

**Assumption:** We assume that $X_1, X_2, \ldots, X_n$ are i.i.d. with shared mean $\mu$ and variance $\sigma^2$. Then

$$\mathbb{E}[\overline{X}] = \frac{1}{n} \sum_{j=1}^{n} \mathbb{E}[X_j] = \frac{n}{n}\mu = \mu$$

$$\mathrm{Var}(\overline{X}) = \frac{1}{n^2} \sum_{j=1}^{n} \mathrm{Var}(X_j) = \frac{n}{n^2}\sigma^2 = \frac{\sigma^2}{n}.$$

So the expectation of the sample mean is always $\mu$: we call it an *unbiased estimator for the mean.* On the other hand, the variance is always smaller than $\sigma^2$, and decreases as we increase $n$. By taking a large enough sample size, we can get as small a variance as we want.

If $n$ is large enough, the sample mean will give an accurate estimate for the true mean $\mu$. This result is called the *Law of Large Numbers*, which says that $\overline{X}$ converges[2] to $\mu$ as $n \to \infty$.

### 1.11.3 The Central Limit Theorem

We know that the sample mean will be quite close to the true value $\mu$ on average. The Central Limit Theorem tells us more about the distribution of the error.

> **Key idea:** The Central Limit Theorem
>
> If $X_1, X_2, \ldots, X_n$ are i.i.d. random variables with shared mean $\mu$ and variance $\sigma^2$, then, for large $n$, the sample mean $\overline{X}$ is approximately normally distributed with mean $\mu$ and variance $\frac{\sigma^2}{n}$; that is, $\overline{X}$ is approximated by $\mathcal{N}(\mu, \frac{\sigma^2}{n})$.
> In other words, for large $n$, the random variable
>
> $$Z = \frac{\overline{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$$
>
> is approximately a standard normal distribution.

Here, when we say that the distribution is approximately normal, we mean that

$$\mathbb{P}(a \leq \overline{X} \leq b) \approx \Phi\left(\frac{b - \mu}{\sigma/\sqrt{n}}\right) - \Phi\left(\frac{a - \mu}{\sigma/\sqrt{n}}\right),$$

whatever the values of $a$ and $b$.

> **Try it out**
>
> - If the random variables $X_1, X_2, \ldots, X_{10}$ are independent, and all are uniformly distributed on the interval $[0, 1]$, use the Central Limit Theorem to estimate $\mathbb{P}(X_1 + X_2 + \cdots + X_{10} > 7)$.
>
> - A manufacturing process is designed to produce bolts with a 0.5cm diameter. Once a day, a random sample of 36 bolts is selected and the diameters recorded. If the average of the 36 values is less than 0.49cm or greater than 0.51cm, then the process is shut down for inspection and adjustment. The standard deviation for individual diameters is 0.02cm. Find approximately the probability that the line will be shut down unnecessarily (i.e., if the true process mean really is 0.5cm).

**Suggested exercises:** Q46–Q50.

---

[2]There's quite a lot of probability theory hiding behind this "converges"!

# 2 Partial derivatives

## 2.1 Functions of several variables

Our course now builds on the calculus that you have learnt in Single Maths B. There you worked with functions of just one variable; in this part of the course we will extend the idea of *differentiation and integration to functions of more than one variable.* This is both great fun and of fundamental importance to many applications of mathematics.

### 2.1.1 Examples of functions of several variables

To orient ourselves, let's just remind ourselves that in ordinary life it is really quite common to have functions of more than one variable.

> **Examples**
>
> 1. **Areas and volumes:**
>
>    a. The volume $V$ of a circular cylinder of radius $r$ (cm) and height $h$ (cm) is $V(r,h) = \pi r^2 h$, which is clearly a function of two variables, $r$ and $h$. The cylinder gets larger if we increase $r$ or $h$.
>
>    b. A rectangular box with sides $x, y, x$ has volume $V(x,y,z) = xyz$, a function of 3 variables.
>
> 2. **Heights above a surface:** The surface of the earth is a two-dimensional sphere; we call this $S^2$. The position of a point on the earth can be specified by two coordinates $(\phi, \theta)$, which roughly map to the concepts of latitude and longitude. The height above sea level can be expressed as $h(\theta, \phi)$.
>
> 3. **Atmospheric temperature:** This is again a function of $\theta, \phi$ but also (quite obviously) time dependent $T(t, \theta, \phi)$. You can clearly have things depending on four variables as well, for example the variation of temperature with height $r$ would be $T(t, r, \theta, \phi)$.

### 2.1.2 Graphs of functions

We understand that an equation $y = f(x)$ describes a curve in a plane. Given $x$ we can compute $y$ if we know the function $f$. The equation $y = f(x)$ describes how $(x, y)$ moves as we vary $x$. The relation may be expressed *implicitly* as $g(x, y) = 0$. Let's discuss a few examples:

1. $x^2 + y^2 = 1$ describes a circle of unit radius.

$$y^2 = 1 - x^2,$$

so $y = \pm\sqrt{1 - x^2}$ for $|x| \leq 1$. Thus in the language above, $g(x, y) = x^2 + y^2 - 1$. We can write two equations of the form $y = f(x)$, one for the lower half of the circle, and one for the upper half.

2. Now let's generalize this to higher dimensions. In three dimensional space, the equation

$$x^2 + y^2 + z^2 = 25$$

describes a sphere of radius 5. This can be inverted for $z$ as a function of $(x, y)$. This gives

$$z^2 = 25 - x^2 - y^2$$

so $z = \pm\sqrt{25 - x^2 - y^2}$ for $x^2 + y^2 \leq 25$. Note that for fixed $y$ it describes circles in $x, z$ plane. (How big are the circles if $y = 4$, i.e. if we take a slice 4 units from the origin? With $y = 4$, we have $z = \pm\sqrt{9 - x^2}$, so we get a circle of radius 3.) Now we have two distinct equations of form $z = f(x, y)$. Then every $x, y$ gives us two points lying on the sphere, one for each equation.

3. An equation of the form $z = f(x, y)$ thus describes a surface. We can think of $z$ as the height of the surface above the $x, y$ plane (or the depth below if $z$ is negative). As in the lower-dimensional case, we can often represent this surface *implicitly* as $g(x, y, z) = 0$ for some choice of $g(x, y, z)$: in the case above we have

$$g(x, y, z) = x^2 + y^2 + z^2 - 25$$

The function $f$ may be defined for all $x, y$ or a restricted set as in the sphere. (What does $x^2 + y^2 > 25$ correspond to?)

Thinking about curves and surfaces in higher dimensions, we have two new helpful ideas:

**Definition**

A *level curve* (or *contour line*) is a curve given by taking a horizontal slice of the surface $g(x, y, z) = 0$; that is, by setting $z$ to some fixed value. For example, if we choose $z = c$ it is the set of all $(x, y)$ such that $g(x, y, c) = 0$. Level curves can be viewed as curves in the $xy$-plane by forgetting about the $z$-direction.

A *section curve* is a curve given by instead taking a vertical slice; that is, by freezing one of the other variables. For example, we set $g(c, y, z) = 0$. This gives a curve in either the $yz$- or the $xz$-plane.

**Examples**

Contour lines on a map represent height above sea level - lines of constant $z = h(x, y)$; isobars on a weather map represent lines of constant atmospheric pressure (at sea level) $p(x, y)$.

### 2.1.3 Examples of graphs of functions

Using this Desmos link, explore the properties of the graphs of some of these functions:

1. $z = 5$: for some range of $x, y$. - a flat roof.

2. $z = x^2 + \sin y$: Defined for finite $x$ and $y$. For fixed $y$ we get a parabola, and for fixed $x$ we have a *sin* curve whose height is shifted.

3. $z = \cos(xy)$: Again the function is defined for all $x, y$. $z = const$ for $xy = const$. Contours of constant height would be hyperbolas. Keeping e.g. $y$ constant gives a cosine in $x$ with a period determined by $y$.

4. $z = \frac{\sin \sqrt{x^2+y^2}}{\sqrt{x^2+y^2}}$ the sombrero: Note the *rotational symmetry* about the $z$ axis through the origin. $z$ has the same value for all $x^2 + y^2 = r^2$ with $r$ constant. This is the equation for a circle in the $x, y$ plane; the rotational symmetry is quite evident in the picture.

5. $z = x^3 - 3xy^2$ the monkey saddle: Cubic curve for fixed $y$ and parabolic for fixed $x$.

6. $z = (x^2 + y^2)/a^2$ satellite dish: This is a circular paraboloid. Parallel rays are focussed onto the focus of the dish at $(0, 0, a)$.

**Suggested questions:** Q1-2.

## 2.2 Partial derivatives

We now embark on a program of extending things that we know from single-variable calculus to multiple variables. The first thing we study is the *partial derivative*.

---

**Definition**

Suppose we have a function $f(x, y)$ of two variables. The *partial derivative of $f$ with respect to $x$* is the slope when we move in the $x$ direction but keep $y$ constant:

$$\frac{\partial f(x, y)}{\partial x} = \lim_{h \to 0} \frac{f(x + h, y) - f(x, y)}{h},$$

while the *partial derivative of $f$ with respect to $y$* is the slope of the function when we move in the $y$ direction but keep $x$ constant:

$$\frac{\partial f(x, y)}{\partial y} = \lim_{k \to 0} \frac{f(x, y + k) - f(x, y)}{k}.$$

---

At a point $(a, b)$, the partial derivative gives the slope of the relevant section curve (i.e. either $z = f(x, b)$ or $z = f(a, y)$). It's called "partial" because we are only differentiating with respect to one of the variables (not all of them), and it is denoted with $\partial$, **not** d. The result is simply the same as taking the derivative with respect to $x$ (or $y$) and treating $y$ (or $x$ respectively) as a constant. For any function $f(x, y)$ we get

$$\frac{\partial f}{\partial x} \text{ by differentiating with respect to x and keeping y constant}$$
$$\frac{\partial f}{\partial y} \text{ by differentiating with respect to y and keeping x constant}$$

Similarly for a function $f(x_1, x_2, \dots, x_n) : \mathbb{R}^n \to \mathbb{R}$, the partial derivative $\frac{\partial f}{\partial x_i}$ is the derivative keeping all but $x_i$ constant.

> **Example**
>
> The volume of a can of radius $r$ and height $h$ is $V(r, h) = \pi r^2 h$.
> If we keep the height fixed what is the rate of change of $V$ relative to the radius $r$? We take the partial derivative with respect to $r$:
> $$\frac{\partial V}{\partial r} = 2\pi r h.$$

Sometimes $\frac{\partial f}{\partial x}$ is written $f_x$ and $\frac{\partial f}{\partial y}$ is written $f_y$. Differentiating again we get

$$\frac{\partial}{\partial x}\left(\frac{\partial f}{\partial x}\right) = \frac{\partial^2 f}{\partial x^2} \quad \text{(also written as } f_{xx}); \qquad \frac{\partial}{\partial y}\left(\frac{\partial f}{\partial x}\right) = \frac{\partial^2 f}{\partial y \partial x} \quad \text{(also written as } f_{xy})$$

$$\frac{\partial}{\partial x}\left(\frac{\partial f}{\partial y}\right) = \frac{\partial^2 f}{\partial x \partial y} \quad \text{(also written as } f_{yx}); \qquad \frac{\partial}{\partial y}\left(\frac{\partial f}{\partial y}\right) = \frac{\partial^2 f}{\partial y^2} \quad \text{(also written as } f_{yy}).$$

This is a crucial (and rather simple) concept for the rest of the course, so before going any further let's work out another example.

> **Example**
>
> Consider $f(x, y) = x^2 y + y^3$.
> We have
> $$\frac{\partial f}{\partial x} = 2xy \qquad \frac{\partial f}{\partial y} = x^2 + 3y^2$$
>
> Let's keep going and work out the *second* partial derivatives:
> $$\frac{\partial^2 f}{\partial x\,\partial y} \equiv f_{yx} = 2x \qquad \frac{\partial^2 f}{\partial y\,\partial x} = f_{xy} = 2x$$
>
> It looks like in this example it did not matter in which order I took the partial derivatives! It turns out this is always true, provided $f$ is a sufficiently nice function.

**Clairault's Theorem (or Schwarz's Theorem):**[1] Consider a function $f : \mathbb{R}^n \to \mathbb{R}$; that is, a real-valued function $f(x_1, \dots, x_n)$ depending on $n$ variables. If the second-order partial derivatives exist and are continuous on a small open disc centred at a point $a = (a_1, a_2, \dots, a_n) \in \mathbb{R}^n$, then

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(a) = \frac{\partial^2 f}{\partial x_j \partial x_i}(a),$$

for all $i, j \in \{1, 2, \dots, n\}$.

Here, "small open disc centred at $a$" means the subset $\{x \in \mathbb{R}^n \mid \|x - a\| < r\} \subseteq \mathbb{R}^n$ of all points in $\mathbb{R}^n$ of distance (strictly) less that some fixed $r > 0$ from $a$. The radius $r$ of the disc can be arbitrarily small: the important thing is that *some* choice of radius works, not precisely which one. Moreover, the continuity condition is automatically satisfied if the second-order partial derivatives are themselves differentiable (i.e. if the third-order partial derivatives exist).

---

[1]...but why don't you like matrices?

Clairault's Theorem holds for all familiar functions: polynomials, trigonometric functions, exponential functions, .... If you'd like to see a relatively simple function for which the second-order partial derivatives do not commute, then check out the function found by Peano here.

We sometimes express Clairault's Theorem in words by saying that "partial derivatives *commute*": the word "commutative" means that the order in which you do the two operations (i.e. partial differentiation with respect to $x_i$ and with respect to $x_j$) does not matter.

### More (Somewhat Tedious) Examples

You probably get the idea by now, but here are a few more examples:

> **Try it out**
>
> Calculate all the first and second partial derivatives of:
>
>    a. $f(x, y) = x^3 - 3xy^2$.
>    b. $f(x, y) = \cos y + \sin(xy)$.
>    c. $f(x, y) = x^2 y + yz + z^2 x$.
>
> **Answers:**
>
>    a. We have
> $$f_x = 3x^2 - 3y^2$$
> $$f_y = -6xy$$
> $$f_{xx} = 6x$$
> $$f_{yy} = -6y$$
> $$f_{xy} = f_{yx} = -6y.$$
>
>    b. Let $f(x, y) = \cos y + \sin(xy)$.
>
> Then
> $$f_x = y \cos(xy)$$
> $$f_y = -\sin y + x \cos(xy)$$
> $$f_{xx} = -y^2 \sin(xy)$$
> $$f_{yy} = -\cos y - x^2 \sin(xy)$$
> $$f_{xy} = f_{yx} = \cos(xy) - xy \sin(xy).$$
>
>    c. When $f(x, y) = x^2 y + yz + z^2 x$, we have
> $$f_x = 2xy + z^2$$
> $$f_y = z + x^2$$
> $$f_z = y + 2xz.$$

**Suggested questions:** Q3-13.

## 2.3 Differentials and Directional Derivatives

How can we find the rate of change of a function $f(x)$ when moving in some arbitrary direction? When we have a single variable, the *gradient* gives us the infinitesimal change of the function $\mathrm{d}f$ with an infinitesimal change in $\mathrm{d}x$, that is we have an equation which looks like this:

$$\mathrm{d}f = \frac{\mathrm{d}f}{\mathrm{d}x}\,\mathrm{d}x. \tag{2.1}$$

The objects representing *small* changes, $\mathrm{d}f, \mathrm{d}x$ are called *differentials.*

Let us now generalise this to the multivariable case. When we have more than one variable we can imagine moving in an arbitrary direction. Let's consider two variables $(x, y)$ for concreteness and imagine that we move along by $\mathrm{d}x$ in the $x$ direction, and $\mathrm{d}y$ in the $y$ direction:

$$(x, y) \mapsto (x + \mathrm{d}x, y + \mathrm{d}y).$$

> **Definition**
>
> Now if we have a function $f(x, y)$, then the change $\mathrm{d}f$ is given by the sum of the contributions from each of the directions, and $\mathrm{d}f$ is called the *total differential*:
>
> $$\mathrm{d}f = \frac{\partial f}{\partial x}\,\mathrm{d}x + \frac{\partial f}{\partial y}\,\mathrm{d}y. \tag{2.2}$$

### Examples

> **Try it out**
>
> 1. Find the total differential of $f(x, y) = x^2 y^3 + \cos xy$.
>
> **Answer:** We have
> $$f_x = 2xy^3 - y\sin xy$$
> $$f_y = 3x^2 y^2 - x\sin xy$$
> so
> $$\mathrm{d}f = f_x\,\mathrm{d}x + f_y\,\mathrm{d}y$$
> $$= \left[2xy^3 - y\sin xy\right]\mathrm{d}x + \left[3x^2 y^2 - x\sin xy\right]\mathrm{d}y.$$

> **Example**
>
> Let's look at an example from physics. Recall the *First Law of Thermodynamics*
>
> $$\mathrm{d}U = T\mathrm{d}S - P\mathrm{d}V$$
>
> where the state variables $S$ and $V$ are the entropy and volume of a closed homogeneous system, respectively, and where $U(S, V)$ is the internal energy, $T(S, V)$ is the temperature and $P(S, V)$ is the pressure. Comparing with Equation 2.2, we see that
>
> $$T \equiv \frac{\partial U}{\partial S} \qquad \text{and} \qquad P \equiv -\frac{\partial U}{\partial V}.$$

(Note: in thermodynamics it is often helpful to use an alternative notation for partial derivatives to make it explicit which variable is being held fixed, but we won't do that here.)

The formal equivalence of the mixed $S, V$ derivatives gives an interesting relation:

$$\frac{\partial T}{\partial V} = \frac{\partial}{\partial V}\left(\frac{\partial U}{\partial S}\right) = \frac{\partial^2 U}{\partial S \partial V} = \frac{\partial^2 U}{\partial V \partial S} = \frac{\partial}{\partial S}\left(\frac{\partial U}{\partial V}\right) = -\frac{\partial P}{\partial S}.$$

The identity $\frac{\partial T}{\partial V} = -\frac{\partial P}{\partial S}$ is a property of any physical thermodynamic system. Those of you who are physicists will likely encounter it again: it is called a *Maxwell relation.*

### 2.3.1 Exact and Inexact differentials

Let us generalise slightly in the case of two variables: the most general differential can be written as

$$\mathrm{d}f = a(x, y)\ \mathrm{d}x + b(x, y)\ \mathrm{d}y. \tag{2.3}$$

There is something slightly misleading about the notation $\mathrm{d}f$: it suggests that all differentials can be written as "d of a function $f$". As it turns out, this is not true, and there are actually two different kind of differentials: *exact* and *inexact*.

> **Definition**
>
> A differential $\mathrm{d}f = a(x, y)\ \mathrm{d}x + b(x, y)\ \mathrm{d}y$ is said to be *exact* if there exists a function $f(x, y)$ such that $\mathrm{d}f$ is the total differential of $f$; that is, if there exists a function $f(x, y)$ such that
>
> $$a(x, y) = f_x \quad \text{and} \quad b(x, y) = f_y.$$
>
> (You should imagine that this means that the functions $a(x, y)$ and $b(x, y)$ above have been correctly chosen so that we can integrate to obtain $f(x, y)$.)
>
> If this *cannot* be done – i.e. if there does not exist any $f$ that makes this possible – then $\mathrm{d}f$ is an *inexact* differential.

It should seem intuitively reasonable that a random choice of functions $a(x, y)$ and $b(x, y)$ will probably result in an inexact differential.

**Examples**

> **Try it out**
>
> Show that $\mathrm{d}f = y\,\mathrm{d}x + x\,\mathrm{d}y$ is an exact differential.
>
> **Answer:**
>
> The question is whether $f_x = y$ and $f_y = x$ can be integrated to find $f$. We have
>
> $$\begin{aligned} f_x = y &\implies f(x, y) = xy + A(y), \\ f_y = x &\implies f(x, y) = xy + B(x). \end{aligned}$$
>
> Note that in the first line $A(y)$ can be an arbitrary function of $y$, but not of $x$. The second line is the opposite, $B(x)$ can be an arbitrary function of $x$ but not of $y$. These two statements together imply

that it must just be a constant, i.e. $A(y) = B(x) = C$. So

$$f(x, y) = xy + C.$$

As we have succeeded in finding $f$, the differential $df$ is indeed exact. (Note that $f$ is not uniquely defined! It will always be ambiguous up to a constant).

An alternative way to find $f$ is the following: from $f_x = y$ we integrate as above to find $f(x, y) = xy + A(y)$. Now, we can use this expression to compute $f_y$ and the outcome must agree with $f_y = x$. That is, we have

$$\begin{aligned} f_y &= x & \text{(from } df) \\ \text{and} \quad f_y &= x + A'(y) & \text{(from differentiating } f(x, y) = xy + A(y)) \, , \end{aligned}$$

so we can conclude that $A'(y) = 0$. By integration, this means that $A(y)$ (a function of one variable!) must be constant, so $f(x, y)$ must be of the form $f(x, y) = xy + C$, as before.

---

**Try it out**

Show that $df = y \, dx - x \, dy$ is an *in*exact differential.

**Answer:**

Note that the differential $df$ differs from the previous example only in the sign of the $dy$ component, so you might initially think that this differential should also be exact. Again, the question is whether $f_x = y$ and $f_y = -x$ can be integrated to find $f$. Suppose that we can do this. Then we have that

$$\begin{aligned} f_x &= y \implies f(x, y) = xy + A(y), \\ f_y &= x \implies f(x, y) = -xy + B(x). \end{aligned}$$

As before, note that $A(y)$ can be an arbitrary function of $y$, but does not depend on $x$, while $B(x)$ can be an arbitrary function of $x$, but does not depend on $y$. These two statements together imply that $B(x) = 2xy + A(y)$, which is nonsense, since we know that the function $B(x)$ on the left-hand side does not depend on $y$, while equality with the right-hand side says that it does. Hence, there is no function $f(x, y)$ whose total differential is $df$.

---

**Try it out**

Show that $df = 3y \, dx + x \, dy$ is an *in*exact differential.

**Answer:**

Let us imagine that $df$ is exact. Then we would have

$$\begin{aligned} f_x &= 3y \implies f = 3xy + A(y) \\ f_y &= x \implies f = xy + B(x) \end{aligned}$$

For similar reasons to the previous example, the two lines are contradictory, so there is no function $f(x, y)$ whose total differential is $df$.

#### 2.3.1.1 Testing for exactness

We may use Clairault's Theorem (i.e. the commutativity of mixed partials) to test for an exact differential. i.e. suppose that we are given a differential

$$\mathrm{d}f = a(x, y)\,\mathrm{d}x + b(x, y)\,\mathrm{d}y.$$

If this is exact, it means that $a \equiv f_x$ and $b \equiv f_y$ for *some* choice of $f(x, y)$. Therefore by Clairault's Theorem

$$a_y = f_{xy} = f_{yx} = b_x. \tag{2.4}$$

Thus we see that if the differential is exact, then this condition on $a, b$ is satisfied. We will not prove the converse statement here, but it turns out the converse is true, i.e. if Equation 2.4 is true, then a function $f(x, y)$ exists[2] so that $a = f_x, b = f_y$. It should make sense to see how this generalises to more dimensions; for example for three dimensions we would write:

$$\mathrm{d}f = a(x, y, z)\,\mathrm{d}x + b(x, y, z)\,\mathrm{d}y + c(x, y, z)\,\mathrm{d}z \,,$$

and then to check for exactness we need to check all pairs, $a_y = b_x$, $a_z = c_x$, $b_z = c_y$. (Check that this is enough if this is not clear).

**Suggested questions:** Q14-19.

## 2.4 The gradient of a function and a first look at vector calculus

Recall that in vector notation, a point is specified by: $x = xi + yj$. Now let's again imagine changing the position as $(x, y) \to (x + \mathrm{d}x, y + \mathrm{d}y)$. Note that we can write this in a vector form as:

$$x \to x + \mathrm{d}x$$

where $\mathrm{d}x$ is the *infinitesimal vector* $\mathrm{d}x = dx\,i + dy\,j$. Note carefully what's happening here – $\mathrm{d}x$ is a vector, and we move a distance $\mathrm{d}x$ in the $i$ direction and $\mathrm{d}y$ in the $j$ direction.

Now recall we had the following expression for the change of the function $f(x, y)$

$$\mathrm{d}f = \frac{\partial f}{\partial x}\,\mathrm{d}x + \frac{\partial f}{\partial y}\,\mathrm{d}y.$$

Let's write the differential in a fancy vector notation:

$$\mathrm{d}f = \frac{\partial f}{\partial x}\,\mathrm{d}x + \frac{\partial f}{\partial y}\,\mathrm{d}y$$
$$= \mathrm{d}x \cdot \left( \frac{\partial f}{\partial x}i + \frac{\partial f}{\partial y}j \right)\,.$$

Here the dot is the usual dot product. In practice, when dealing with vectors we typically suppress the $i$ and $j$ notation, writing instead $x = (x, y)$, $\mathrm{d}x = (dx, dy)$, $\frac{\partial f}{\partial x}i + \frac{\partial f}{\partial y}j = \left( \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$, etc.

The expression Equation 2.5 now suggests that we define a new vector object $\nabla f$.

---

[2]Actually this is generally only true *locally*, i.e. the function $f(x, y)$ that we build might have some problems if we try to define it in all space. In fact this works only when the spaces we are considering are "simple"; if they have holes in them etc. then we cannot globally define $f$. This is thus a connection between "topology" (i.e. global properties of spaces) and calculus. It is also a surprisingly important subject to physicists: google "de Rham cohomology" to find out more.

> **Definition**
>
> The "*gradient* of *f*", (or "*grad f*" or "*del f*" for short) is given by
>
> $$\nabla f = \left( \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right).$$
>
> It is a *vector (field)* whose components are the partial derivatives of *f*. Note that the gradient can also be written in terms of basis vectors as
>
> $$\nabla f = \frac{\partial f}{\partial x} i + \frac{\partial f}{\partial y} j.$$

It is also convenient to think of $\nabla$ as an object in its own right – it is a *vector differential operator* $\nabla = (\frac{\partial}{\partial x}, \frac{\partial}{\partial y})$ and is called *del, grad,* or *nabla.*

In terms of the gradient, we can write the expression for the differential of the function *f* as:

$$\mathrm{d}f = \mathrm{d}x \cdot \nabla f.$$

Let's introduce one more bit of terminology.

> **Definition**
>
> The *directional derivative* of *f* in the direction of a vector *u* is
>
> $$\nabla_u f = \frac{1}{|u|}\, u \cdot \nabla f.$$
>
> It is a scalar quantity, representing the rate of change of *f* in the direction of **u**.

### 2.4.1 What do $\nabla f$ and $\nabla_u f$ mean?

We now have quite a lot of formalism. Let us work out an example. Consider the "bowl" function

$$f(x, y) = x^2 + y^2.$$

What is $\nabla f$? We have

$$\nabla f(x,y) = \left( \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right) = (2x, 2y) \qquad \text{or} \qquad \nabla f = 2x\, i + 2y\, j.$$

Let us draw a picture; we see that the gradient points *outwards*, which is also the direction in which the function $f(x, y)$ gets bigger and bigger. This is always true – the gradient always points **in the direction of greatest increase.**

Let us now prove this with equations. For any direction $u$, the directional derivative is

$$\begin{aligned} \nabla_u f &= \frac{1}{|u|} u \cdot \nabla f \\ &= 1 \cdot |\nabla f| \cos\theta, \end{aligned}$$

where $\theta$ is the angle between the two vectors $u$ and $\nabla f$. The directional derivative is therefore a maximum when $\cos\theta = 1$; that is, when $\theta = 0$, meaning that $u$ and $\nabla f$ point in the same direction – in other words,

the gradient always points in the direction of greatest increase, as claimed. Notice, also, that the rate of change of $f$ in the direction of $\nabla f$ (i.e.the maximal rate of change) is

$$\nabla_{\nabla f} f = \frac{1}{|\nabla f|} \nabla f \cdot \nabla f = \frac{|\nabla f|^2}{|\nabla f|} = |\nabla f|.$$

The function $f(x, y)$ is, by definition, constant along the *level curves* (contours) of the surface described by $z = f(x, y)$. This means that in directions tangent to the level curves the directional derivative should be 0 (as $f(x, y)$ is not changing in this direction); that is, we should have $u \cdot \nabla f = 0$ (i.e. $\cos \theta = 0$) – in other words **the level curves are at right angles to $\nabla f$**.

Moreover, if we look at the level curves (contours) of the surface on the $xy$-plane, then the fact that $\nabla f$ is the direction of greatest increase means that it should always point "**up the slope**".

---

**Examples**

Let $f(x, y) = x^2 + y^2$ as before and consider the paraboloid or "bowl" described by $z = x^2 + y^2$. The gradient of $f$ is

$$\nabla f(x, y) = (2x, 2y).$$

The level curves to the surface are all of the form $x^2 + y^2 = c$. That is, if $c > 0$ then the level curve is a circle, if $c = 0$ then the level curve consists of a single point (the origin), while if $c < 0$ then the level curve is empty. Viewed in the $xy$-plane, the level curves look like a collection of concentric circles around the origin (with $c \geq 0$ being the square of the radius). Let's focus on the interesting case where the level curves are circles, i.e. where $c > 0$. As we just learned, the gradient $\nabla f = (2x, 2y)$ must always be perpendicular to the level curves (circles), and this is clear in this case from drawing a simple sketch. Moreover, the gradient points in the direction (outwards, perpendicular to the level curves) of greatest increase of the function $f$; that is, the level curves are the circles $f(x, y) = c$ and, hence, $f$ increasing is the same as the radius of the circles increasing.

Can we write down unit vectors $\hat{t}$ tangent to the level curves? By *inspection* (which means guessing from looking at the pictures, but sounds fancier), we can see that suitable unit vectors are given by

$$\hat{t} = \frac{1}{\sqrt{x^2 + y^2}} (-y, x).$$

---

**Try it out**

Find the rate of change of $f(x, y) = y^4 + x^2 y^2 + x$ at $(0, 1)$ in the direction of the vector $i + 2j$.
**Answer:**
First find $\nabla f$. We have

$$f_x = 1 + 2xy^2 = 1$$
$$f_y = 4y^3 + 2yx^2 = 4$$
$$\nabla f = (1 + 2xy^2)i + (4y^3 + 2yx^2)j,$$

so $\nabla f(0, 1) = i + 4j$. Now we need the *unit vector* in the direction of $i + 2j$. This is

$$\hat{n} = (i + 2j)/(1 + 2^2)$$
$$= \frac{1}{\sqrt{5}}(i + 2j)$$

Therefore, the rate of change of $f$ at $(0, 1)$ in the direction $i + 2j$ is

$$\hat{n} \cdot \nabla f(0, 1) = \frac{1}{\sqrt{5}}(i + 2j) \cdot (i + 4j) = \frac{9}{\sqrt{5}}.$$

**Try it out**

The temperature on a metal plate is $T(x, y) = x^2 e^{-y}$. At the point $(2, 1)$, in what direction does the temperature increase most rapidly?

**Answer:**

First find $\nabla T$. We have

$$T_x = 2xe^{-y} = 4/e\,,$$
$$T_y = -x^2 e^{-y} = -4/e\,,$$
$$\nabla T = 2xe^{-y}i - x^2 e^{-y}j,$$

so $\nabla T(2, 1) = \frac{4}{e}(i - j)$. Therefore, (a unit vector in) the direction of greatest increase at $(2, 1)$ is

$$\frac{1}{\sqrt{2}}(i - j),$$

and the rate of increase in this direction at $(2, 1)$ is

$$|\nabla T(2, 1)| = \frac{4\sqrt{2}}{e}\,.$$

**Try it out**

Find the level curve of $f(x, y) = y^4 + x^2 y^2 + x$ through the point $(0, 1)$ and verify that its tangent at this point is orthogonal to $\nabla f$.

**Answer:**

The level curves are defined by $c = f(x, y) = y^4 + x^2 y^2 + x$. At $(0, 1)$ you can easily verify that $f(0, 1) = 1$, so we must have $c = 1$. Thus, the equation of the level curve is $y^4 + x^2 y^2 + x = 1$. A tangent vector to this level curve at to this point $(0, 1)$ has "slope" $\frac{dy}{dx}(0, 1)$ in the $xy$-plane. Hence, differentiating we find

$$\frac{dy}{dx}(4y^3 + 2yx^2) + 2xy^2 + 1 = 0,$$

$$\text{and, thus,} \quad \frac{dy}{dx}(0, 1) = -\frac{1}{4}.$$

So the corresponding unit vector in the $xy$-plane (i.e. having this slope) is

$$\hat{p} = \frac{1}{\sqrt{17}}(4, -1).$$

Next $\nabla f$ is given by

$$\nabla f = (2xy^2 + 1, \ 4y^3 + 2yx^2),$$

so $\nabla f(0, 1) = (1, 4)$ and, therefore, $\nabla f(0, 1) \cdot \hat{p} = 0$, as desired.

### 2.4.2 $\nabla f$ in arbitrary dimensions

We've been working in two dimensions, but of course all of the concepts generalise to any number of dimensions. Let $f(x_1, x_2, \ldots, x_n)$ be a function depending on $n$ variables and let $e_1, e_2, \ldots, e_n$ be the standard basis of $\mathbb{R}^n$ (i.e. one unit vector along each coordinate axis). Then the gradient of $f$ is a function $\nabla f : \mathbb{R}^n \to \mathbb{R}^n$ such that

$$\nabla f = \sum_{i=1}^n \frac{\partial f}{\partial x_i}\, e_i = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \ldots, \frac{\partial f}{\partial x_n} \right), \tag{2.5}$$

and the directional derivative (i.e. rate of change) of $f$ in the direction of a unit vector $\hat{u} = \sum_{i=1}^n u_i\, e_i = (u_1, u_2, \ldots, u_n)$ is given by

$$\nabla_{\hat{u}} f = \hat{u} \cdot \nabla f \, .$$

For example, the temperature in this room can be written $T(x, y, z)$ and we can define the rate of change of $T$ when moving in some arbitrary direction $\hat{u}$ in the same fashion.

Finally, note that the gradient operator satisfies the following two properties when acting on scalar functions:

1. Distributivity: $\nabla(f + g) = \nabla f + \nabla g$.

2. Product rule: $\nabla(fg) = (\nabla f)g + f\nabla g$.

Both of these properties should be familiar from ordinary derivatives: they follow from the definition of the gradient in a particular basis, i.e. Equation 2.5.

**Suggested questions:** Q20-23.

# 3 Applications of Partial Derivatives

In this chapter we will discuss a few applications of partial derivatives.

## 3.1 The chain rule in multiple variables

Suppose we have functions $f(x)$ and $g(x)$, each depending on a single variable. Then we can *compose* them to get either a function $(f \circ g)(x) = f(g(x))$ (first do $g$, then do $f$) or a function $(g \circ f)(x) = g(f(x))$ (first do $f$, then do $g$), each of which also depends on only one variable. For example, the functions $f(x) = x^2$ and $g(x) = \sin(x)$ can be composed in to give either $(f \circ g)(x) = f(g(x)) = (\sin(x))^2$ (first $g$, then $f$) or $(g \circ f)(x) = g(f(x)) = \sin(x^2)$ (first $f$, then $g$).

Recall that the chain rule (in one variable) tells us how to differentiate compositions of functions (of one variable). More precisely, recall that

$$\frac{\mathrm{d}}{\mathrm{d}x}(f \circ g)(x) = \frac{\mathrm{d}}{\mathrm{d}x}f((g(x))) = f'(g(x))\, g'(x).$$

You might be more familiar with thinking about the chain rule in the following (equivalent) way, which is perhaps more aesthetically pleasing (because we can imagine "canceling the '$\mathrm{d}x$'"): suppose you have a function $f(x)$ and suppose that the variable $x$ also depends on another variable $t$ (so that we have function $x(t)$). Then we can write the chain rule as

$$\frac{\mathrm{d}f}{\mathrm{d}t} = \frac{\mathrm{d}f}{\mathrm{d}x}\frac{\mathrm{d}x}{\mathrm{d}t}.$$

We can derive this formula by noting that the total differentials of $f(x)$ and $x(t)$ are given by

$$\mathrm{d}f = \frac{\mathrm{d}f}{\mathrm{d}x}\,\mathrm{d}x \qquad \text{and} \qquad \mathrm{d}x = \frac{\mathrm{d}x}{\mathrm{d}t}\,\mathrm{d}t,$$

which can then be combined to yield

$$\mathrm{d}f = \frac{\mathrm{d}f}{\mathrm{d}x}\frac{\mathrm{d}x}{\mathrm{d}t}\,\mathrm{d}t.$$

Dividing across by $\mathrm{d}t$ now gives us the chain rule

$$\frac{\mathrm{d}f}{\mathrm{d}t} = \frac{\mathrm{d}f}{\mathrm{d}x}\frac{\mathrm{d}x}{\mathrm{d}t}.$$

If we want to generalise the chain rule to functions $f(x_1, x_2, \ldots, x_n)$ of several variables, there are two scenarios we need to deal with:

1. all the variables $x_1, x_2, \ldots, x_n$ are functions of a single variable $t$; i.e. we have a composition $f(x_1(t), x_2(t), \ldots, x_n(t))$ depending on a single variable.

2. all the variables $x_1, x_2, \ldots, x_n$ are functions of several variables $u_1, u_2, \ldots, u_m$; i.e. we have a composition $f(x_1(u_1, u_2, \ldots, u_m), x_2(u_1, u_2, \ldots, u_m), \ldots, x_n(u_1, u_2, \ldots, u_m))$ depending on several variables.

In the first case, the chain rule will give the (ordinary) derivative $\frac{\mathrm{d}f}{\mathrm{d}t}$, while, in the second case, we will obtain the partial derivatives $\frac{\partial f}{\partial u_1}, \frac{\partial f}{\partial u_2}, \ldots, \frac{\partial f}{\partial u_m}$.

### 3.1.1 The chain rule for dependence on only one variable

Consider a function $f(x, y)$, where and $x(t)$ and $y(t)$ are both functions of just a single variable $t$.

You should imagine that we move through a curve $(x(t), y(t))$ in $\mathbb{R}^2$ which is parametrised by $t$, and we evaluate $f(x, y)$ at our instantaneous position, i.e. we compute $f(x(t), y(t))$. We then would like to ask how the combined function $f(x(t), y(t))$ changes as a function of $t$. We begin with the change in $f$ (the total differential), which is given by Equation 2.2:

$$\mathrm{d}f = \frac{\partial f}{\partial x}\,\mathrm{d}x + \frac{\partial f}{\partial y}\,\mathrm{d}y.$$

We also have

$$\mathrm{d}x = \frac{\mathrm{d}x}{\mathrm{d}t}\,\mathrm{d}t \qquad \text{and} \qquad \mathrm{d}y = \frac{\mathrm{d}y}{\mathrm{d}t}\,\mathrm{d}t.$$

So, just as in the case of the usual chain rule, we can now combine these expressions to obtain

$$\mathrm{d}f = \frac{\partial f}{\partial x}\frac{\mathrm{d}x}{\mathrm{d}t}\,\mathrm{d}t + \frac{\partial f}{\partial y}\frac{\mathrm{d}y}{\mathrm{d}t}\,\mathrm{d}t. \tag{3.1}$$

By dividing across by $\mathrm{d}t$, we find the final expression for the chain rule in two dimensions.

> **Key idea**
>
> For a composed function $f(x(t), y(t))$ depending on a single variable $t$, the chain rule is
>
> $$\frac{\mathrm{d}f}{\mathrm{d}t} = \frac{\partial f}{\partial x}\frac{\mathrm{d}x}{\mathrm{d}t} + \frac{\partial f}{\partial y}\frac{\mathrm{d}y}{\mathrm{d}t} = \left(\frac{\mathrm{d}x}{\mathrm{d}t}, \frac{\mathrm{d}y}{\mathrm{d}t}\right) \cdot \nabla f.$$

Let's understand what this formula is saying: as we change $t$ and we move along the path, there are two ways in which $f(x, y)$ can change: that arising from the change in $x$, and that arising from the change in $y$. That is why there are two terms. Note that it is quite easy to generalise to the $n$-variable case, as follows.

> **Key idea**
>
> For a composed function $f(x(t))$, where $x(t) = (x_1(t), ...x_n(t))$, the chain rule (describing how $f$ changes as $t$ varies) is
>
> $$\frac{\mathrm{d}f}{\mathrm{d}t} = \frac{\partial f}{\partial x_1}\frac{\mathrm{d}x_1}{\mathrm{d}t} + \frac{\partial f}{\partial x_2}\frac{\mathrm{d}x_2}{\mathrm{d}t} + \cdots + \frac{\partial f}{\partial x_n}\frac{\mathrm{d}x_n}{\mathrm{d}t} = \frac{\mathrm{d}x}{\mathrm{d}t} \cdot \nabla f,$$
>
> where $\frac{\mathrm{d}x}{\mathrm{d}t} = \left(\frac{\mathrm{d}x_1}{\mathrm{d}t}, \frac{\mathrm{d}x_2}{\mathrm{d}t}, ..., \frac{\mathrm{d}x_n}{\mathrm{d}t}\right)$.

As usual, we now work out some examples:

> **Try it out**
>
> Consider a cylinder of radius $x$ and height $y$. Suppose that the cylinder changes its size as $x = 3t$ and $y = 4 + t^2$. What is the rate of change of $V$ with respect to $t$?
> **Answer:**

*Method 1*: Direct substitution, i.e. just write everything in terms of $t$.

$$V = \pi x^2 y$$
$$= \pi 9 t^2 (4 + t^2).$$

Then

$$\frac{\mathrm{d}V}{\mathrm{d}t} = 72\pi t + 36\pi t^3.$$

*Method 2*: Use the chain rule,

$$\frac{\mathrm{d}V}{\mathrm{d}t} = \frac{\partial V}{\partial x}\frac{\mathrm{d}x}{\mathrm{d}t} + \frac{\partial V}{\partial y}\frac{\mathrm{d}y}{\mathrm{d}t}$$
$$= 2\pi x y \cdot 3 + \pi x^2 \cdot 2t$$
$$= \pi 18 t (4 + t^2) + \pi 18 t^3$$
$$= 72\pi t + 36\pi t^3.$$

---

### Try it out

For $f = \sin(xy)$ find $\frac{\mathrm{d}f}{\mathrm{d}t}$ along the curve parametrised by $x = t^2$, $y = t^3$.

**Answer:**

We have

$$\frac{\mathrm{d}f}{\mathrm{d}t} = \frac{\partial f}{\partial x}\frac{\mathrm{d}x}{\mathrm{d}t} + \frac{\partial f}{\partial y}\frac{\mathrm{d}y}{\mathrm{d}t}$$
$$= y\cos(xy)(2t) + x\cos(xy)(3t^2)$$
$$= 5t^4\cos(t^5).$$

Note that we could have again substituted and used $f(t) = \sin(t^5)$.

---

### Try it out

For $f(x, y, z) = 3xe^{yz}$ find the value of $\frac{\mathrm{d}f}{\mathrm{d}t}$ at the point on the curve $x = \cos t$, $y = \sin t$, $z = t$ where $t = 0$.

**Answer:**

We have

$$\frac{\mathrm{d}f}{\mathrm{d}t} = \frac{\partial f}{\partial x}\frac{\mathrm{d}x}{\mathrm{d}t} + \frac{\partial f}{\partial y}\frac{\mathrm{d}y}{\mathrm{d}t} + \frac{\partial f}{\partial z}\frac{\mathrm{d}z}{\mathrm{d}t}$$
$$= e^{yz}(-3\sin t + 3xz\cos t + 3yx).$$

Then

$$\left.\frac{\mathrm{d}f}{\mathrm{d}t}\right|_{t=0} = -3,$$

since at $t = 0$, $(x, y, z) = (1, 0, 0)$.

---

### Try it out

Let $f = f(x, t)$ where $x = x(t)$. What is $\frac{\mathrm{d}f}{\mathrm{d}t}$?

**Answer:**

The chain rule tells us that for functions $f(x, y)$ we have

$$\frac{\mathrm{d}f}{\mathrm{d}t} = \frac{\partial f}{\partial x}\frac{\mathrm{d}x}{\mathrm{d}t} + \frac{\partial f}{\partial y}\frac{\mathrm{d}y}{\mathrm{d}t}.$$

For this example we can take $y(t) = t$. Then we have

$$\frac{\mathrm{d}f}{\mathrm{d}t} = \frac{\partial f}{\partial x}\frac{\mathrm{d}x}{\mathrm{d}t} + \frac{\partial f}{\partial t}.$$

since $\frac{\mathrm{d}t}{\mathrm{d}t} = 1$. Note the crucial difference between the two operations $\frac{\mathrm{d}}{\mathrm{d}t}$ and $\frac{\partial}{\partial t}$; the first one *also* takes into account the implicit change arising from the fact that $f$ depends on $x$ which depends on $t$, whereas the second only takes into account the explicit change arising from the direct dependence on $t$ in the second argument.

### 3.1.2 The chain rule for dependence on several variables

Up till now everything ultimately depended only on a single variable $t$.

Suppose instead that as before we consider a function $f(x, y)$, where $x = x(u, v)$ and $y = y(u, v)$ are functions of *two* other variables $u$ and $v$. We may then consider the following composite function $f(x(u, v), y(u, v))$, and we might be interested in computing the partial derivatives $\frac{\partial f}{\partial u}$ and $\frac{\partial f}{\partial v}$.

**Key idea**

By the chain rule,

$$\frac{\partial f}{\partial u} = \frac{\partial f}{\partial x}\frac{\partial x}{\partial u} + \frac{\partial f}{\partial y}\frac{\partial y}{\partial u} \ ,$$
$$\frac{\partial f}{\partial v} = \frac{\partial f}{\partial x}\frac{\partial x}{\partial v} + \frac{\partial f}{\partial y}\frac{\partial y}{\partial v} \ .$$

(3.2)

**Try it out**

Work through the calculations in terms of $\mathrm{d}f$s, as in Equation 3.1, to derive Equation 3.2.

*Note that to find $\frac{\partial f}{\partial x}$ we hold $y$ constant, but to find $\frac{\partial x}{\partial u}$, we need to hold $v$ constant.*

The generalization to $n$ variables $x_1, x_2, \dots, x_n$ which depend on $m$ other variables $u_1, u_2, \dots, u_m$ is straightforward: for each $i \in \{1, 2, \dots, n\}$ we have a function $x_i = x_i(u_1, u_2, \dots, u_m)$, so for each variable $u_j$, $j \in \{1, 2, \dots, m\}$, we have a chain rule

$$\frac{\partial f}{\partial u_j} = \sum_{i=1}^{n} \frac{\partial f}{\partial x_i}\frac{\partial x_i}{\partial u_j} \qquad \text{for each } j \in \{1, 2, \dots, m\}.$$

For those of you who like matrices, this is a good time to note that we can write all $m$ of these chain rules concisely as a single *matrix equation*

$$\left( \frac{\partial f}{\partial u_1}, \frac{\partial f}{\partial u_2}, \dots, \frac{\partial f}{\partial u_m} \right) = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right) \left[ \frac{\partial x}{\partial u} \right] = \nabla f \left[ \frac{\partial x}{\partial u} \right],$$

where we consider the object $\left[\frac{\partial x}{\partial u}\right]$ as an $(n \times m)$-matrix with $(i,j)^{\text{th}}$ entry $\frac{\partial x_i}{\partial u_j}$. For those of you who *don't* like matrices[1], you may happily ignore this for now.

**Examples**

> **Try it out**
>
> Let $f(x, y) = xy$ where $x = u \cos v$ and $y = u \sin v$. Compute $\frac{\partial f}{\partial u}$ both from direct substitution and using the chain rule.
> **Answer:**
> Substituting gives $f(u, v) = u^2 \sin v \cos v$ and $f_u = u \sin 2v$. So we should expect $f_u = 2u \sin v \cos v$. Let's check: the chain rule gives
> $$f_u = f_x x_u + f_y y_u$$
> $$= y \cos v + x \sin v$$
> $$= 2u \sin v \cos v$$
> as required.

> **Try it out**
>
> If $f$ is a function of $x, y$ where $x = u^2 - v^2$ and $y = 2uv$ show that $u f_u + v f_v = 2x f_x + 2y f_y$ and that $f_{uv} = 2f_y - 2y f_{xx} + 4x f_{xy} + 2y f_{yy}$.
> **Answer:**
> Note that since we do not know the explicit form of $f$, there is no possibility of computing its partial derivatives $f_u$ and $f_v$ via a substitution. Nevertheless, from
> $$\frac{\partial f}{\partial u} = \frac{\partial f}{\partial x}\frac{\partial x}{\partial u} + \frac{\partial f}{\partial y}\frac{\partial y}{\partial u} \quad \text{and} \quad \frac{\partial f}{\partial v} = \frac{\partial f}{\partial x}\frac{\partial x}{\partial v} + \frac{\partial f}{\partial y}\frac{\partial y}{\partial v}$$
> $$= 2u f_x + 2v f_y \qquad\qquad = -2v f_x + 2u f_y$$
> we conclude that
> $$u f_u + v f_v = \left(2u^2 f_x + 2uv f_y\right) + \left(-2v^2 f_x + 2uv f_y\right)$$
> $$= 2x f_x + 2y f_y.$$
> Now, to compute $f_{uv}$, observe first that
> $$f_{uv} = \frac{\partial}{\partial v}(2u f_x + 2v f_y) = 2u \left(f_x\right)_v + 2f_y + 2v \left(f_y\right)_v.$$
> From the chain rule (together with Clairault's Theorem) we have
> $$\left(f_x\right)_v = f_{xx} x_v + f_{xy} y_v = -2v f_{xx} + 2u f_{xy},$$
> $$\left(f_y\right)_v = f_{yx} x_v + f_{yy} y_v = -2v f_{xy} + 2u f_{yy}.$$
> Therefore, substituting these into $f_{uv}$ we find that
> $$f_{uv} = 2u \left(-2v f_{xx} + 2u f_{xy}\right) + 2f_y + 2v \left(-2v f_{xy} + 2u f_{yy}\right)$$
> $$= 2f_y - 4uv f_{xx} + 4(u^2 - v^2) f_{xy} + 4uv f_{yy}$$
> $$= 2f_y - 2y f_{xx} + 4x f_{xy} + 2y f_{yy}.$$

---

[1] ...but why don't you like matrices?

I will highlight one application of the chain rule: recall that there are multiple coordinate systems we can use for $\mathbb{R}^2$; we can use the regular Cartesian $(x, y)$, or the polar coordinates $(r, \theta)$. These are related by

$$x = r \cos$$
$$y = r \sin$$

and you have seen in the earlier part that the unit vectors are also related by

$$e_r = \cos\ i + \sin\ j$$
$$e_\theta = -\sin\ i + \cos\ j.$$

Now let's think about the *gradient*; up till now, we have only discussed the gradient of a scalar function $f$ in Cartesian coordinates:

$$\nabla f(x, y) = \frac{\partial f}{\partial x}i + \frac{\partial f}{\partial y}j.$$

What happens in polar coordinates? It turns out that it is possible to express everything above in polar coordinates – we do this by using the chain rule to replace $\frac{\partial f}{\partial x}$ with $\frac{\partial f}{\partial r}$ and so on. The derivation is spelled out at the end of the lecture notes if you are interested; when you do it you find the **gradient in polar coordinates**:

$$\nabla f(r, ) = \frac{\partial f}{\partial r}e_r + \frac{1}{r}\frac{\partial f}{\partial \theta}e. \tag{3.3}$$

This is basically what you would expect except for the interesting factor of $\frac{1}{r}$ on the last term – do you understand why this is there? I will leave you to ponder the geometry and figure out what this is saying.

**Suggested questions:** Q1-8

## 3.2 Multivariate Taylor expansions

In this section we will learn how to do a Taylor expansion in multiple variables, as well as understanding how to classify the different sorts of **critical points** that can happen for a function of multiple variables.

### Recap: The single-variable case

In principle this is a recap, but in practice it may very well be the first time you see this. Let us understand the idea of a **Taylor series expansion**. Suppose that we have a smooth function of a single variable $x$ (that is infinitely differentiable at a point $a$).

The Taylor series expansion tries to find a *polynomial expression* that approximates the function in the neighbourhood of $a$. The higher the order of polynomial we choose the better the approximation can be, and the further we can get from $x = a$ while still having a reasonable approximation. In order to derive the general form for these polynomials, suppose that such a thing exists and has the form

$$f(x) \approx P_n(x) = c_0 + c_1(x - a) + c_2(x - a)^2 + c_3(x - a)^3 + ... + c_n(x - a)^n.$$

That is we are taking a polynomial of order $n$ to approximate the function $f(x)$. If $|x - a| \ll 1$ then the approximation should improve as we increase $n$ (the condition for this is known as Taylor's Theorem which we won't discuss in this course).

Before showing you how to find the $c$'s, I feel I should address the basic philosophical question: **why on earth** would you want to do this? In full honesty this is one of the most useful things we will learn in

this course. The reason is that generally if you are trying to solve a real-life problem of **any** sort, the kinds of $f(x)$ that you get are just insanely hideous and impossible to work with. On the other hand, if it can be well approximated by a polynomial, usually you can make some progress. It is not a terrible oversimplification to say that the vast majority of physics consists of solving a system where only $c_1$ and $c_2$ are nonzero (which you can usually do), and then spending your entire career trying to figure out how to put back $c_3$.

Returning to mathematics: what are the values of the coefficients $c_k$? If $P_n(x)$ is to be an approximation to $f(x)$ near $x = a$, then the very least we might expect is that $P_n(a) = f(a)$, i.e. that they agree when $x = a$. But $P_n(a) = c_0$, so we set the constant $c_0 = f(a)$. In a similar way, it is reasonable to expect that, for $P_n(x)$ to be a good approximation to $f(x)$ near $x = a$, all of the derivatives (up to the $n^{\text{th}}$) of these functions must agree when $x = a$. For the first derivative, this means that

$$c_1 = P_n'(a) = f'(a),$$

and continuing in this fashion we find that

$$c_k = \frac{1}{k!}\frac{\mathrm{d}^k P_n}{\mathrm{d}x^k}\bigg|_{x=a} = \frac{1}{k!}\frac{\mathrm{d}^k f}{\mathrm{d}x^k}\bigg|_{x=a}, \quad k \in \{1, 2, \ldots, n\}.$$

Note that it is often convenient to write $f^{(k)}(x)$ as a shorthand notation for the $k^{\text{th}}$ derivative of $f(x)$ (since the $'$ notation gets messy for higher-order derivatives), and with this notation we can write $c_k = \frac{1}{k!}f^{(k)}(a)$ for all $k \in \{1, 2, \ldots, n\}$.

---

**Definition**

If a function $f(x)$ is $n$-times differentiable at a point $a \in \mathbb{R}$, then the degree-$n$ polynomial

$$P_n(x) = f(a) + f'(a)(x - a) + \frac{f''(a)}{2}(x - a)^2 + \frac{f'''(a)}{3!}(x - a)^3 + \cdots + \frac{f^{(n)}(a)}{n!}(x - a)^n$$

$$= f(a) + \sum_{k=1}^{n}\frac{f^{(k)}(a)}{k!}(x - a)^k$$

is called the *Taylor polynomial of degree $n$* for $f(x)$ around $a \in \mathbb{R}$ and, for $x$ near $a$ (i.e. $|x - a| \ll 1$), we have (Taylor's Theorem)

$$f(x) \approx P_n(x).$$

If $f(x)$ is infinitely differentiable at $x = a$, then the *Taylor series expansion* of $f(x)$ around $a \in \mathbb{R}$ is the infinite series

$$T_{f,a}(x) = f(a) + \sum_{k=1}^{\infty}\frac{f^{(k)}(a)}{k!}(x - a)^k.$$

---

It is tempting to believe that, by letting $n$ go to infinity and, hence, getting better and better approximations, we should end up in a situation where the Taylor series expansion is equal to the function itself (near $a$), i.e. that we should obtain an equality $f(x) = T_{f,a}(x)$ near $a$. While this is true in many familiar settings, it is not true in general, as we shall see below. Functions $f(x)$ which are equal to their Taylor series expansions (around a point $a$) are said to be *real analytic*.

The Taylor series expansion $T_{f,a}(x)$ around $x = a$ of a function $f(x)$ will always converge for at least one value of $x$, since $T_{f,a}(a) = f(a)$. If there exists some number $R > 0$ such that $T_{f,a}(x)$ converges for every $x \in (a - R, a + R)$ (i.e. for every $x \in \mathbb{R}$ such that $|x - a| < R$), then we say that $T_{f,a}(x)$ has *radius of convergence $R$*.

(**NB:** All the coefficients in the Taylor polynomial and Taylor series are numbers obtained from evaluating the derivatives of $f(x)$ at $x = a$. The end result should involve only summations of scalar multiples of terms of the form $(x - a)^k$. If you ever find yourself writing anything other than such terms (e.g. $e^x$ or $\sin x$) when computing the Taylor polynomial or Taylor series, then this is **not** the right idea.)

---

**Examples**

We would like to approximate the infinitely differentiable function $f(x) = e^x$ around $x = 0$ (i.e. $a = 0$). Observe first that, since $f^{(n)}(x) = e^x$ for every $n \in \mathbb{N} = \{1, 2, 3, ... \}$, we have $f^{(n)}(0) = 1$ for all $\in \mathbb{N}$. Then the Taylor polynomial of degree $n$ for $e^x$ around $x = 0$ is

$$P_n(x) = 1 + x + \frac{x^2}{2} + \frac{x^3}{3!} + \cdots + \frac{x^n}{n!},$$

while the Taylor series expansion of $e^x$ around $x = 0$ is

$$T_{f,0}(x) = \sum_{k=0}^{\infty} \frac{x^k}{k!}.$$

It is a fact (and, indeed, in some places this is used as the definition of the exponential function) that, in this case, we have

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

for all $x \in \mathbb{R}$ (so $e^x$ is a real analytic function and its Taylor series has radius of convergence $\infty$).

---

**Try it out**

Let $f(x) = \sin x$ around $x = 0$. Then $f(0) = \sin 0 = 0$, $f'(0) = \cos 0 = 1$, $f''(0) = -\sin 0 = 0$, $f'''(0) = -\cos 0 = -1$, $f^{(4)}(0) = \sin 0 = 0$ and so on. Thus, for each $n \in \mathbb{N}$ and for $x$ near 0 we have

$$\sin x \approx x - \frac{x^3}{3!} + \frac{x^5}{5!} - \cdots + (-1)^n \frac{x^{2n+1}}{(2n+1)!}.$$

It can once again be shown that $\sin x$ is equal to its Taylor series expansion around 0, so that

$$\sin x = \sum_{k=0}^{\infty} \frac{(-1)^k x^{2k+1}}{(2k+1)!}$$

for all $x \in \mathbb{R}$.

To illustrate how the Taylor polynomials around $x = 0$ approximate $\sin x$, here is an image showing the Taylor polynomial of degree 61 for $\sin x$ (the red curve), together with $\sin x$ itself (the blue dashed curve).
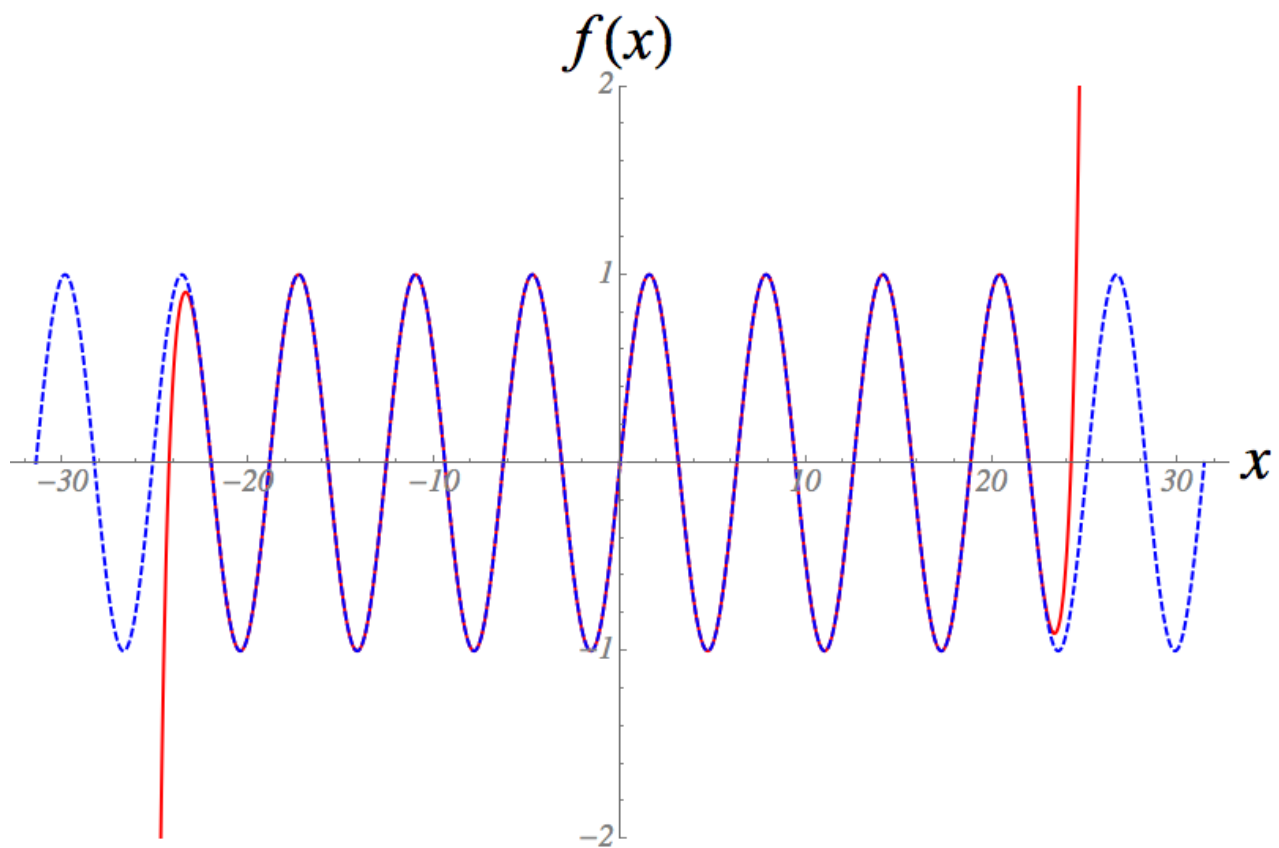
Figure 3.1: image

Similarly, it can be shown that the following functions equal their Taylor series expansions:

$$\cos x = \sum_{k=0}^{\infty} \frac{(-1)^k x^{2k}}{(2k)!} \qquad \text{for all } x \in \mathbb{R},$$

$$\log(1+x) = \sum_{k=1}^{\infty} \frac{(-1)^{k+1} x^k}{k} \qquad \text{for all } x \in \mathbb{R} \text{ with } |x| < 1,$$

$$\frac{1}{1-x} = \sum_{k=0}^{\infty} x^k \qquad \text{for all } x \in \mathbb{R} \text{ with } |x| < 1.$$

In the latter two cases, the radius of convergence of the Taylor series is 1 (think about what happens if you choose a particular $x \in \{\pm 1\}$, i.e. $|x| = 1$, in each case), which coincides with the fact that the functions themselves are undefined for certain values of $x$.

### Try it out

The function

$$f(x) = \begin{cases} e^{-\frac{1}{x^2}}, & \text{if } x \neq 0, \\ 0 & \text{if } x = 0, \end{cases}$$

is infinitely differentiable at $x = 0$, with $f^{(k)}(0) = 0$ for all $k \in \mathbb{N}$. Therefore, the Taylor series

expansion of $f(x)$ around $x = 0$ is the constant zero function $T_{f,0}(x) \equiv 0$. Since $f(x) \neq 0$ whenever $x \neq 0$, by definition, it follows that

$$f(x) \neq T_{f,0}(x),$$

except at $x = 0$. Nevertheless, if you look at the graph of $f(x)$, you will see that the Taylor series expansion $T_{f,0}(x) \equiv 0$ is a very, very good approximation of $f(x)$ near $x = 0$.

## Critical points in 1 dimension

A "critical point" is the generic name for an extreme point in the system.

> **Definition**
>
> Let $f(x)$ be a differentiable function of one variable. A point $a \in \mathbb{R}$ is a *critical point* of $f(x)$ if
>
> $$\frac{\mathrm{d}f}{\mathrm{d}x}(a) = 0 \qquad \text{(i.e. } f'(a) = 0\text{)}.$$

A critical point $a$ can be of one of three types:

- a (local) minimum (and stable) if $f(x) > f(a)$ for all $x$ near $a$ $(x \neq a)$;

- a (local) maximum (and unstable) if $f(x) < f(a)$ for all $x$ near $a$ $(x \neq a)$;

- an inflection point, if it is neither a (local) maximum nor a (local) minimum.

We can use the Taylor polynomial/series to help determine which types of critical point we have.

Before doing this, perhaps it is helpful to indicate where the terminology "stable" and "unstable" comes from – this is from physics.

> **Examples**
>
> Let $V(x)$ be the *potential energy* at position $x \in \mathbb{R}$, and suppose that $x(t)$ is the position of a particle at time $t$. Recall that the force $F$ due to the potential energy satisfies the equation
>
> $$F(x) = -\frac{\mathrm{d}V}{\mathrm{d}x}$$
>
> and, from Newton's Second Law, we have that $x(t)$ satisfies
>
> $$m\frac{\mathrm{d}^2 x}{\mathrm{d}t^2} = F(x(t)).$$
>
> Consider the graph of the potential energy $V$. We can think of the position $x(t)$ as telling us how we move along the $x$-axis, and then we can imagine following the graph of $V$ to understand the behaviour of the potential energy of the particle as time evolves. Observe now that the two equations above tell us that, at each time $t$, there is a force $F(x(t))$ on the particle induced by the potential energy $V(x(t))$ which points *down the slope* of the graph of $V$.
>
> Let $x_0 = x(t_0)$ be a critical point of $V$ (so that $F(x_0) = 0$) and consider what happens if the position of the particle is perturbed slightly away from $x_0$. If $x_0$ is a local maximum of $V$, then this perturbation will cause the corresponding point on the graph to roll "down the hill" away from the top. If we look

at the $x$-axis while this is happening, we see that the force $F$ induced on the particle pushes it away from the position $x_0$, never to return; that is, the position $x_0$ is an *unstable* equilibrium point. On the other hand, if $x_0$ is a local minimum of $V$, then the perturbation will cause the corresponding point on the graph to oscillate a little and eventually return to its starting point. If we again look at the $x$-axis while this is happening, we see that the force $F$ induced on the particle keeps pushing it back towards the position $x_0$ until it comes to rest; that is, the position $x_0$ is a *stable* equilibrium point.

> **Try it out**
>
> Let $f(x) = \cos x$. It is easy to check that $x = 0$ is a critical point of $\cos x$. If we move a little away from $x = 0$ to $x = 0 + h$, we have (from the quadratic Taylor approximation to $\cos x$)
>
> $$\cos(0 + h) = \cos h \approx 1 - \frac{h^2}{2} .$$
>
> Clearly, the $h^2$ term is always negative. Thus, if we go a little away from the critical point $x = 0$, the value of $\cos x$ will always be less than 1, and we can conclude that $\cos x$ has a local maximum at $x = 0$.

In general, if $f(x)$ has a critical point at $x = a$, then we have

$$f(x) \approx f(a) + \frac{f''(a)}{2}(x - a)^2,$$

and we conclude that $x = a$ is a local maximum if $f''(a) < 0$ (even if we cannot be bothered to plot $f$), and that $x = a$ is a local minimum if $f''(a) > 0$. Therefore, we can discover and classify the critical points of a differentiable function $f(x)$ as follows.

- Find all points $a$ where $f'(a) = 0$.

- Find the numerical value of $P = f''(a)$.

- If $P < 0$ it is a maximum, if $P > 0$ it is a minimum. If $P = 0$ we cannot conclude what type of critical point it is, and we need more information.

An example where $P = 0$ occurs at the cricial point $x = 0$ of the function $f(x) = x^4$. Although $f(x) = x^4$ has a minimum at $x = 0$, we have $P = f''(0) = 0$. On the other hand, $x = 0$ is also a critical point of the function $g(x) = x^3$ and, again, $P = g''(0) = 0$, but in this case $g(x)$ has an inflection point at $x = 0$.

Back to SMB Term 2....

### 3.2.1 Multivariate Taylor expansions

Let me first slightly rephrase the Taylor series as a function of $h$, the displacement from $x = a$:

$$f(a + h) = f(a) + hf'(a) + \frac{h^2}{2!}f''(a) + ...$$

We can also write it as an *operator equation* using the fact that, as we just saw, $e^A = 1 + A + \frac{A^2}{2!} + ...$

$$f(a + h) = e^{h\frac{\mathrm{d}}{\mathrm{d}x}} f(x)|_{x=a} .$$

This last equation makes it obvious how to generalise: a function of two variables expanded about $(x, y) = (a, b)$ can be found by first expanding about $x = a$ and then about $y = b$; doing this explicitly we first get

$$f(a + h, b + k) = f(a, b + k) + hf_x(a, b + k) + \frac{h^2}{2!}f_{xx}(a, b + k) + ...$$

Next we approximate $f(a, b + k)$ and the derivatives in $x$, by in turn expanding them as Taylor series in $y$ about $y = b$. That is

$$f(a, b + k) = f(a, b) + kf_y(a, b) + \frac{k^2}{2!}f_{yy}(a, b) ...$$

$$f_x(a, b + k) = f_x(a, b) + kf_{xy}(a, b) + ...$$

$$f_{xx}(a, b + k) = f_{xx}(a, b) + kf_{xxy}(a, b) + ...$$

Keeping terms up to quadratics (obviously we can extend but it gets a bit laborious) we have:

---

**Definition**

The Taylor expansion of $f(x, y)$ in two dimensions, up to quadratic terms, is:

$$f(a + h, b + k) = f + hf_x + kf_y + \frac{1}{2}(h^2 f_{xx} + 2hk f_{xy} + k^2 f_{yy}) + ... \qquad (3.4)$$

with the understanding that $f$ and all its derivatives on the right-hand side are evaluated at $(a, b)$.

---

I now present an alternative way to arrive at this formula, which may or may not help you in understanding it. This form uses the operator understanding of the Taylor expansion. That is

$$\begin{aligned}
f(a + h, b + k) &= e^{h\partial_x}e^{k\partial_y}f(x, y)|_{x=a, y=b} \\
&= e^{h\partial_x + k\partial_y}f(x, y)|_{x=a, y=b} \\
&= \left(1 + h\partial_x + k\partial_y + \frac{1}{2}(h\partial_x + k\partial_y)^2 + ...\right)f(x, y)|_{x=a, y=b} \\
&= f + hf_x + kf_y + \frac{1}{2}(h^2 f_{xx} + 2hk f_{xy} + k^2 f_{yy}) + ...
\end{aligned}$$

Note that we have arrived at the same result. This approach is powerful but we will not use it too much in these lectures. As an aside, it is worth mentioning that the reason that we can combine the exponentials trivially, $e^{h\partial_x}e^{k\partial_y} \equiv e^{h\partial_x + k\partial_y}$, is that the operators in the exponent "commute", by which we mean that it doesn't matter which way round they go, namely $\partial_x\partial_y = \partial_y\partial_x$. This is thanks to Clairault's theorem again.

Note that in a space of variables $\mathbf{x} = (x_1, ..., x_n)$ we can write the expansion of $f$ at $\mathbf{x} = \mathbf{x_0} + \mathbf{h}$ as

$$f(\mathbf{x_0} + \mathbf{h}) = e^{\mathbf{h}.\nabla}f|_{\mathbf{x_0}} .$$

---

**Try it out**

Compute the Taylor polynomial of $\cos(x + y)$ about $(0,0)$ up to and including quadratic terms
**Answer:**

---

In the above we have $a = 0 = b$ and $h = x$ and $k = y$: so we have

$$f = \cos(x + y)$$
$$f_x = f_y = -\sin(x + y)$$
$$f_{xx} = f_{xy} = f_{yy} = -\cos(x + y).$$

Then

$$f(x, y) = f + xf_x + yf_y + \frac{1}{2}(x^2 f_{xx} + 2xy f_{xy} + y^2 f_{yy})|_{0,0}$$
$$= 1 - \frac{1}{2}(x + y)^2 + \dots$$

Note this is slightly trivial since we could have just expanded $z = x + y$.

---

**Try it out**

For $f(x, y) = \log(x + 2y)$, find the Taylor expansion about $(1,0)$.
In the above notation we can take $a = 1$, $b = 0$, $h = x - 1$ and $k = y$. We have

$$f|_{(1,0)} = \log(x + 2y)|_{(1,0)} = 0$$
$$f_x|_{(1,0)} = \frac{1}{x + 2y}|_{(1,0)} = 1$$
$$f_{xx}|_{(1,0)} = -\frac{1}{(x + 2y)^2}|_{(1,0)} = -1$$
$$f_y|_{(1,0)} = \frac{2}{x + 2y}|_{(1,0)} = 2$$
$$f_{yy}|_{(1,0)} = -\frac{4}{(x + 2y)^2}|_{(1,0)} = -4$$
$$f_{xy}|_{(1,0)} = -\frac{2}{(x + 2y)^2}|_{(1,0)} = -2,$$

so that

$$f(x, y) = (x - 1) + 2y - \frac{1}{2}((x - 1)^2 + 4(x - 1)y + 4y^2) + \dots$$
$$= z - \frac{z^2}{2} + \dots \quad [z = x + 2y - 1].$$

Again, we could have done this more simply. This is because we can combine all of the functional dependence of the function into a single variable $z$.

---

**Try it out**

Given $f(x, y) = ye^{xy}$, find the Taylor expansion of $f$ about $(2,3)$.
**Answer:**
This finally cannot be done so simply. Here we have $a = 2$, $b = 3$, $h = x - 2$ and $k = y - 3$.

Hence we have

$$f|_{(2,3)} = ye^{xy} = 3e^6$$
$$f_x|_{(2,3)} = y^2 e^{xy}|_{(2,3)} = 9e^6$$
$$f_{xx}|_{(2,3)} = y^3 e^{xy}|_{(2,3)} = 27e^6$$
$$f_y|_{(2,3)} = e^{xy} + xye^{xy}|_{(2,3)} = 7e^6$$
$$f_{yy}|_{(2,3)} = \left(2xe^{xy} + x^2 ye^{xy}\right)|_{(2,3)} = 16e^6$$
$$f_{xy}|_{(2,3)} = \left(2ye^{xy} + xy^2 e^{xy}\right)|_{(2,3)} = 24e^6 \ .$$

So then

$$f(x,y) = e^6\left[3 + 9(x-2) + 7(y-3)\right.$$
$$+ \qquad\qquad \frac{1}{2}\left(27(x-2)^2 + 48(x-2)(y-3) + 16(y-3)^2\right) \ .$$

**Suggested questions:** Q10-12.

## 3.3 Critical points

### Recap of 1-dimensional case

The critical points of a function $f : \mathbb{R} \to \mathbb{R}$ are all the points with $f_x = 0$. If $f_{xx} > 0$ it is a local minimum. If $f_{xx} < 0$ it is a local maximum. If $f_{xx} = 0$ (e.g. $f = x^4$ at $x = 0$) more analysis is needed.

### 3.3.1 2-dimensional case

We wish to generalise this to find critical points, and say whether they are local maxima, minima, or saddle-points in 2 or more dimensions. A critical point is a point at which both $f_x = 0$ and $f_y = 0$. Or equivalently, $\nabla f = \mathbf{0}$.

**Examples**

1: $f(x,y) = x^2 + y^2$.
The graph of $z = f(x,y)$ is a parabolic cylinder.

$$f_x = 2x \ ; \ f_y = 2y$$

The point (0,0) is the only critical point. Before getting into fancy definitions, it already seems clear that this critical point is a minimum.
2: $f(x,y) = x^2 - y^2$.
The graph of $z = f(x,y)$ is a saddle – it is a critical point, but it is neither a maximum, nor a minimum!

$$f_x = 2x \ ; \ f_y = -2y$$

The point (0,0) is the only critical point. It is quite clear from the functional form that it increases away from origin for fixed $y$ but decreases for fixed $y$.

**Distinguishing local maxima and minima using the Taylor expansion:**

Let us now be somewhat more formal.

> **Definition**
>
> - A point $(a, b)$ is said to be a local *maximum* if $f(a, b) > f(x, y)$ for all points $(x, y)$ in a sufficiently small neighbourhood surrounding $(a, b)$.
>
> - A point $(a, b)$ is said to be a local *minimum* if $f(a, b) < f(x, y)$ for all points $(x, y)$ in a sufficiently small neighbourhood surrounding $(a, b)$.
>
> Critical points where neither of the two above criteria are true – i.e. a critical point that is neither a maximum nor a minimum – are called "saddle points", based on the intuition above.

We can use the Taylor expansion about $(a, b)$ to tell us about the nature of the point there. To simplify things call $h = x - a$ and $k = y - b$ and call

$$P = f_{xx}(a, b)$$
$$Q = f_{xy}(a, b)$$
$$R = f_{yy}(a, b).$$

Then using the Taylor expansion, we can write

$$f(x, y) = f(a, b) + h f_x(a, b) + k f_y(a, b) + \frac{1}{2}(h^2 P + 2hkQ + k^2 R) + \ldots$$

*A necessary condition for a local maximum, local minimum or saddle point is that $f_x = f_y = 0$.*

The test for what sort of critical point it is:

> **Key idea**
>
> Let $M = PR - Q^2 = f_{xx}(a, b) f_{yy}(a, b) - f_{xy}(a, b)^2$.
>
> - If $M > 0$ and $P = f_{xx}(a, b) > 0$ then we have a local minimum.
>
> - If $M > 0$ and $P = f_{xx}(a, b) < 0$ then we have a local maximum.
>
> - If $M < 0$ then we have a saddle point.
>
> - If $M = 0$ then the test is inconclusive.

*Proof:* From the Taylor expansion, the value of the function near $(a, b)$ can be approximated by a quadratic polynomial (whose linear term vanishes because $(a, b)$ is a critical point):

$$f(a + h, b + k) \approx f(a, b) + \frac{1}{2P}(h^2 P^2 + 2hkQP + k^2 RP)$$

$$= f(a, b) + \frac{1}{2P}((hP + kQ)^2 + k^2 M).$$

If $P > 0$ and $M > 0$ then $f(a + h, b + k) - f(a, b) > 0$ and $f(a, b)$ is a minimum. If $P < 0$ then the reverse is true. If $M < 0$ then for some values of $h, k$ $f(a + h, b + k) - f(a, b)$ is positive and for others it's negative; thus we have a saddle point.

Find and classify the critical point(s) of $f(x,y) = x^2 + y^2$.
**Answer:**
The graph of $z = f(x,y)$ is a bowl, with a minimum at $x = y = 0$. The partial derivatives are

$$f_x = 2x, f_y = 2y,$$

so

$$f_{xx} = 2 \ f_{yy} = 2 \ f_{xy} = 0.$$

Now $M = 2 \times 2 - 0 = 4$ with $P = 2$ - so we do have a minimum!

Find and classify the critical point(s) of $f(x,y) = x^2 - y^2$
**Answer:**
Looking at the graph, we should find that $(0,0)$ is a saddle. We have

$$f_{xx} = 2 \ f_{yy} = -2 \ f_{xy} = 0,$$

and then $M = -4$ - as we'd expect!

Find and classify the critical point(s) of $f(x,y) = x^2 - y^2 + y^4 + x^2 y^2$.
**Answer:**
First find where $f_x = f_y = 0$.
$$f_x = 2x(1 + y^2)$$
$$f_y = 2y(x^2 + 2y^2 - 1)$$
$$f_{xx} = 2(1 + y^2)$$
$$f_{yy} = 2(x^2 + 6y^2 - 1)$$
$$f_{xy} = 4xy$$

Solving $f_x = f_y = 0$ gives $x = 0$ and $y = 0, \pm\frac{1}{\sqrt{2}}$ (since $y^2 + 1 \geq 1$ then $f_x$ can only give $x = 0$) so have

$$(x,y) = (0,0) \ or \ (0, \frac{1}{\sqrt{2}}) \ or \ (0, -\frac{1}{\sqrt{2}})$$

which I'll label $A, B$ and $C$.

- At $A$ we have $P = 2$, $R = -2$, $Q = 0$ so $M = -4$ and $A$ is a saddle

- At $B$ and $C$ we have $P = 3$, $R = 4$, $Q = 0$ so $M = 12 > 0$. Also $P > 0$ so that $B, C$ are minima.

Investigate the critical point(s) of $f(x,y) = x^2 + y^2 + xy - x + y$.
**Answer:**

First we find where $f_x = f_y = 0$:

$$f_x = 2x + y - 1$$
$$f_y = 2y + x + 1.$$

We must have $y = 1 - 2x$, so $2 - 3x + 1 = 3 - 3x = 0$; $x = 1$ and $y = -1$.
At the point $(1, -1)$ we have

$$M = f_{xx}f_{yy} - f_{xy}^2$$
$$= 2 \times 2 - 1^2 = 3 > 0.$$

Since $P > 0$ we have a local minimum.

---

**Try it out**

You are a box manufacturer. You need to make a rectangular box *open at the top* with a volume of $32m^3$. What are the dimensions in order to make the surface area as small as possible?
**Answer:**
First write the expressions for the volume and surface area if the base width height are $x, y, z$;

$$V(x, y, z) = xyz$$
$$\tilde{S}(x, y, z) = xy + 2xz + 2yz.$$

Given $V(x, y, z) = 32$, we determine that $z = \frac{32}{xy}$. Under this constraint, the surface area $\tilde{S}$ can be rewritten as a function of two variables

$$S(x, y) = \tilde{S}(x, y, \tfrac{32}{z}) = xy + \frac{64}{y} + \frac{64}{x}.$$

The extrema of $S(x, y)$ occur where $S_x = S_y = 0$: that is, where

$$0 = S_x = y - \frac{64}{x^2},$$
$$0 = S_y = x - \frac{64}{y^2}.$$

Solving this gives $x = y$ and then $x^3 = 64$. Hence, $x = y = (64)^{\frac{1}{3}} = 4m$ and $z = \frac{1}{2}(2V)^{\frac{1}{3}} = 2m$, which yield a surface area $S(4, 4) = 48m^2$.
Our goal was to find the minimal surface area, so we need to look at the second derivatives of $S(x, y)$. We have $S_{xx} = \frac{2(64)}{x^3}$, $S_{yy} = \frac{2(64)}{y^3}$ and $S_{xy} = 1$, so $P = S_{xx}(4, 4) = 2$, $Q = S_{xy}(4, 4) = 1$ and $R = S_{yy}(4, 4) = 2$. Therefore, $M = PR - Q^2 = 3 > 0$ and $P > 0$, so the surface area $S(4, 4) = 48m^2$ is indeed the minimum.

---

### 3.3.2 $n$-dimensional case

**Not examined: just for completeness**.

We can generalise this to any number of dimensions. Let $f : \mathbb{R}^n \to \mathbb{R}$. Critical points are given by $\nabla f(a) = \mathbf{0}$. To determine their nature, define the ***Hessian:***

$$H_{ij} = \frac{\partial^2 f}{\partial x_i \, \partial x_j}.$$

The Taylor expansion in $n$ dimensions near that point $\mathbf{x} = \mathbf{a}$ is given by the following expression, which is a generalization to multiple variables of Equation 3.4:

$$f(\mathbf{x}) = \sum_{i=1}^{n}(x_i - a_i)\frac{\partial f}{\partial x^i}(\mathbf{a}) + \frac{1}{2}\sum_{i,j=1}^{n} H_{ij}(\mathbf{a})(x_i - a_i)(x_j - a_j)$$

At a critical point the first term vanishes. So to figure out what happens we need to understand the term with the $H_{ij}$. We now need to know a little bit of linear algebra – consider the $n$ eigenvalues of $H_{ij}$. If they are all positive (i.e. $H_{ij}$ is positive definite) then it is a local minimum. All negative (i.e. $H_{ij}$ is negative definite) it is a maximum. If there are both positive and negative eigenvalues it is a saddle. Note that in two dimensions we have $H_{ij} = \begin{pmatrix} P & Q \\ Q & R \end{pmatrix}$, and positive or negative definiteness is guaranteed by $M = \det H > 0$, giving precisely the criteria specified above.

**Suggested questions:** Q8, Q13.