

# Essential Mathematics for Neuroscience

Fabian Sinz, Jakob Macke & Philipp Lies



October 21, 2009

# Contents

<b>1</b>	<b>Basics</b>	<b>2</b>
1.1	Essential Functions . . . . .	3
1.1.1	Polynomials and Powers . . . . .	4
1.1.2	Linear Functions . . . . .	5
1.1.3	Trigonometric Functions . . . . .	6
1.1.3.1	The geometric view . . . . .	7
1.1.3.2	The Periodic Signal View . . . . .	10
1.1.4	The $e$ -function and the Logarithm . . . . .	13
1.1.4.1	The Exponential Function . . . . .	13
1.1.4.2	Logarithms . . . . .	15
1.1.5	Lines (Affine Functions) . . . . .	18
1.1.6	Piecewise Defined Functions . . . . .	18
1.1.7	Sketching Functions . . . . .	19
1.1.7.1	Adapting Functions . . . . .	19
1.1.7.2	Compositions of Functions . . . . .	21
1.2	Basic Calculus . . . . .	23
1.2.1	Derivatives . . . . .	23
1.2.2	Higher-Order Derivatives . . . . .	33
1.2.3	Finding Maxima/Minima of a Function . . . . .	34
1.2.4	Approximating Functions Locally by Lines and Polynomials	39
<b>2</b>	<b>Appendix</b>	<b>47</b>
2.1	Notation and Symbols . . . . .	48

# Chapter 1

## Basics

This chapter serves as an introduction to a few basic elements that will be needed throughout the course. We begin by reviewing basic families of functions like *linear functions*, *polynomials* and *trigonometric functions*, as well as some of their properties. Afterwards we will look at some elementary calculus on those types of functions.

## 1.1 Essential Functions

A function is a rule that relates to sets of quantities, the *inputs* and the *outputs*. Each input  $x$  is deterministically related to an output  $f(x)$ . For example,  $f(x)$  might temperature on day  $x$ , or the firing rate of a neuron in response to a stimulus  $x$ . Thus, functions can be used as mathematical models of processes in which one quantity is transformed into another in a deterministic way. Even when the process of transformation is not deterministic, usually an underlying deterministic process corrupted by random noise can be used. In the example above, the firing rate of the neuron could be  $f(x) + \varrho$ , where  $\varrho$  is a noise-term. In contrast to a deterministic function,  $f(x) + \varrho$  denotes a whole set of values for a given  $x$  since the random term  $\varrho$  can take different values for each trial. Therefore,  $f(x) + \varrho$  is not a function in the strict sense. The reason is that, functions—by definition—are rules how to assign elements  $x$  of one set to *unique* elements  $f(x)$  of another set. Only if the target elements are unique, the assignment rule is called *function*. When defining a function, we have to specify the two *sets* between the function is mapping and the *rule* that transforms an element of the target set to an element of the input set. For example, if we want to define a function  $f$  that is transforming elements of a set  $A$  into elements of a set  $B$  according to the rule  $r$ , we would write this as

$$\begin{aligned} f: \quad A &\rightarrow B \\ a &\mapsto r(a). \end{aligned}$$

Here,  $a$  is an element of  $A$  (written  $a \in A$ ) and  $r(b)$  is an element of  $B$  (i.e.  $r(b) \in B$ ). The set  $A$  is usually called *domain of  $f$*  while  $B$  is called *the co-domain of  $f$* . The arrow “ $\rightarrow$ ” is used to denote the mapping between the two sets, while “ $\mapsto$ ” denotes the mapping from an element of the domain to an specific element of the co-domain. This means that “ $\rightarrow$ ” tells us what kind of objects are mapped into another and “ $\mapsto$ ” specifies the assignment rule.

The rule  $r$  can be anything that can be done with elements of  $A$ . For example, if the function  $f$  simply doubles any real number, we would write

$$\begin{aligned} f: \quad \mathbb{R} &\rightarrow \mathbb{R} \\ x &\mapsto 2 \cdot x \quad x \in \mathbb{R}. \end{aligned}$$

In most cases, the inputs and outputs of function will be numbers, but this does not necessarily have to be the case (i.e. the elements of the domain  $A$  and the co-domain  $B$  do not need to be numbers).

Although, in principle, there are infinitely many functions on the real numbers, knowing only a few of them is usually enough to get along well in most natural sciences. The reason is that most complex functions are built by adding, multiplying, or composing simpler ones. It is important that you get comfortable with those simpler functions since they are your toolbox to understand and build more complex functions. Once you have an intuition how those simple functions behave, it is often not too difficult to get a feeling for a more complicated one. In this section we will review the most important simple functions and present their most important properties.

### 1.1.1 Polynomials and Powers

Polynomials is a very common class of functions. The two most widely known kinds of polynomials are the parabola  $f(x) = x^2$  and the more general quadratic function  $f(x) = ax^2 + bx + c$ . In general, polynomials consist of a sum of positive integer powers  $k$  of  $x$  with coefficients  $a_k$ :

$$f(x) = a_n x^n + \dots + a_1 x + a_0.$$

The single terms in the sum are called *monomials*. The *degree* of the polynomial is the largest exponent of its monomials. The polynomial above has a degree of  $n$ . Polynomials have nice properties like e.g. the *derivatives* and *anti-derivatives* of polynomials are easy to calculate and yield polynomials again. One frequent use of polynomials is to approximate any function at a certain location. This approximation is called *Taylor-Expansion*. We will discuss the Taylor-Expansion and many properties of polynomials in later chapters.

This is a good point to introduce the notation for sums over several elements: Instead of indicating the entire sum by three dots “...” we use the greek uppercase letter sigma  $\Sigma$  (like sum) to indicate a sum over all terms directly after the sigma. These terms are usually indexed and the range of the index is written below and above the  $\Sigma$ . Since  $x^0 = 1$  for all  $x \in \mathbb{R}$  we write the polynomial from above as

$$\begin{aligned} f(x) &= a_n x^n + \dots + a_1 x + a_0 \\ &= \sum_{k=0}^n a_k x^k. \end{aligned}$$

While polynomials have exponents  $k \in \mathbb{N}_0$  (where  $\mathbb{N}_0$  denotes the set of natural number including 0), exponents can in principle be in  $\mathbb{R}$  as well. There are two most important cases: when the exponent is negative and when it is a rational number (i.e. a number that can be written as a fraction). A negative exponent of a number is merely a shortcut for  $x^{-a} = \frac{1}{x^a}$ . In many cases, for example when calculating derivatives, the notation with negative exponent is useful. A fraction in the exponent is another way of writing roots. For example the square root  $\sqrt{x}$  is equivalently written as  $x^{\frac{1}{2}}$ . In general, the  $n$ th root of  $x$  can be written as  $\sqrt[n]{x} = x^{\frac{1}{n}}$ .

We conclude this section by stating a few calculation rules for powers for  $x, a \in \mathbb{R}$ . You should know all of them by heart and be able to use them effortlessly.

**Calculation Rules for Powers**

The following rules apply to any  $x, a \in \mathbb{R}$ :

1. Anything to the power of zero is one:  $x^0 = 1$
2. Multiplying two terms with the same basis is equivalent to adding their exponents:  $x^a \cdot x^b = x^{a+b}$
3. Dividing two terms with the same basis is equivalent to subtracting their exponents:  $\frac{x^a}{x^b} = x^a \cdot x^{-b} = x^{a-b}$
4. Exponentiating a term is equivalent to multiplying its exponents:  $(x^a)^b = x^{a \cdot b}$
5. A special case of rule 3. is given by  $\frac{1}{x^a} = x^{-a}$
6. The  $a^{th}$  root of  $x$  is given by  $\sqrt[a]{x} = x^{\frac{1}{a}}$  for  $x \geq 0$ .

**1.1.2 Linear Functions**

Linear functions are among the simplest functions one can imagine. You can imagine a linear function as a line (plane, or hyperplane) through the origin. Algebraically, their key property is that the function value of a sum  $x + y$  of elements  $x, y$  equals the sum of their function values  $f(x) + f(y)$ . The same is true for multiples of input elements, i.e. the function value of some multiple  $a \cdot x$  of an element  $x$  from the domain is the multiple of the function value  $a \cdot f(x)$ . If any function fulfills these two properties, it is linear by definition.

**Definition (Linear Function)** A function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is said to be *linear* if it fulfills the following two properties:

$$f(x + y) = f(x) + f(y) \quad \text{for all } x, y \in \mathbb{R} \quad (1.1)$$

$$f(a \cdot x) = a \cdot f(x) \quad \text{for all } a, x \in \mathbb{R}. \quad (1.2)$$

◇

These properties have remarkable consequences. While for general functions, a single input-output pair of values  $(x, f(x))$  does not tell anything about the value of the function at other locations  $y \neq x$ , a single such pair is enough to know the value of a linear function at any location: Assume we are given the input-output pair  $(x, f(x))$  and we know that  $f$  is linear. In order to calculate the value of  $f$  at another location  $y$ , we search for a scalar  $a$  that scales  $x$  into  $y$ , i.e.  $y = a \cdot x$ . Clearly, this scalar is easily given by  $a = \frac{y}{x}$ . Once we know  $a$ , we can compute  $f(y)$  via

$$\begin{aligned} f(y) &= f(a \cdot x) \\ &= a \cdot f(x). \end{aligned}$$

**Example (Mathematical Modelling of Receptive Fields)** For some neurons, it is often assumed that their responses, i.e. the spike rate  $r(x)$ , depends linearly on the stimulus  $x$ .

Assume our cell responds to a visual image  $I_1$  with a spike rate of  $r_1 = 20$  spikes per second and to another image  $I_2$  with  $r_2 = 60$  spikes per second. What spike rate would we expect to the mean of  $I_1$  and  $I_2$ , if our neuron was truly linear in the stimuli? The answer is easy to calculate. Let  $r : \mathcal{I} \rightarrow \mathbb{R}$  denote the function from images (denoted by  $\mathcal{I}$ ) to spike rate. We already know  $r(I_1) = r_1 = 20 \frac{sp}{s}$  and  $r(I_2) = r_2 = 60 \frac{sp}{s}$ . Then the response to the mean of the two images is

$$\begin{aligned} r\left(\frac{1}{2}I_1 + \frac{1}{2}I_2\right) &= r\left(\frac{1}{2}I_1\right) + r\left(\frac{1}{2}I_2\right) \\ &= \frac{1}{2}r(I_1) + \frac{1}{2}r(I_2) \\ &= \frac{1}{2}r_1 + \frac{1}{2}r_2 \\ &= 10 \frac{sp}{s} + 30 \frac{sp}{s} \\ &= 40 \frac{sp}{s} \end{aligned}$$

This property does not only hold for two input stimuli. It holds for an arbitrary number of stimuli. If the rate function  $r$  realized of our neuron is linear, then response to the mean of  $n$  images is just the mean response to the single images.

$$\begin{aligned} r\left(\frac{1}{n} \sum_{k=1}^n I_k\right) &= \frac{1}{n} \sum_{k=1}^n r(I_k) \\ &= \frac{1}{n} \sum_{k=1}^n r_k \end{aligned}$$

**Question:**

?

Of course, real neurons are not truly linear. If a neuron was indeed linear, for inputs, this would lead to some very unrealistic conclusions. Name two of them!

**Answer:**

For some stimuli, the spike rate would be negative. Also, for stimuli with very high input, the spike-rate would be arbitrarily large, i.e. the neuron's rate would not saturate.

<

### 1.1.3 Trigonometric Functions

Trigonometric functions are functions of an angle  $\vartheta$ . The most common trigonometric functions are  $\sin(\vartheta)$ ,  $\cos(\vartheta)$ ,  $\tan(\vartheta)$  and  $\cotan(\vartheta)$ .

In general, there are two natural ways to think about trigonometric functions: the geometrical view quantities and the periodic signal view.

### 1.1.3.1 The geometric view

In the geometric view,  $\cos(\vartheta)$  and  $\sin(\vartheta)$  represent the  $x$ -coordinate and the  $y$ -coordinate of a point on the intersection between a circle with radius one centered at the origin, and a line through the origin that encloses an angle of  $\vartheta$  with the  $x$ -axis. In this view,  $\tan(\vartheta)$  and  $\cotan(\vartheta)$  have a natural interpretation as well, namely the length of the line, touching the circle, between the upper leg of the angle and the  $x$ -axis or the  $y$ -axis, respectively. Alternatively,  $\tan(\theta)$  is the ratio between the  $x$ - and the  $y$ -coordinate (see Figure 1.1).

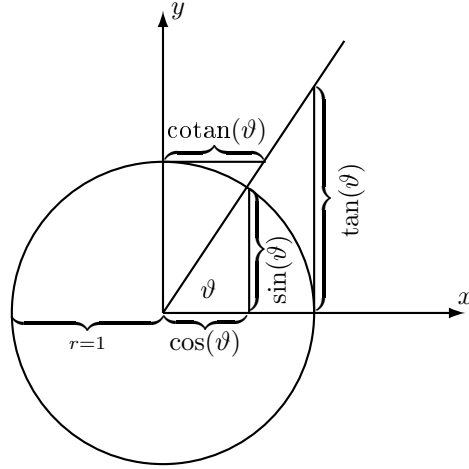
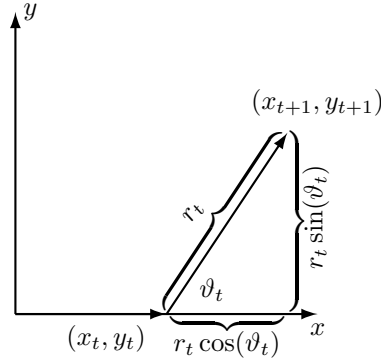


Figure 1.1: Geometrical view of  $\sin(\vartheta)$ ,  $\cos(\vartheta)$ ,  $\tan(\vartheta)$  and  $\cotan(\vartheta)$ . For a given angle  $\vartheta$ ,  $(\cos(\vartheta), \sin(\vartheta))$  are the coordinates of the point on the intersection between a circle of radius  $r = 1$  and the upper leg of the angle.  $\tan(\vartheta)$  is the length of the line between the  $x$ -axis and the upper leg of the angle, that "touches" the unit circle (therefore the name *tangens* from *lat. tangere = to touch*).  $\cotan(\vartheta)$  is defined analogously for the line touching the circle from above.

**Example (Path Integration)** The term *path integration* denotes the ability of moving organisms (such as ants) to remember the direction and length of the vector to its home while moving in the environment. We will now just look at a special case of updating the home vector in world coordinates, i.e. a global fixed coordinate system. This means that the home base is assigned the coordinates  $(0,0)$  and the moving organism stores its position according to a global coordinate frame.



Imagine you are an ant living in a completely flat world. For the sake of simplicity we further imagine that you have a compass and you know the length of your steps, so that you can measure the angles you turn and the distance you walked. Now imagine that you already explored your environment for a while and that you are now standing at position  $(x_t, y_t)$  looking along the  $x$ -axis. If you know turn by an angle of  $\vartheta_t$  and move into the new direction by a distance of  $r_t$ , what is your new position  $(x_{t+1}, y_{t+1})$ ?



Looking again at figure 1.1 and the figure above shows that the new direction in which you are walking is  $(\cos(\vartheta_t), \sin(\vartheta_t))$ . Since you are walking along that direction for a distance of  $r_t$ , the total displacement is  $r \cdot (\cos(\vartheta_t), \sin(\vartheta_t))$ . If we add this displacement to the old position, we get the new position

$$(x_{t+1}, y_{t+1}) = (x_t, y_t) + r_t \cdot (\cos(\vartheta_t), \sin(\vartheta_t))$$

◁

Due to the geometrical property of  $\sin(\vartheta)$  and  $\cos(\vartheta)$  we can derive a very useful equality by employing Pythagoras theorem.

**Theorem (Pythagoras)** For a right angle triangle with side lengths  $a, b$  and  $c$ , where  $c$  is the length of the longest leg, while the legs with length  $a$  and  $b$  enclose an angle of  $90^\circ$  (or  $\frac{\pi}{2}$  in radians), the following equality holds:

$$a^2 + b^2 = c^2$$

◇

In the case of  $\sin(\vartheta)$  and  $\cos(\vartheta)$ , we know the value of  $c$  for any given  $\vartheta$ . It is simply  $c = 1$ , since the point  $p = (\cos(\vartheta), \sin(\vartheta))$  lies on the circle of radius  $r = 1$ . Therefore, we know that

$$\cos(\vartheta)^2 + \sin(\vartheta)^2 = 1$$

for all angles  $\vartheta$ .

There is another useful equality that we can read of figure 1.1. Assume we want to know the scaling factor  $s$  that transforms  $\sin(\vartheta)$  into  $\tan(\vartheta)$ , i.e.  $\sin(\vartheta) \cdot s = \tan(\vartheta)$ . From looking at figure 1.1 we know that it is the same scaling factor that scales  $\cos(\vartheta)$  into 1, i.e.  $\cos(\vartheta) \cdot s = 1$ . In this case,  $s$  is easy to calculate: It is simply  $s = \frac{1}{\cos(\vartheta)}$ . Therefore we have found a formula how to calculate  $\tan(\vartheta)$  from  $\sin(\vartheta)$  and  $\cos(\vartheta)$ :

$$\tan(\vartheta) = \frac{\sin(\vartheta)}{\cos(\vartheta)}.$$

In an analogous manner we can also derive the equality

$$\cotan(\vartheta) = \frac{\cos(\vartheta)}{\sin(\vartheta)}.$$

Since  $\tan(\vartheta)$  and  $\cotan(\vartheta)$  can be written in terms of a quotient, the radius of the circle does not even have to be one. Assume we want to know the tangens between the  $x$ -axis and the leg  $(0, p)$  for a point  $p = (x, y)$  that lies on a circle with radius  $r$ . Since we can write  $p$  equivalently as  $p = (x, y) = r \cdot (\cos(\vartheta), \sin(\vartheta))$  for an appropriate value of  $r$ , we have that

$$\frac{y}{x} = \frac{r \cdot \sin(\vartheta)}{r \cdot \cos(\vartheta)} = \tan(\vartheta).$$

Therefore,  $\tan(\vartheta)$  is simply the quotient of the opposite leg and the adjacent leg of a right angle triangle. Similarly we can get

$$\frac{x}{y} = \frac{r \cdot \cos(\vartheta)}{r \cdot \sin(\vartheta)} = \cotan(\vartheta).$$

Further, sine and cosine can also be defined in terms of quotients for circles with an arbitrary radius. According to the intercept theorem the ratio of  $\cos(\vartheta)$  to 1 is equal to the ratio of  $x$  to  $r$  for every point  $p = (x, y)$  on a circle with radius  $r$ . Equally, the ratio of  $\sin(\vartheta)$  to 1 is equal to the ratio of  $y$  to  $r$ . Therefore, we get

$$\begin{aligned} \sin(\vartheta) &= \frac{y}{r} \\ \cos(\vartheta) &= \frac{x}{r} \end{aligned}$$

Other useful properties that can also be read off from the geometric view. Among them are the symmetry properties of  $\sin(\vartheta)$ ,  $\cos(\vartheta)$ ,  $\tan(\vartheta)$  and  $\cotan(\vartheta)$ :

$$\begin{aligned} \cos(-\vartheta) &= \cos(\vartheta) \\ \sin(-\vartheta) &= -\sin(\vartheta) \\ \tan(-\vartheta) &= -\tan(\vartheta) \\ \cotan(-\vartheta) &= -\cotan(\vartheta) \end{aligned}$$

Sometimes we do not want to compute the sine or cosine of an angle but the angle itself. For example, we take point  $p = (x, y)$  and want to know what

the angle between the x-axis and the line from the origin to  $p$  is. We know that  $\tan(\vartheta) = \frac{y}{x}$ , but what is  $\vartheta$ ? To solve for  $\vartheta$ , we need the inverse trigonometric functions. For every trigonometric function there is an inverse trigonometric function:  $\arccos(x)$ ,  $\arcsin(x)$ ,  $\arctan(x)$ ,  $\text{arccotan}(x)$ . In some texts they are confusingly written as  $\cos^{-1}(\vartheta)$  whereas  $\cos(\vartheta)^{-1}$  means  $1/\cos(\vartheta)$ , so you should always use the arc-notation. Now, by using the  $\arctan$  function we can compute the angle between  $p$  and the x-axis as  $\vartheta = \arctan(\frac{y}{x})$ .

**Remark** There are two units for measuring angles which are commonly used and get quite often mixed up: *radians* and *degrees*. A full circle of  $360^\circ$  corresponds to  $2\pi$  radians. In computer programs, the default unit is radians. Therefore to compute the sine of  $45^\circ$  one has to compute  $\sin(\pi/4)$ . Likewise if you computed an angle by using  $\arccos(x)$  your output is in radians. To convert from radians to degrees, you have to multiply your result by  $\frac{180}{\pi}$  and similarly converting from degrees to radians by multiplying with  $\frac{\pi}{180}$ .

&lt;

### 1.1.3.2 The Periodic Signal View

The periodic signal view of thinking about  $\sin(\vartheta)$  and  $\cos(\vartheta)$  is especially useful when dealing with signals such as e.g. membrane potentials of neurons or stripes of natural images. It is also closely related to techniques such as Fourier- or Spectral Analysis, which are indispensable tools for the analysis of neurophysiological signals. In order to get from the geometric to the periodic signal view, just imagine a point on the unit circle that is moving with constant speed counterclockwise along the circle. If we plot the time  $t$  against the point's  $x$ -coordinate  $x(t) = \cos(\vartheta(t))$ , we get a strongly periodic function. Same applies to the  $y$ -coordinate  $y(t) = \sin(\vartheta(t))$ . Figure 1.2 shows the graphs of the two functions. For the moment we wrote  $\vartheta$  as a function of time in order to be able to say that the point is moving with constant speed. The faster the point is moving along the circle the more cycles we can get in one fixed time interval. This is expressed by the frequency  $\omega$  of the sine or cosine, respectively. It tells us how many cycles our point does in one unit time interval. We can therefore equivalently write  $x(t) = \cos(\omega \cdot t)$  and  $y(t) = \sin(\omega \cdot t)$ . From now on we will drop the dependence on time. The frequency then tells us how many cycles of our point fit in the interval  $[0, 2\pi]$ .

**Example (Sine and Cosine Gratings)** Assume that you want measure the orientation selectivity a V1-cell you are recording from with an electrode. A simple experiment would be to present gratings with different orientations and varying amount of bars in one unit interval, i.e. different spatial frequency. The most common way of producing such patterns is to use *sine* and *cosine gratings*. Figure 1.3 shows such a grating.

Since an image can be seen as a function, that assigns a graylevel value to each pixel ( $f(x)$  is the gray value of pixel  $x$ ) those gratings are simply sine or

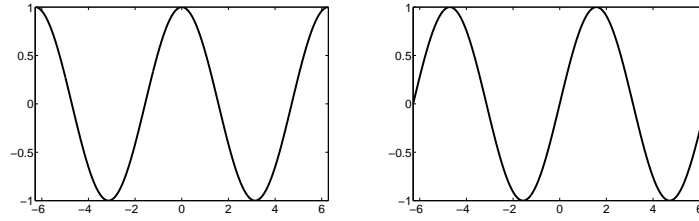


Figure 1.2:  $\cos(\omega t)$  (left) and  $\sin(\omega t)$  (right) in the interval of  $t \in [-2\pi, 2\pi]$  with a frequency of  $\omega = 1$ .

cosine functions over  $\mathbb{R}^2$ . In order to get to the graylevel values at the different pixels, the function is discretized, i.e. is only evaluated at certain locations that correspond to the pixels. At the moment we shall only look at how to produce vertical or horizontal gratings. We will see how to produce gratings of arbitrary orientation later.

For a vertical grating, we know that the graylevel value along the vertical axis must be constant. If  $I(x, y)$  denotes the function that represents the image, i.e. the functions that assigns a graylevel value to a position  $(x, y)$ , we can produce a vertical grating by  $I(x, y) = \sin(\omega x)$ . Instead of using the sine we could as well have used the cosine function. A horizontal grating can be generated in exactly the same way by replacing  $x$  by  $y$  inside the sine function. Here is the matlab code to produce a horizontal grating of frequency  $\omega = 2$ :

```
>> [X,Y] = meshgrid([-2*pi:0.01:2*pi]); % get the sample points
>> omega = 2; % set the frequency
>> imagesc(sin(omega*X)) % display grating
>> colormap(gray) % set colormap to gray values
>> axis off % switch off the axes
```

At the moment all our gratings have a fixed contrast, since  $-1 \leq \sin(x), \cos(x) \leq 1$  for all  $x \in \mathbb{R}$ . However, we can vary the contrast by varying the amplitude of the sine. This is done by premultiplying an appropriate scaling factor  $I(x, y) = A \cdot \sin(\omega x)$ . In order to build that into the matlab code above, you must specify the maximum and the minimum gray value when calling the function `imagesc`, since it automatically scales the gray values otherwise

```
>> [X,Y] = meshgrid([-2*pi:0.01:2*pi]); % get the sample points
>> omega = 2; % set the frequency
>> A = 3; % set the contrast
>> maxA = 10; % set maximal contrast
>> imagesc(A*sin(omega*X), [-maxA,maxA]) % display grating
>> colormap(gray) % set colormap to gray values
>> axis off % switch off the axes
```

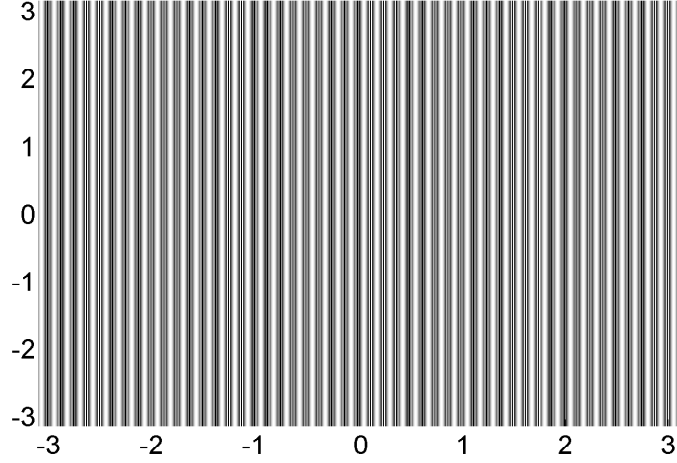


Figure 1.3: Example of a vertical sine grating with a spatial frequency of approx. 8Hz along the  $x$ -axis. The graylevel value at each position  $(x, y)$  is given by  $I(x, y) = \sin(50 \cdot x)$ .

Up to now we saw how to control the frequency and the amplitude of sine and cosine functions. In order to complete this subsection we will also see how to shift the sine and the cosine functions along the  $x$ -axis. Let us look at a sine function with a certain frequency  $\sin(\omega t)$ . At  $t = 0$  also  $\sin(\omega t) = 0$ . If we want to shift the sine function along the  $x$ -axis by an offset of  $\phi$ , we must ensure that the shifted version is zero at  $t = \phi$  and not at  $t = 0$ . However, this will be the case, if we subtract  $\phi$  from the argument of the sine. Therefore, a sine that is shifted by an offset of  $\phi$  along the  $x$ -axis is given by  $\sin(\lambda t - \phi)$ . Cosine functions are shifted in exactly the same way. This offset  $\phi$  is called *phase* of the signal. Now we are able to write down the general form of a sine or cosine signal with a given amplitude  $A$  and phase  $\phi$ : It is

$$A \cdot \sin(\omega t - \phi) \quad \text{and} \quad A \cdot \cos(\omega t - \phi).$$

Before finishing this section about trigonometric functions we just want to mention two equalities that are useful when calculating with sine and cosine. These equalities are called the *Addition Theorems*.

**Theorem (Addition Theorems)** The following equalities hold for all  $x, y \in \mathbb{R}$ :

$$\begin{aligned} \cos(x + y) &= \cos(x) \cos(y) - \sin(x) \sin(y) \\ \sin(x + y) &= \sin(x) \cos(y) + \cos(x) \sin(y) \end{aligned}$$

◇

Using the symmetry properties of  $\sin(x)$  and  $\cos(x)$  one can also derive similar expressions for  $\cos(x - y)$  and  $\sin(x - y)$ .

**Exercise**

E

Write down the the corresponding expressions for  $\cos(x - y)$  and  $\sin(x - y)$  .

◁

**Important Rules for Trigonometric Functions**

The following rules apply to any  $x, y, \vartheta \in \mathbb{R}$ :

- Pythagoras's Theorem:  $\cos(\vartheta)^2 + \sin(\vartheta)^2 = 1$
- Symmetry Properties:

$$\begin{aligned}\cos(-\vartheta) &= \cos(\vartheta) \\ \sin(-\vartheta) &= -\sin(\vartheta) \\ \tan(-\vartheta) &= -\tan(\vartheta) \\ \cotan(-\vartheta) &= -\cotan(\vartheta)\end{aligned}$$

- Addition Theorems

$$\begin{aligned}\cos(x + y) &= \cos(x)\cos(y) - \sin(x)\sin(y) \\ \sin(x + y) &= \sin(x)\cos(y) + \cos(x)\sin(y)\end{aligned}$$

**1.1.4 The  $e$ -function and the Logarithm****1.1.4.1 The Exponential Function**

The exponential or  $e$ -function  $f(x) = e^x = \exp(x)$  is one of the most frequently occurring and important function in the everyday life of a natural scientist. The number denoted by “ $e$ ” is an irrational number called *Euler's number*. Its first digits are  $e \approx 2.7183$ . You should remember the approximate value of its inverse  $\frac{1}{e} \approx 0.37$  because it is used to define time constants of neural signal transduction.

The exponential function appears in many probability distribution in statistics, in solution of differential equations and will appear in many mathematical model of neural processes. The exponential function has some very nice properties. One of them is that it is its own derivative  $f'(x) = (e^x)' = e^x$ .

The following examples show three cases, in which the exponential function naturally occurs.

**Examples**

1. One central probability density in statistics is the *Normal* or *Gaussian Distribution*  $\mathcal{N}(\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ . In many experiments that involve

noisy measurements, the noise is assumed to be Gaussian with mean  $\mu = 0$ . One possible justification for this assumption is that sums of random quantities with other probability distributions tend to be Gaussian if the total number of this quantities increases. Since we think of the noise in an experiment as the superposition of many other processes that we are not interested in, the assumption that there are a large number of them which are linearly superimposed motivates the Gaussian noise assumption. Apart from that, the Gaussian distribution is frequently used because it has a lot of properties that make it possible to calculate analytical solutions of the respective statistical problems.

2. The potential change  $\Delta V_m(t)$  over time at a passive neuron membrane after applying a rectangular current pulse can be described by the following equation

$$\Delta V_m(t) = I_m R \left(1 - e^{-\frac{t}{\tau}}\right).$$

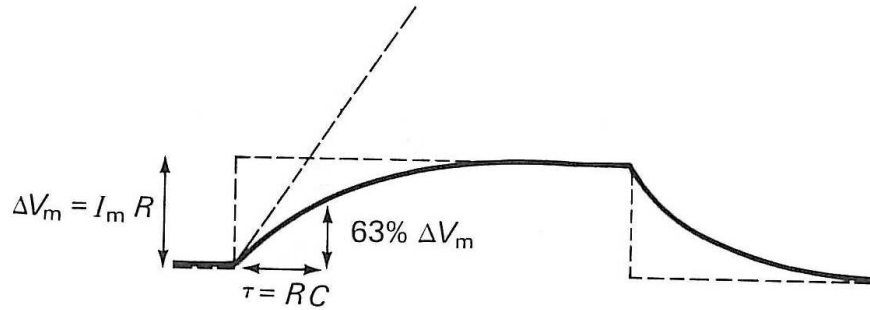


Figure 1.4: Potential change  $\Delta V_m(t)$  (solid line) over time at a passive neuron membrane after applying a rectangular current pulse (dashed line) [Figure from [1]].

Figure 1.4 shows the time course of the potential. Here  $I_m$  is the current, that has been injected,  $R$  is the membrane resistance and  $\tau$  is the *time constant* of the membrane. It tells us how fast the membrane potential follows the rectangular pulse. The greater  $\tau$  is, the longer it takes until  $\Delta V_m(t) = I_m R$ , where  $I_m R$  is the potential change induced by injecting the current  $I_m$ . In some sense  $\tau$  defines a time scale for that membrane.  $\tau$  is the time needed until the potential change reaches  $0.63 \cdot I_m R = (1 - 0.37) \cdot I_m R = (1 - \frac{1}{e}) I_m R$ , i.e. 63% of the potential change induced by the rectangular pulse. One can show that  $\tau = RC$ , where  $C$  is the membrane capacitance, i.e. its ability to buffer charge.

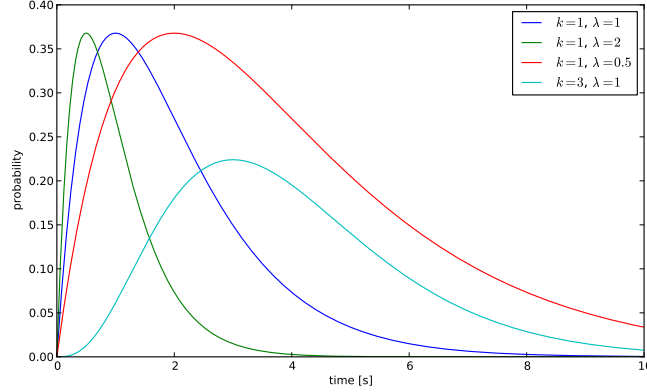


Figure 1.5: Poisson Process with 4 different sets of parameters.

3. A simple stochastic model for spike generation by a neuron is the *Poisson Process*. In this model, time is divided into a large number of bins. In each bin a spike occurs with probability  $p$  independent of whether a spike occurred in the bin before or not. If  $p$  is sufficiently small and the average spiking rate is  $\lambda$ , the probability of observing exactly  $k$  spikes in a time window of length  $\Delta t$  is given by the distribution of the Poisson Process with rate  $\lambda$ :

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

Figure 1.5 gives an example how a Poisson Process looks like for different sets of parameters. The symbol expression " $k!$ " denotes the *factorial function*  $k! = k \cdot (k-1) \cdot \dots \cdot 2 \cdot 1$ . We can use a notation similar to the  $\Sigma$  for sums to denote a product of several components: The uppercase Greek letter  $\Pi$  (for product) denotes a product of all elements following it. The indices are written in the same way as for sums (see also Appendix 2.1). Therefore, we can write the factorial function as  $k! = \prod_{n=1}^k n$ .

4. If a spike train is generated by a *Poisson Process*, then the distribution of the *inter-spike-intervals (ISIs)* is an *exponential distribution*. That is, if we observe a spike at time 0, then the probability (density) of observing the next spike at time  $s$  is given by  $p(s) = \mu e^{-\mu s}$ , for  $\mu = 1/\lambda$  in the example above.

&lt;

#### 1.1.4.2 Logarithms

Logarithms and their derivatives often occur in statistics. Estimating the parameters of a statistical model is often done via maximum likelihood estimation.



This involves taking the derivative of the log of the likelihood function.

Here we introduce logarithms. More advanced examples will be discussed in later chapters. For now, we start with a small example:

**Example** In this example, we look again at the time course of the membrane potential after the application of a rectangular current pulse  $\Delta V_m(t) = I_m R \left(1 - e^{-\frac{t}{\tau}}\right)$ . Assume that we want to measure the time constant of a certain membrane. For this purpose we excited the membrane with a rectangular current pulse and measured  $\Delta V_m(t)$  at several points  $t_k$ ,  $k = 1, \dots, n$  in time. Now we want to solve  $\Delta V_m(t) = I_m R \left(1 - e^{-\frac{t}{\tau}}\right)$  for  $\tau$  at each time step  $t_k$  and get  $n$  values  $\tau_k$  that we average to get your final estimation  $\hat{\tau} = \frac{1}{n} \sum_{k=1}^n \tau_k$  of the membrane's time constant. How do we solve for  $\tau_k$ ? The first step is easy: You rearrange the the terms to get

$$\begin{aligned} \Delta V_m(t) = I_m R \left(1 - e^{-\frac{t}{\tau_k}}\right) &\Leftrightarrow I_m R - \Delta V_m(t_k) = I_m R \cdot e^{-\frac{t_k}{\tau_k}} \\ &\Leftrightarrow \frac{I_m R - \Delta V_m(t_k)}{I_m R} = e^{-\frac{t_k}{\tau_k}}. \end{aligned}$$

Now we somehow have to extract  $-\frac{t_k}{\tau_k}$  from the exponent. This means that you are searching for a number  $a$ , such that  $e^a = \frac{I_m R - \Delta V_m(t_k)}{I_m R}$ . This is exactly the definition of the *natural logarithm*  $\ln(x)$ : It is the number that you have to put in the exponent of  $e$  in order to obtain  $x$ . Now you can solve for  $\tau$ :

$$\begin{aligned} \frac{I_m R - \Delta V_m(t_k)}{I_m R} = e^{-\frac{t_k}{\tau_k}} &\Leftrightarrow \ln \left( \frac{I_m R - \Delta V_m(t_k)}{I_m R} \right) = -\frac{t_k}{\tau_k} \\ &\Leftrightarrow -\frac{\ln \left( \frac{I_m R - \Delta V_m(t_k)}{I_m R} \right)}{t_k} = \frac{1}{\tau_k} \\ &\Leftrightarrow -\frac{t_k}{\ln \left( \frac{I_m R - \Delta V_m(t_k)}{I_m R} \right)} = \tau_k \end{aligned}$$

◁

We just saw that the natural logarithm is the function that cancels the exponential function, i.e.  $\ln(e^x) = e^{\ln(x)} = x$ . In general, a function  $g$  that cancels another function  $f$  is called *inverse function* of  $f$  and is denoted by  $g = f^{-1}$ . Not all functions have inverses. Some of them only have an inverse on a restricted range. We will discuss inverse function in more detail in the chapter about analysis.

So far we have seen how to solve  $e^x$  for  $x$  by using the natural logarithm. What if we want to solve an equation like  $2^x$  for  $x$ ? Here we cannot use the natural logarithm since  $\ln(x)$  is only the inverse function of  $e$ , not 2. Here we must use a logarithm that fits to 2. This logarithm is denoted by  $\log(x)$  and has the property that  $2^{\log(x)} = \log(2^x) = x$ . In general there is a logarithm for

any number  $b$  that is the inverse of the function  $f(x) = b^x$ . The number  $b$  is called *base* of the logarithm. If the base is not  $b = 2$  or  $b = e$ , we indicate the base in the subscript of  $\log_b(x)$ . However, the only frequently used logarithms are  $\log(x) = \log_2(x)$  and  $\ln(x) = \log_e(x)$ .

**Remark** It happens quite often that the base of the logarithm is not specified, either because it does not matter or it is clear from the context. In this case, the notation  $\log(x)$  is usually used. Usually this means that the base is  $e$  (e.g. Physics), sometimes it means that the base is 2 (e.g. in Information Theory) and sometimes it is used for base 10 (e.g. Economics). In the remaining part of the script we simply use  $\log(x)$  to denote any logarithm. The base should always be clear from the context or it does not matter. If we want to emphasize a certain base we will write it in the subscript or use the explicit notation  $\ln(x) = \log_e(x)$  or  $\lg(x) = \log_{10}(x)$ . !

&lt;

There is a neat trick how to calculate logarithms of arbitrary bases by using logarithms of another base:

$$\log_b(x) = \frac{\ln(x)}{\ln(b)} = \frac{\log(x)}{\log(b)}.$$

Until now we skipped an important detail of logarithms. When only dealing with real numbers, the logarithm is only defined on the strictly positive part of  $\mathbb{R}$ . We denote this set by  $\mathbb{R}^+$ . The reason for this restriction is easy to see. If we remember that  $\log_b x$  is the number that has to be put in the exponent of  $b$  in order to obtain  $x$ . If  $x$  is negative, there can generally be no such real number since  $b^x$  is positive. (The concept of logarithm can be extended to negative and imaginary numbers, but does lead to some complications, so we will not cover it here.)

We conclude this section with a few calculation rules. Most of them follow directly from the calculation rules of powers or the definition of the logarithm.

#### Calculation Rules for Logarithms

The following rules apply to any logarithm:

- $\log_b(b^x) = x$
- $\log_b(x \cdot y) = \log_b(x) + \log_b(y)$  for  $x, y \in \mathbb{R}^+$
- $\log_b\left(\frac{x}{y}\right) = \log_b(x) - \log_b(y)$  for  $x, y \in \mathbb{R}^+$
- $\log_b(x) = \frac{\ln(x)}{\ln(b)} = \frac{\log(x)}{\log(b)}$  for  $x, b \in \mathbb{R}^+$
- $\log_b(x^y) = y \cdot \log_b(x)$  for  $x \in \mathbb{R}^+, y \in \mathbb{R}$

### 1.1.5 Lines (Affine Functions)

Most people think about lines when they think about linear functions. However, lines are not generally linear functions. Only the lines that include the origin are strictly speaking linear functions. Lines versions of linear functions that are shifted along the  $y$ -axis. The general equation for a line is

$$f(x) = mx + t,$$

where  $m$  is called the *slope* of  $f$ . It is the first derivative or, equivalently, the amount about  $f$  changes if we increase  $x$  by one, i.e.  $m = f(x+1) - f(x)$ . The value of  $t$  determines the  $y$  coordinate of the point where  $f$  cuts through the  $y$ -axis. This can easily be seen by evaluating  $f$  at  $x = 0$ . Obviously, the function value  $f(0) = t$  must be the location on the  $y$ -axis where  $f$  hits it.

From the general form of lines we can also see why they are not strictly linear if  $t \neq 0$ . If they were linear,  $f(x+y) = f(x) + f(y)$  would have to hold for all  $x, y \in \mathbb{R}$ . However, it is easy to check that this is not the case:

$$\begin{aligned} f(x) + f(y) &= mx + t + my + t \\ &= m(x+y) + 2t \\ &\neq m(x+y) + t \\ &= f(x+y). \end{aligned}$$

Therefore, lines with  $t \neq 0$  are not linear. But we can always make them linear by subtracting  $t$ . This yields a line with the same slope  $m$ , which is shifted along the  $y$ -axis such that it cuts through  $(0,0)$ , which make it a truly linear function.

Functions of the form  $f(x) = mx + t$  are also called *affine functions*.

### 1.1.6 Piecewise Defined Functions

Sometimes, it is convenient to define a function by using two or more other functions. This is useful if we want to change the behaviour of the function on certain parts of the domain. Achieving a certain behaviour with a single expression might be difficult. It is then usually easier to use several expressions, one for each part of the domain. These functions are called *piecewise defined functions*. We just mention a few important examples here.

#### Examples

1. The Heaviside function is defined as:

$$f(x) = \begin{cases} 0 & x \leq 0 \\ 1 & x > 0 \end{cases}.$$

2. The absolute value is given by

$$f(x) = \begin{cases} -x & x \leq 0 \\ x & x > 0 \end{cases}$$

3. The maximum-function is defined as

$$f(x, y) = \begin{cases} y & x \leq y \\ x & x > y \end{cases}$$

◁

Piecewise defined functions can for example be used to define a more realistic model of neurons.

**Example: Neurons that are linear over some range.**

Earlier, we saw that neurons are linear for some stimuli  $x$ , but clearly not for all: Firstly, the firing rate of a neuron can not be negative, and secondly, there is a maximal firing rate that can not be exceeded. Let us suppose that a neuron responds to a one-dimensional stimulus  $x$  with firing rate  $f(x)$ .

$$f(x) = \begin{cases} 0 & x < 0 \\ x & 0 \leq x \leq f_{max} \\ f_{max} & x > f_{max} \end{cases}$$

◁

Piecewise functions are just like any other function. Within each interval, the functions can be differentiated and integrate like other functions. Nevertheless, there is some care needed at the points at which the intervals meet. In particular, it is often (but not always) the case that piecewise functions are not continuous or differentiable at these points.

### 1.1.7 Sketching Functions

Nowadays, computers are around almost everywhere allowing us to plot functions whenever needed. Nevertheless, being able to imagine how functions look like and sketch them is a useful ability because it gives us a better intuition for what those functions do. There are some simple tricks for imagining and drawing functions which we briefly present in this section. They can basically be classified into two categories. The first is for functions that are transformations of certain basic functions which one usually knows by, like the exponential function, sine and cosine function or easy polynomials like the parabola. The second is for functions that are compositions of known basis functions.

#### 1.1.7.1 Adapting Functions

Everyone knows how to sketch the parabola  $f(x) = x^2$ : It is opened upwards, symmetric, equals one for  $x = \pm 1$  and diverges to infinity for  $x \rightarrow \pm\infty$ . But what about the function  $f(x) = -\frac{1}{2}(x-2)^2 + 5$ ? We will see that it is easy to adapt  $f(x) = x^2$  to make it look like  $f(x) = -\frac{1}{2}(x-2)^2 + 5$ . The first question, we have to answer, is in which order we want to introduce the changes to  $x^2$  to

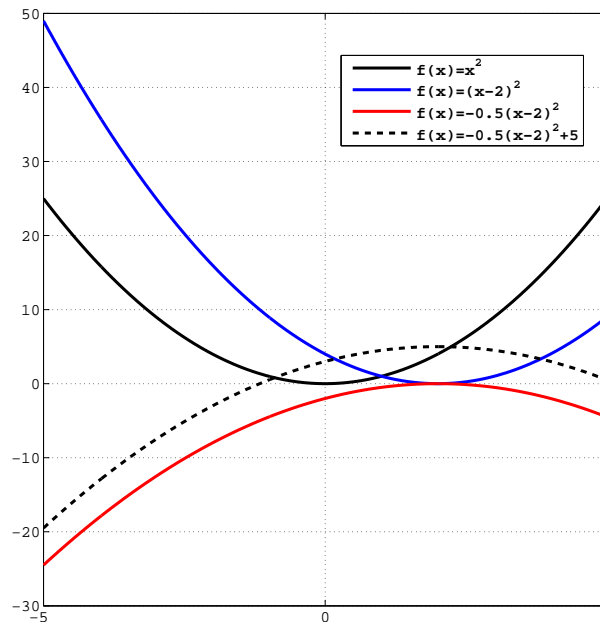


Figure 1.6: Different steps of transforming  $f(x) = x^2$  into  $f(x) = -\frac{1}{2}(x-2)^2 + 5$ .

transform it into  $-\frac{1}{2}(x-2)^2 + 5$ . The answer is: We do that in exactly that order in which we would compute the result of  $-\frac{1}{2}(x-2)^2 + 5$ . This means, we first subtract 2, then square the result, then premultiply  $-\frac{1}{2}$  and finally add 5. If we would not do that we would end up with a different function since we would violate calculation rules at some point in the process.

So let us start to transform  $x^2$ . You can follow the different steps graphically in Figure 1.6. As we just mentioned, the first step is to transform  $x^2$  into  $(x-2)^2$ . Here, we get our first rule: The graph of  $f(x-a)$  is simply the graph of  $f(x)$  shifted by  $a$  to the right. Of course, if  $a$  is negative, shifting by  $a$  becomes a shift to the left. Why is that the case? If you think about it,  $x-a$  can also be seen as shifting the whole  $x$ -axis by  $-a$ , i.e.  $a$  to the left. This, however, is equivalent to shifting  $f$  to the right by  $a$ . Applied to your example, this means that we have to shift  $x^2$  to the right by 2 in order to obtain the graph of  $(x-2)^2$ .

The next step is to include the factor  $-\frac{1}{2}$ . We already know the graph of  $(x-2)^2$ , how does the graph of  $-\frac{1}{2}(x-2)^2$  look like? Well, premultiplying  $-1$  surely reflects the graph along the  $x$ -axis. The factor  $\frac{1}{2}$  squeezes the result. For example,  $y$ -values that used to be  $-1$  are now  $-\frac{1}{2}$ ,  $y$ -values that used to be  $-2$  are now  $-1$ , and so on.

Now, we are almost there. The last step is to include the additive constant  $+5$ . This is easy: It simply shifts the graph of  $-\frac{1}{2}(x-2)^2$  upwards by 5. This is it! We arrived at the graph of  $-\frac{1}{2}(x-2)^2 + 5$  by a few simple adaptation rules for the graph of  $x^2$ . With this few simple rules, you can already sketch a

decent amount of functions.

**Rules for adapting functions**

- The graph of  $f(x - a)$  is simply the graph of  $f(x)$  shifted by  $a$ .
- The graph of  $-f(x)$  is the graph of  $f(x)$  flipped along the  $x$ -axis.
- The graph of  $a \cdot f(x)$  is the graph of  $f(x)$  stretched ( $a > 1$ ) or squeezed ( $0 \leq a < 1$ ) by  $a$ .
- The graph of  $f(x) + b$  is the graph of  $f(x)$  shifted by  $b$  along the  $y$ -axis.

### 1.1.7.2 Compositions of Functions

Sketching compositions of functions is a little bit more art, but for a lot of examples it is not so difficult. As an example, we use the function  $f(x) = \exp(-(x-2)^2)$  which is just a nicer way of writing  $f(x) = e^{-(x-2)^2}$ . The rules from above are not sufficient to sketch this function. There is, however, one rule that we can use: If we know the graph of  $f(x) = \exp(-x^2)$ , we know that we arrive at  $f(x) = \exp(-(x-2)^2)$  by shifting the graph to the right by 2. Therefore, we look at how to sketch the graph of  $f(x) = \exp(-x^2)$  in the following.

The first step, you can always do is to check whether there are function values which are easy to compute and help drawing the graph. Usually, it is a good idea to look at the behavior of  $f(x)$  at  $x = 0$ ,  $x = \pm 1$  and what  $f(x)$  does if  $x$  goes to  $\pm\infty$ . In our case, the interesting cases are  $x = 0$  and  $x \rightarrow \pm\infty$ . The position  $x = 0$  is interesting since anything to the power of 0 equals 1. Therefore  $f(0) = 1$ . When we ask how  $f(x)$  behaves for  $x \rightarrow \pm\infty$  we can first observe that the behavior will be the same at both sides since the squaring operation cancels out the sign. So let us look at what happens when  $x \rightarrow \infty$ . Let us advance step by step, just as before. First, when  $x \rightarrow \infty$ , surely we will get  $x^2 \rightarrow \infty$  as well. By that, we can immediately see that  $-x^2 \rightarrow -\infty$ . Therefore, we only need to know what happens to  $\exp(z)$  if its argument  $z$  assumes a very large negative value. We can rewrite the problem a bit by using one of the calculation rules for exponentials:  $e^{-x} = \frac{1}{e^x}$ . Now, the answer should be easy. If  $-x^2 \rightarrow -\infty$ , there will be a large value in the denominator and, therefore,  $f(x) = \exp(-x^2)$  will approach zero. We can even say a bit more, namely that it will approach zero from above since  $\exp(z)$  can never become negative if  $z \in \mathbb{R}$ .

In a similar manner, you can sketch functions of the form  $f(x) = g(x) + h(x)$  or  $f(x) = \frac{g(x)}{h(x)}$ . First, find a few points where the function value is easy to compute. Then check what happens if  $x$  approaches points, where one of the functions goes to zero or infinity. Then the question is usually, which of the functions “wins”, i.e. which approaches zero or infinity faster. For example,  $f(x) = x^2 - x$  will definitely diverge to infinity for  $x \rightarrow \infty$  since  $x^2$  grows much faster than  $-x$  is able to drag it into the negative side. Similarly  $f(x) = \frac{x}{\exp(x)}$  will approach zero for  $x \rightarrow \infty$  since  $\exp(x)$  grows faster than  $x$  does.

Drawing compositions of functions takes a bit of practice, but is a useful tool for understanding how functions look like and what they do.

## 1.2 Basic Calculus

In this section we will cover basic rules for calculating derivatives and simple integrals. Along with introducing the different rules we will also introduce the derivatives of all the functions covered in the section before. In order to keep the equations simple, we will from now on leave out the brackets for functions like  $\sin$ ,  $\cos$ ,  $\log$ , ... as long as the argument of the function is clear from the context.

### 1.2.1 Derivatives

The derivative  $f'(x_0) = \frac{df}{dx}(x_0)$  (denoted with a “’”,  $\frac{df}{dx}$ , or  $\frac{d}{dx}f$ ) of a function  $f(x)$  has two intuitive meanings:

1. It measures the rate of change of a function at a certain location  $x_0$ .
2. It is the slope of the line touching the function  $f$  at the point  $(x_0, f(x_0))$ . This line is called *tangent line* or simply *tangent*.

Using the first intuition,  $f'(x_0)$  is an approximation of how the function value of  $f(x)$  changes when going from  $x_0$  to  $x_0 + 1$ . If  $f$  is a linear function, the approximation will be exact, that means  $f$  will change exactly by  $f'(x_0)$ . If  $f$  is not linear we will make some error, but we still can use  $f'(x)$  to construct the best linear approximation of  $f$  at  $x_0$ . But first, we will start with an example of the exact case.

**Example** Consider the linear function  $f(x) = 3x$ . According to the first intuition, the derivative  $f'(x)$  is the rate of change of  $f$ , i.e. the change in the function value of  $f$  divided by the change in the value of  $x$ . Consider two arbitrary points  $x_0$  and  $x_1$ . The rate of change is then given by:

$$\begin{aligned} f'(x) &= \frac{f(x_1) - f(x_0)}{x_1 - x_0} \\ &= \frac{3x_1 - 3x_0}{x_1 - x_0} \\ &= \frac{3(x_1 - x_0)}{(x_1 - x_0)} \\ &= 3. \end{aligned}$$

Since any linear function can be written as  $f(x) = ax$ , we just showed that the first derivative  $f'(x)$  of a linear function does not depend on  $x$ . This means that it is the same everywhere. This is what we expect intuitively from a line. Secondly we verified the second intuition for linear functions. The first derivative at a point  $x_0$  is the slope of the tangent line of  $f$  at  $(x_0, f(x_0))$ . Since the tangent line is simply the linear function itself, the first derivative  $f'(x)$  is the slope of the linear function.



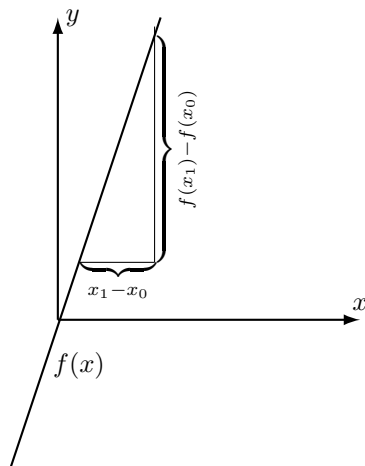


Figure 1.7: Geometrical picture for calculating derivatives of linear functions.

&lt;

The situation changes when we consider arbitrary functions  $f(x)$ . In that case the value of the rate of change, given by the quotient  $\frac{f(x_1) - f(x_0)}{x_1 - x_0}$ , will depend on the choices of  $x_1$  and  $x_0$ . This raises the question how we could define the rate of change in a meaningful way? The second intuition can help us here: In order to get the first derivative at a point  $x_0$  we approximate  $f(x)$  at  $x_0$  with a line and define the first derivative to be its slope. Since we are merely interested in computing the slope of that line, we do not need to compute the full line equation but start with the slope right away. Remember, given a line  $g(x) = ax + t$ , we can compute its slope via  $\frac{g(x_1) - g(x_0)}{(x_1 - x_0)} = \frac{ax_1 + t - ax_0 - t}{(x_1 - x_0)} = a$  where  $x_1$  and  $x_0$  are two arbitrary points. Now imagine we have a line  $g(x) = ax + t$  that contains the two points  $(x_0, f(x_0))$  and  $(x_1, f(x_1))$  (see Figure 1.8). In order to compute its slope we do not need the full line equation. Instead we can simply use the quotient from above and compute

$$\begin{aligned} a &= \frac{g(x_1) - g(x_0)}{x_1 - x_0} \\ &= \frac{f(x_1) - f(x_0)}{x_1 - x_0}. \end{aligned}$$

The slope  $a$  of this line is not yet the first derivative since  $g$  contains  $(x_0, f(x_0))$  and  $(x_1, f(x_1))$ . This means, that (in most cases) it will intersect with  $f$  and

not *touch* it at  $(x_0, f(x_0))$ . However, we can achieve this goal by moving  $x_1$  close to  $x_0$ . Once it is infinitely, called *infinitesimally*, close to  $x_0$ ,  $a = \frac{f(x_1) - f(x_0)}{x_1 - x_0}$  will be the slope of the tangent line, i.e.  $a = f'(x_0)$ . Mathematically this is expressed in terms of a limit. We do not go into the details of limits here, but merely demonstrate, how the derivative of a function  $f$  at  $x_0$  is defined. The rate of change with a infinitesimal close point  $x_1 = x_0 + h$  can be written in terms of limits as

$$f'(x) = \lim_{h \rightarrow 0} \frac{\overbrace{f(x_0 + h) - f(x_0)}^{=x_1}}{\underbrace{h}_{=x_1 - x_0}}.$$

The "lim" expresses that we let  $h$  come infinitesimally close to zero and therefore  $x_1 = x_0 + h$  infinitesimally close to  $x_0$ . The expression  $\lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}$  is called *differential quotient* of  $f$ . Note that, in order to obtain a unique notion of a derivative at a point  $x_0$ , the direction from which  $x_1 = x_0 + h$  approaches  $x_0$  should not matter. This means that  $h$  could be negative ( $x_1$  approaches  $x_0$  from the left) or positive ( $x_1$  approaches  $x_0$  from the right).

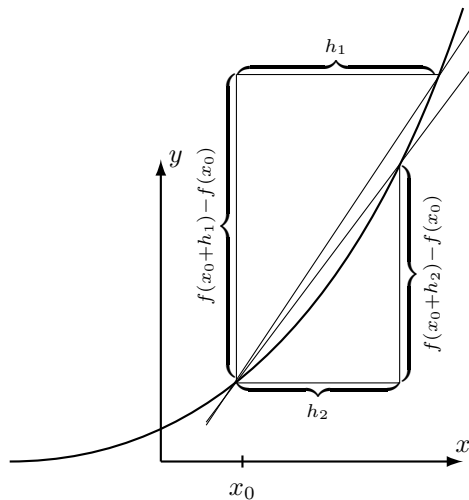


Figure 1.8: Geometrical picture of the differential quotient for arbitrary functions. As  $h$  becomes smaller, slope of the line through  $(x_0, f(x_0))$  and  $(x_0 + h, f(x_0 + h))$  converges to the derivative of  $f$  at  $x_0$ .

Looking at Figure 1.8, we can see that not every function has a derivative at any point  $x_0$ . If the function makes a step at  $x_0$  such that there is a gap

between the function value at  $x_0$  and the function value at  $x_0$  plus or minus an infinitesimal  $h$ , then the value of the differential quotient at  $x_0$  depends on the direction from which  $x_0 + h$  approaches  $x_0$  and the derivative would not be unique. Therefore, functions are only differentiable at points where the function does not make such a step. This property is expressed in the concept of *continuity*.

**Definition (Continuous Function)** A function  $f$  is said to be *continuous at a point*  $x_0$  if

$$\lim_{h \rightarrow 0} f(x_0 + h) = \lim_{h \rightarrow 0} f(x_0 - h) = f(x_0),$$

i.e. if getting infinitesimal close to  $x_0$  (no matter from which side) implies getting infinitesimal close to  $f(x_0)$ .

A function  $f$  is said to be *continuous*, if it is continuous in every point of its domain.

◇

If a function is continuous, we can calculate the value of the differential quotient  $\lim_{h \rightarrow 0} \frac{f(x_0+h)-f(x_0)}{h}$ . If the function is not continuous, the derivative at that point is classified as *not defined*. In a similar fashion, functions that have a kink at  $x_0$  give rise to an undefined derivative at  $x_0$ . The reason is that because of the kink approaching  $x_0$  from the left yields a different slope than approaching  $x_0$  from the right. For example  $f(x) = |x| = \text{abs}(x)$  is not differentiable at  $x_0 = 0$ . At all other points, however, it is. We will see how to differentiate  $f(x) = |x|$  in an example below.

However, if we can find a linear function that with its slope equal to the value of the differential quotient, no matter if  $h$  approaches zero from the left ( $h < 0$ ) or from the right ( $h > 0$ ), the function is called *differentiable*.

**Definition (Differentiable Function)** A function  $f$  is said to be *differentiable at a point*  $x_0$  if there exists a linear function  $L(x)$  such that

$$L(x_0) - \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h} = 0,$$

no matter from which side  $h$  approaches zero.

A function  $f$  is said to be *differentiable*, if it is differentiable in every point of its domain.

◇

Fortunately, we do not need to go through the tedious process of calculating the value of the differential quotient every time. There are easier ways to compute the slope of the tangent at a point  $x_0$ . In the following we will review the

most important rules how to compute derivatives. In most cases those rules are sufficient to compute derivatives of functions that you are dealing with.

We start with the simplest rule: Calculating the derivatives of powers. If  $f(x) = ax^b$ , then the derivative of  $f$  is given by

$$f'(x) = abx^{b-1}.$$

As we can see, the rule is fairly easy: We just premultiply the exponent  $b$ , subtract 1 from the exponent of  $x$  and leave the factor  $a$  untouched. Before looking at a small example, we just introduce another rule: The derivative of a sum equals the derivatives of the single terms:

$$f(x) = g_1(x) + g_2(x) \Rightarrow f'(x) = g_1'(x) + g_2'(x)$$

Or in general:

$$f(x) = \sum_{k=1}^n g_k(x) \Rightarrow f'(x) = \left( \sum_{k=1}^n g_k(x) \right)' = \sum_{k=1}^n g_k'(x).$$

Here,  $g_k$  are arbitrary differentiable functions.

**Example** Assume that  $f(x) = \frac{1}{2}x^2 + 5x + 10$ . In order to apply the rule, we must rewrite the constant 10. Since we already know that  $x^0 = 1$  we can write  $10 = 10x^0$  and  $f$  becomes  $f(x) = \frac{1}{2}x^2 + 5x + 10x^0$ . Now we can apply our rules step by step:

$$\begin{aligned} f'(x) &= \left( \frac{1}{2}x^2 + 5x + 10x^0 \right)' \\ &= \left( \frac{1}{2}x^2 \right)' + (5x)' + (10x^0)' \\ &= x + 5. \end{aligned}$$

As we can see, the last term vanishes since premultiplying 0 cancels the whole term. Since we can write any constant  $b$  that does not depend on  $x$  as  $b = bx^0$ , constants always vanish when calculating the derivative.

◁

We can generalize our example to arbitrary polynomials:

$$f(x) = \sum_{k=0}^n a_k x^k \Rightarrow f'(x) = \sum_{k=1}^n a_k k x^{k-1}.$$

Note, that the index  $k$  of  $f'(x)$  starts at one instead of zero. This is because the last term  $a_0 x^0 = a_0$  vanishes when differentiating.

We can also use this rule to calculate the derivative of roots and ratios.

**Examples**

1. Let  $f$  be  $f(x) = \sqrt{x}$ . Since we can write  $f$  as  $f(x) = x^{\frac{1}{2}}$ , the derivative of  $f$  is given by  $f'(x) = \frac{1}{2}x^{-\frac{1}{2}} = \frac{1}{2\sqrt{x}}$ .
2. Let  $f$  be  $f(x) = \frac{1}{x^a}$ . Since we can write  $f$  as  $f(x) = x^{-a}$ , the derivative is given by  $f'(x) = -ax^{-a-1} = -\frac{a}{x^{a+1}}$

&lt;

As those examples show, we can already calculate the derivative of a fair amount of functions. However, so far we cannot differentiate functions that are compositions of other functions. For example, we cannot get  $f'(x)$  for function  $f(x) = \sqrt{\frac{1}{2}x^2 + 5x + 10}$  with our current set of rules, since  $f(x) = g_2(g_1(x))$  is the composition of  $g_2(y) = \sqrt{y}$  and  $y = g_1(x) = \frac{1}{2}x^2 + 5x + 10$ . Therefore, we introduce a new rule, called *chain rule*: Let  $f$  be  $f(x) = g_2(g_1(x))$  with two arbitrary differentiable functions  $g_1$  and  $g_2$ , then the derivative  $f'(x)$  is given by

$$f'(x) = g_2'(g_1(x)) \cdot g_1'(x).$$

According to this rule we first differentiate  $g_2(y)$  while treating  $y = g_1(x)$  as a variable on its own. After that, we multiply the result with the derivative  $g_1'(x)$  of  $g_1$  with respect to  $g_1$ . Let us illustrate this rule by differentiating  $f(x) = \sqrt{\frac{1}{2}x^2 + 5x + 10}$ .

**Example** Let  $f$  be  $f(x) = \sqrt{\frac{1}{2}x^2 + 5x + 10}$ . As already mentioned before,  $f$  is the composition of  $g_2(y) = \sqrt{y}$  and  $y = g_1(x) = \frac{1}{2}x^2 + 5x + 10$ . We already calculated the derivative of those functions:

$$\begin{aligned} g_2'(y) &= \frac{1}{2\sqrt{y}} \\ g_1'(x) &= x + 5. \end{aligned}$$

Applying the chain rule therefore yields

$$\begin{aligned} f'(x) &= g_2'(g_1(x)) \cdot g_1'(x) \\ &= \frac{1}{2\sqrt{\frac{1}{2}x^2 + 5x + 10}} \cdot g_1'(x) \\ &= \frac{1}{2\sqrt{\frac{1}{2}x^2 + 5x + 10}} \cdot (x + 5). \end{aligned}$$

&lt;

Now, there is only one rule left: How to differentiate products  $f(x) = g_1(x) \cdot g_2(x)$  of two differentiable functions  $g_1$  and  $g_2$ . If  $f$  is the product of two functions  $g_1$  and  $g_2$ , the derivative is given by

$$f'(x) = g_1'(x) \cdot g_2(x) + g_1(x) \cdot g_2'(x).$$

This rule is called *product rule*. It tells us to first differentiate  $g_1(x)$  and multiply the result with  $g_2(x)$ , then do it the other way round and sum the results in the end. Again, an example will illustrate this rule.

**Example** Let  $f$  be  $f(x) = (x+1)(\frac{1}{2}x^2 + 5)$ . Apparently,  $f$  is the product of  $g_1(x) = (x+1)$  and  $g_2(x) = (\frac{1}{2}x^2 + 5)$ . It is easy to calculate their derivatives:

$$\begin{aligned} g_1'(x) &= 1 \\ g_2'(x) &= x. \end{aligned}$$

Therefore, by applying the product rule we get the derivative of  $f$ :

$$\begin{aligned} f'(x) &= (x+1)' \cdot \left(\frac{1}{2}x^2 + 5\right) + (x+1) \cdot \left(\frac{1}{2}x^2 + 5\right)' \\ &= \frac{1}{2}x^2 + 5 + (x+1)x \\ &= \frac{3}{2}x^2 + x + 5. \end{aligned}$$

In this particular case, we can check the rule by expanding

$$\begin{aligned} f(x) &= (x+1)\left(\frac{1}{2}x^2 + 5\right) \\ &= \frac{1}{2}x^3 + \frac{1}{2}x^2 + 5x + 5. \end{aligned}$$

Using the rule for polynomials, yields:

$$\begin{aligned} f'(x) &= \left(\frac{1}{2}x^3 + \frac{1}{2}x^2 + 5x + 5\right)' \\ &= \frac{3}{2}x^2 + x + 5, \end{aligned}$$

which is the same results as the one from the product rule.

◁

In most cases, a fourth rule for differentiating quotients of functions is specified. However, we can derive the *quotient rule* from the rules we already know. Before stating the general quotient rule, we look at a small example.

**Example** Let  $f$  be  $f(x) = \frac{(x+1)}{(\frac{1}{2}x^2+5)}$ . For calculating the derivative of  $f$  we reformulate it first into  $f(x) = (x+1)(\frac{1}{2}x^2+5)^{-1}$ . Now we can apply the product rule and the chain rule to differentiate  $f$ . Let us look at the calculation step by step:

$$\begin{aligned}
 f'(x) \quad \text{Product Rule} \quad & (x+1)' \left( \frac{1}{2}x^2+5 \right)^{-1} + (x+1) \left( \left( \frac{1}{2}x^2+5 \right)^{-1} \right)' \\
 \text{Chain Rule} \quad & (x+1)' \left( \frac{1}{2}x^2+5 \right)^{-1} - (x+1) \left( \frac{1}{2}x^2+5 \right)^{-2} \left( \frac{1}{2}x^2+5 \right)' \\
 = \quad & \frac{(x+1)'}{\left( \frac{1}{2}x^2+5 \right)} - \frac{(x+1) \left( \frac{1}{2}x^2+5 \right)'}{\left( \frac{1}{2}x^2+5 \right)^2} \\
 = \quad & \frac{(x+1)' \left( \frac{1}{2}x^2+5 \right) - (x+1) \left( \frac{1}{2}x^2+5 \right)'}{\left( \frac{1}{2}x^2+5 \right)^2} \\
 = \quad & \frac{-\frac{1}{2}x^2 - x + 5}{\left( \frac{1}{2}x^2+5 \right)^2}.
 \end{aligned}$$

&lt;

In general, it might be easier to calculate all the derivatives first. In this example, however, the derivatives were only resolved at the end to show an important aspect: When looking at the fourth line of the calculation, we see that it consists exclusively of terms that already appeared in the product. We can generalize this, to get the quotient rule: Let  $f$  be a quotient of two functions  $f(x) = \frac{g_1(x)}{g_2(x)}$ . By rewriting  $f(x) = g_1(x) \cdot g_2(x)^{-1}$  and applying the product rule, the chain rule and the rule for powers, we arrive at

$$f'(x) = \frac{g_1'(x)g_2(x) - g_1(x)g_2'(x)}{g_2(x)^2}.$$

Here is a summary of all the rules we just saw:

**Differentiation Rules**

- Derivatives of constant functions: The derivative of any constant function  $f(x) = a$  is  $f'(x) = 0$ .
- Summation Rule: Let  $f(x) = \sum_{k=1}^n g_k(x)$  a sum of arbitrary differentiable functions. Then  $f'(x)$  is given by:

$$f'(x) = \sum_{k=1}^n g'_k(x).$$

- Power Rule: Let  $f$  be  $f(x) = ax^b$ , then  $f'(x)$  is give by:

$$f'(x) = abx^{b-1}.$$

- Chain Rule: Let  $f(x) = g_2(g_1(x))$  be a composition of arbitrary differentiable functions. Then  $f'(x)$  is given by:

$$f'(x) = g'_2(g_1(x)) \cdot g'_1(x).$$

- Product Rule: Let  $f(x) = g_1(x)g_2(x)$  be a product of arbitrary differentiable functions. Then  $f'(x)$  is given by:

$$f'(x) = g'_1(x) \cdot g_2(x) + g_1(x) \cdot g'_2(x).$$

- Quotient Rule: Let  $f(x) = \frac{g_1(x)}{g_2(x)}$  be a quotient of two arbitrary differentiable functions. Then  $f'(x)$  is given by:

$$f'(x) = \frac{g'_1(x)g_2(x) - g_1(x)g'_2(x)}{g_2(x)^2}.$$

We conclude this section by stating derivatives of important functions that you cannot compute using the rules above. Since they occur quite often, you should know them by heart.

**Derivatives of important functions**

- $f(x) = \sin(x) \Rightarrow f'(x) = \cos(x)$
- $f(x) = \cos(x) \Rightarrow f'(x) = -\sin(x)$
- $f(x) = e^x \Rightarrow f'(x) = e^x$
- $f(x) = \ln(x) = \log_e(x) \Rightarrow f'(x) = \frac{1}{x}$



**Example**

Let us calculate the derivative of a little bit more advanced case:  $f(x) = |x| = \text{abs}(x)$ . Figure 1.9 show the graph of the function.

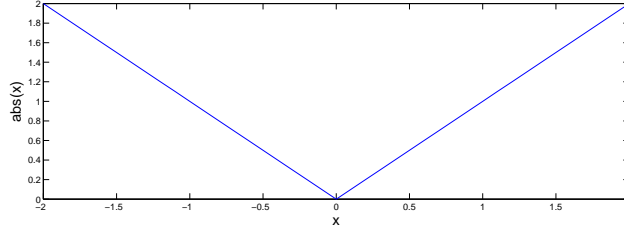


Figure 1.9: Graph of the function  $f(x) = |x|$ .

We can either use the definition of  $|x|$  as a piecewise defined function, or use the fact that  $f(x) = \sqrt{x^2}$ . In both cases the derivative is  $f'(x) = -1$  for  $x < 0$  and  $f'(x) = 1$  for  $x > 0$ . But what happens at  $x_0 = 0$ . If we use the piecewise definition

$$|x| = \begin{cases} -x & x \leq 0 \\ x & x > 0 \end{cases},$$

then the value of the differential quotient  $\lim_{h \rightarrow 0} \frac{|x_0+h| - |x_0|}{h}$  is  $-1$  if we approach  $x_0 = 0$  from the left (i.e.  $h < 0$ ) and  $1$  if we approach  $x_0 = 0$  from the right (i.e.  $h > 0$ ). Since there cannot be a linear function which has the slope  $-1$  and  $1$  at the same time, the derivative is not defined at  $x_0 = 0$ . A similar thing happens when we use  $|x| = \sqrt{x^2}$ . Since  $(\sqrt{x^2})' = \frac{x}{\sqrt{x^2}}$  we cannot choose  $x$  to be zero since this would make the denominator zero and the fraction would not be defined. Therefore the derivative of  $f(x) = |x|$  is given by

$$f'(x) = \begin{cases} -1 & x < 0 \\ \text{undefined} & x = 0 \\ 1 & x > 0 \end{cases}.$$

◁

**Example**

There is an easy numerical way to check and compute derivatives with matlab based on the differential quotient:

$$f'(x) = \lim_{h \rightarrow 0} \frac{\overbrace{f(x_0 + h)}^{=x_1} - f(x_0)}{\underbrace{h}_{=x_1 - x_0}}.$$

The idea is that  $\frac{\overbrace{f(x_0 + h) - f(x_0)}^{=x_1}}{\underbrace{h}_{=x_1 - x_0}}$  already sufficiently approximates  $f'(x)$  for re-

ally small  $h$ . This means, that you choose at a set of closely spaced points  $x_1, \dots, x_i, x_{i+1}, \dots, x_n$  and compute the function values  $f(x_1), \dots, f(x_i), f(x_{i+1}), \dots, f(x_n)$  at those points and compute  $f'(x_i) \approx \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i}$ . Fortunately, there is a built-in matlab function that takes the differences for you. Here is an example for how to compute the numerical derivative of  $f(x) = \sin(x)$ :

```
>> h = .001; x = 0:h:2*pi; % define h and the base points
>> f = sin(x); % sample the function at x
>> df = diff(f)/h; % compute the differences
>> plot(x,f,'k'), hold on; % plot the function
>> plot(x(1:end-1),df,'r'); % plot the approximation
>> plot(x(1:end-1),cos(x(1:end-1)),'g'); % plot the true derivative
```

You surely noticed that we were only able to compute  $f'(x)$  for all but the last  $x$ . The reason is, that there is no  $x_{n+1}$  and  $f(x_{n+1})$  that we could have used for computing  $f'(x_n)$ . However, there is also a little more involved matlab function that computed the derivatives at all base points. Using that function the example becomes:

```
>> h = .001; x = 0:h:2*pi; % define h and the base points
>> f = sin(x); % sample the function at x
>> dfdh = gradient(f,h); % compute the first derivative
>> plot(x,f,'k'), hold on; % plot the function
>> plot(x,dfdh,'r'); % plot the approximation
>> plot(x,cos(x),'g'); % plot the true derivative
```

You can (and should) use this function to check derivatives that you computed analytically.

◁

## 1.2.2 Higher-Order Derivatives

The term *higher-order derivatives* comprises the result of iterated differentiation. Since we can treat  $f'(x)$  as a function on its own, there is really nothing new here. If  $f'(x)$  is differentiable, then  $f(x)$  twice differentiable. This can be extended to higher-orders: The third order derivative which is written  $f'''(x)$  or sometimes  $f^{(3)}(x)$  is just the derivative of the second order derivative.

The second derivative, however, has a special meaning. Geometrically, it is the slope of the slope, i.e. how do the linear approximations at  $f(x)$  vary with  $x$ . This give us a notion of *curvature* of  $f(x)$ . Intuitively, if the linear approximations at  $f(x_0)$  bend away very quickly in the area around  $x_0$ , then  $f(x)$  must have a high curvature at  $x_0$ .

The geometric intuition behind higher-order derivatives is somewhat harder. In particular, our visual system is not very good at judging higher-order derivatives. By visual inspection, it is pretty hard to say if e.g. the fourth derivative is positive or negative.

### Example

Suppose that the function  $f(t)$  is the position of an object in space at time  $t$ . Then,  $f'(t)$  gives us the rate at which the object changes its position, i.e. the object's speed at time  $t$ . The second order derivative  $f''(t)$  gives us the rate at which the object changes its speed, which is simply the acceleration of the object. Similarly, if  $f(t)$  is the concentration of  $Ca^{2+}$ -ions in a neuron at time  $t$ , then  $f'(t)$  gives you the rate at which ions enter or leave the cell at time  $t$ , and  $f''(t)$  indicates whether this rate is getting bigger or smaller.

◁

### 1.2.3 Finding Maxima/Minima of a Function

Derivatives can be used to find global minima or maxima of functions. If  $f(x)$  has a local maximum or minimum at a point  $x_0$ , the tangent line is horizontal. This means that the slope of the tangent and, therefore, the derivative of the function at that point  $x_0$  must be  $f'(x_0) = 0$ . Mathematically speaking,  $f'(x_0) = 0$  is a *necessary condition* for the function  $f$  having a maximum or a minimum, i.e. “ $x_0$  is maximum/minimum”  $\Rightarrow f'(x_0) = 0$ . However, it can happen, that  $f'(x_0) = 0$  but  $x_0$  is not a maximum or a minimum. Those points are called saddle-points. They arise for example, when  $f(x)$  increases, becomes more and more flat when approaching  $x_0$ , is completely flat at  $x_0$ , and increases again afterwards. Therefore, since  $f'(x_0) = 0$  does not imply a maximum or a minimum, we must find a condition that does that, i.e. we need to find a *sufficient condition*, one for which “condition”  $\Rightarrow$  “ $x_0$  is maximum/minimum”.

Since we already know that we only need to look at points  $x_0$  where  $f'(x_0) = 0$ , we can think about what must happen for  $x_0$  to be—say—a maximum. If  $x_0$  is a maximum, then the slope of tangent lines of points to the right of  $x_0$  must become more and more negative. At the same time the slope of tangent lines at points left to  $x_0$  must become less and less negative and approach zero at  $x_0$ . If you think about it, this is exactly the same as requiring  $f''(x_0) < 0$ . Similarly,  $f'(x_0) = 0$  and  $f''(x_0) > 0$  imply a minimum of  $f$  at  $x_0$ . If  $f'(x_0) = 0$ ,  $f''(x_0) = 0$ , but  $f'''(x_0) \neq 0$  then  $f$  has a saddle-point at  $x_0$ . In general, if the first  $n$  derivatives are zero and the  $(n+1)$ th derivative  $f^{(n+1)}(x_0)$  is not equal to zero, then  $x_0$  is a minimum, if  $n+1$  is even and  $f^{(n+1)}(x_0) > 0$ , a maximum if  $n+1$  is even and  $f^{(n+1)}(x_0) < 0$ , and a saddle-point if  $n+1$  is odd and  $f^{(n+1)}(x_0) \neq 0$ .

**Finding Maxima, Minima, and Saddle-Points**

Let  $f$  be a  $(n + 1)$  times differentiable function. Then  $f$  has a

- *minimum* at  $x_0$  if and only if  $f^{(k)}(x_0) = 0$  for  $k = 1, \dots, n$  and  $f^{(n+1)}(x_0) > 0$  with  $(n + 1)$  even.
- *maximum* at  $x_0$  if and only if  $f^{(k)}(x_0) = 0$  for  $k = 1, \dots, n$  and  $f^{(n+1)}(x_0) < 0$  with  $(n + 1)$  even.
- *saddle-point*  $x_0$  if and only if  $f^{(k)}(x_0) = 0$  for  $k = 1, \dots, n$  and  $f^{(n+1)}(x_0) \neq 0$  with  $(n + 1)$  odd.

It should be pointed out that this procedure only yields *local* maxima, and not necessarily *global* ones.

**Example**

Consider the function  $f(x) = (x - 1)^2$ . Then  $f'(x) = 2(x - 1)$  and  $f''(x) = 2$ . In order to find candidates for a maximum or a minimum, we set  $f'(x) = 0$  and solve for  $x$ :

$$\begin{aligned} f'(x) = 0 &\Leftrightarrow 2(x - 1) = 0 \\ &\Leftrightarrow x = 1. \end{aligned}$$

Since  $f''(x) > 0$  for all  $x$  (therefore also for  $x = 1$ ),  $x = 1$  must be a minimum of  $f$ .

&lt;

**Example: Estimating the rate of a Poisson distribution**

A probabilistic model for a neuron that generates completely random spike trains with no temporal structure is the homogeneous *Poisson process*. Since there is no temporal structure, it serves as a baseline model that other—more advanced—models for spike trains can compare to. You can generate spike trains from that model as follows: Bin the time axis into sufficiently small bins (say  $1ms$ ). For each bin you randomly place a spike with a certain small probability. The matlab code looks like this

```
>> t = [0:0.001:10];
>> p = 0.05; % spiking probability in each bin
>> ind = find( rand(size(t)) <= p); % sample spike times
>> s = zeros(size(t)); % generate spike train
>> s(ind) = 1; % insert spikes
>> plot(t,s,'k')
```

of observing  $k$  spikes in a time window of length  $\Delta t = 1$  is given by the Poisson distribution

$$p(k \text{ spikes}) = \frac{\lambda^k e^{-\lambda}}{k!}$$

or simply

$$p(k) = \frac{\lambda^k e^{-\lambda}}{k!}.$$

$\lambda$  is called the rate of the Poisson distribution. It tells us how many spike we expect in one second on average. The question that we want to solve in this example is how to estimate  $\lambda$  from a number of observed spike trains. But lets first look at our matlab example again and check, whether our spike counts in one second really follow a Poisson distribution. For that reason, we generate a large number of spike trains of length 1s, count the spikes in each spike train and look at the empirical histogram. The matlab code looks like this:

```
>> t = [0:0.001:1];
>> p = 0.05; % probability of spike in each bin
>> m = 5000; % number of spike trains
>> U = rand(m,length(t));
>> S = zeros(size(U)); % generate spike train matrix (trains X time)
>> S(U <= p) = 1;
>> figure
>> imagesc(1-S); colormap gray;
>> xlabel('time bins'); ylabel('spike trains')
>> figure C = sum(S,2); % get spike counts
>> K = [1:max(C)];
>> hist(C,K); % plot histogram
>> hold on
>> lambda = p*length(t);
>> plot(K, m*poisspdf(K,lambda),'r','LineWidth',2); hold off
```

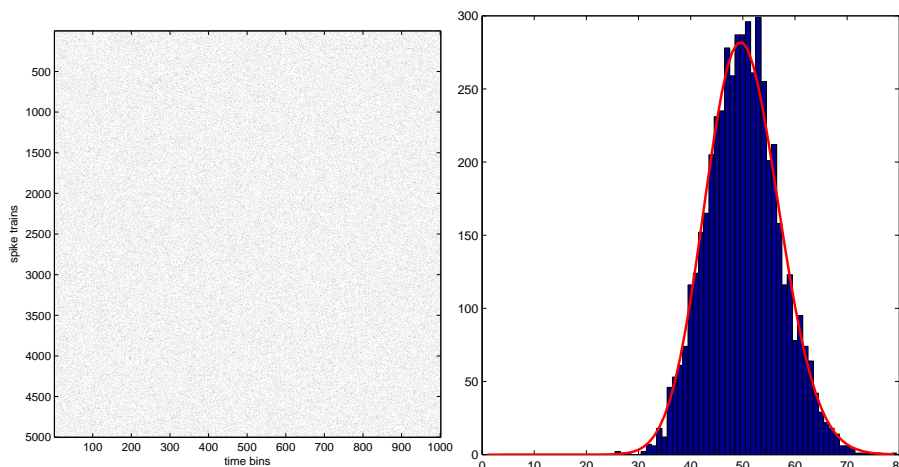


Figure 1.10: Poisson spikes with rate  $\lambda = 50.05$ .

A typical output of the matlab code fragment is shown in Figure 1.10. Before we see how to estimate  $\lambda$  from the data, let us first spend a thought about what values we expect. The answer is not too difficult. Since we see a spike with probability  $p = 0.05$  in each bin, all bins are independent, we expect to see a total number of  $\lambda = n_{bins} \cdot p$  spikes per second where  $n_{bins}$  is the number of bins per second. Strictly speaking, the rate  $\lambda$  is defined as the number of bins times the probability of spiking per bin in the limit of infinite  $n_{bins}$ . The idea is, that the larger the number of bins gets, the lower the probability of observing a spike in a specific bin becomes, i.e. it goes to zero for  $n_{bins}$  to infinity. Therefore, the product of both number can be something finite, i.e. the rate  $\lambda$ .

Let us now turn to the question of how to estimate  $\lambda$ . Assume that you are only given the number of spikes  $k_1, \dots, k_m$  in each of the spike trains of length  $1s$ . Certainly you also do not know  $p$ . However, you do know (or do assume) that the spikes have been generated as described above and therefore the spike counts  $k$  must be Poisson distributed. We further assume that each spike train has been generated independently of the others.

Given a distribution and a number of observations, one principled way of estimating the parameters of that distribution (in our case the rate  $\lambda$ ) is the *maximum likelihood* method. The idea behind that method is simple: We choose  $\lambda$  such that the probability of our observations is maximized:

$$\hat{\lambda} = \operatorname{argmax}_{\lambda \in \mathbb{R}^+} p(k_1, \dots, k_m | \lambda).$$

Since we assumed that the observations are independent from each other, the probability of observing all spike counts jointly, is the product of the probabilities of observing each single spike count:

$$p(k_1, \dots, k_m | \lambda) = \prod_{i=1}^m p(k_i | \lambda).$$

For each single spike count, we know that  $k_i$  is Poisson distributed and therefore

$$\begin{aligned} p(k_1, \dots, k_m | \lambda) &= \prod_{i=1}^m p(k_i | \lambda) \\ &= \prod_{i=1}^m \frac{\lambda^{k_i} e^{-\lambda}}{k_i!}. \end{aligned}$$

Now, we only need to find that maximum of  $f(\lambda) = \prod_{i=1}^m \frac{\lambda^{k_i} e^{-\lambda}}{k_i!}$  with respect to  $\lambda$ . Unfortunately, taking the derivative of  $f(\lambda) = \prod_{i=1}^m \frac{\lambda^{k_i} e^{-\lambda}}{k_i!}$  is very complicated. Luckily, there is a simple trick that we can apply: By taking the log of both sides changes the function values, but leaves the position of the maxima unchanged (the reason for that is, that log is a strictly increasing continuous function). However, taking the log makes the right hand side considerably easier

to deal with:

$$\begin{aligned}
 \log f(\lambda) &= \log \left( \prod_{i=1}^m \frac{\lambda^{k_i} e^{-\lambda}}{k_i!} \right) \\
 &= \sum_{i=1}^m \log \frac{\lambda^{k_i} e^{-\lambda}}{k_i!} \\
 &= \sum_{i=1}^m (\log \lambda^{k_i} + \log e^{-\lambda} - \log k_i!) \\
 &= \sum_{i=1}^m (k_i \log \lambda - \lambda - \log k_i!).
 \end{aligned}$$

Let us now compute the maximum of  $\sum_{i=1}^m k_i \log \lambda - \lambda - \log k_i!$ . For that reason, we first compute the first derivative

$$\begin{aligned}
 \frac{d}{d\lambda} \left( \sum_{i=1}^m k_i \log \lambda - \lambda - \log k_i! \right) &= \sum_{i=1}^m \left( k_i \frac{d}{d\lambda} \log \lambda - \frac{d}{d\lambda} \lambda - \frac{d}{d\lambda} \log k_i! \right) \\
 &= \sum_{i=1}^m \left( k_i \frac{1}{\lambda} - 1 \right) \\
 &= -m + \frac{1}{\lambda} \sum_{i=1}^m k_i.
 \end{aligned}$$

Setting the first derivative to zero and solving for  $\lambda$  yields

$$\begin{aligned}
 -m + \frac{1}{\lambda} \sum_{i=1}^m k_i = 0 &\Leftrightarrow \frac{1}{\lambda} \sum_{i=1}^m k_i = m \\
 &\Leftrightarrow \lambda = \frac{1}{m} \sum_{i=1}^m k_i.
 \end{aligned}$$

In order to check that it is really a maximum, we compute the second derivative

$$\frac{d}{d\lambda} \left( -m + \frac{1}{\lambda} \sum_{i=1}^m k_i \right) = -\frac{1}{\lambda^2} \sum_{i=1}^m k_i.$$

Since rates are always positive (i.e.  $\lambda > 0$ ) the second derivative is always smaller than zero and, therefore,  $\hat{\lambda} = \frac{1}{m} \sum_{i=1}^m k_i$  is really a maximum. It also has a very intuitive interpretation. The maximum likelihood estimate for the rate of a Poisson distribution is the mean of the observed spike counts. In that sense it is also not surprising, that the mean of a Poisson distribution is  $\lambda$ .

You should include the maximum likelihood estimate in the matlab example from above and compare it to the expected rate  $\lambda = p \cdot n_{bins}$ . The code fragment, you have to add, is:

```
>> lambda_est = mean(C);
```

◁

### 1.2.4 Approximating Functions Locally by Lines and Polynomials

**Example: Finding a linear approximation to the function  $\sin(x)$  for  $x$  near 0**

Suppose  $f(x) = \sin(x)$ . If we want to approximate  $f(x)$  at  $x_0 = 0$  by a line, i.e. a function of the form  $g(x) = wx + b$ , such that  $g(x)$  is a good approximation for  $x$  near 0. From the definition of the derivative, we saw that the derivative of a function at  $x_0$  is really the slope of a linear approximation of  $f(x)$  at  $x_0$ . Therefore we can just compute  $f'(x) = \cos(x)$  and evaluate it at  $x_0 = 0$ . This yields the value of  $w$ , since  $w$  is the derivative of  $g(x)$  and we want the derivatives of the function  $f$  and its linear approximation  $g$  to be the same. In our case,

$$w = f'(x_0) = \cos(0) = 1.$$

Now, we must ensure, that  $g(x_0) = f(x_0)$ . We do that by adjusting  $b$ . The offset  $b$  shifts the function horizontally. If we want  $g$  to have the same value at  $x_0$  as  $f$ , we just need to add  $f(x_0)$ . Therefore  $b = f(x_0)$ . Since  $f(x_0) = \sin(0) = 0$ , we set  $b = 0$  and get  $g(x) = wx$  as linear approximation of  $f(x)$  at  $x_0$ . Figure 1.11 plots the linear approximation.

◁

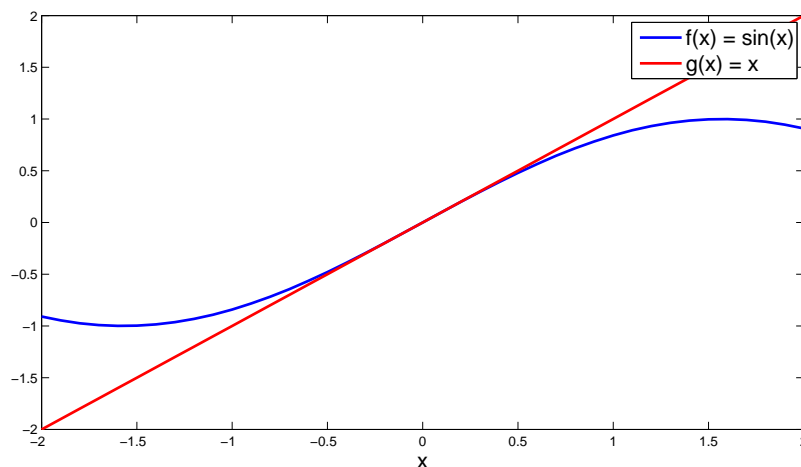


Figure 1.11: Graph of the function  $f(x) = \sin(x)$  and its linear approximation  $g(x) = x$  at  $x_0 = 0$ .

In the example, we implicitly used that the function is approximated at  $x_0 = 0$  when computing the offset  $b$ . In general, if  $x_0 \neq 0$ , we cannot simply use  $b = f(x_0)$ , since



$$\begin{aligned}
g(x_0) &= wx_0 + b \\
&= f'(x_0)x_0 + f(x_0) \\
&\neq f(x_0) \text{ for } w, x_0 \neq 0.
\end{aligned}$$

How can we compute the offset  $b$  in the general case  $x_0 \neq 0$ ? We can answer this question by some geometrical thoughts. First of all, imagine we set  $b = 0$ . The approximating function  $g(x) = wx = f'(x_0) \cdot x$  would be a true linear function that had the same slope as  $f$  at  $x_0$  but not the same function value  $g(x_0) = wx_0 \neq f(x_0)$ . Since we know that  $b$  lifts  $g$  along the  $y$ -axis we need to find out by how much we must lift it in order to obtain  $g(x_0) = f(x_0)$ . The answer is: We need to lift it by the difference between the true function value  $f(x_0)$  and the linear function  $g(x_0) = wx_0$  that has the same slope but the wrong offset, i.e.  $b = f(x_0) - wx_0$ . By substituting  $w = f'(x_0)$  and  $b = f(x_0) - wx_0$  in the general line equation, we end up with

$$\begin{aligned}
g(x) &= f'(x_0)x + f(x_0) - f'(x_0)x_0 \\
&= f(x_0) + f'(x_0)(x - x_0)
\end{aligned}$$

as the best “linear” approximation (it is not linear, it is a line) of  $f$  at  $x_0$ .

The quality of the approximation depends, of course, on the properties of  $f(x)$ . Clearly, in the example above, the approximation is really bad for e.g.  $x = \pi$ .

### Example

The linear approximation to  $f(x) = e^x$  at  $x_0 = 0$  is given by

$$\begin{aligned}
g_0(x) &= f(x) + f'(x_0)(x - x_0) \\
&= e^0 + e^0 \cdot x \\
&= 1 + x.
\end{aligned}$$

The approximation at  $x_0 = 1$  is given by:

$$\begin{aligned}
g_1(x) &= f(x) + f'(x_0)(x - x_0) \\
&= e^1 + (x - 1) \cdot e^1 \\
&= xe.
\end{aligned}$$

Figure 1.12 shows the graphs of these functions.

◁

As we can see in figure 1.12, the linear approximations to  $e^x$  are not very good. Can we find a better approximation by using a quadratic approximation, i.e. one that uses a polynomial of degree 2?

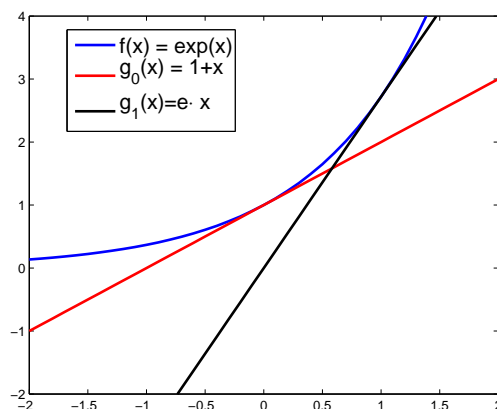


Figure 1.12: Graph of the functions  $f(x) = \exp(x)$  and its linear approximation  $g_0(x) = 1 + x$  and  $g_1(x) = ex$  at  $x_0 = 0$  and  $x_0 = 1$ .

The answer is yes. We simply add a quadratic term to our approximating function  $g$ . So far  $g$  had the form  $g(x) = w(x - x_0) + b$ . To turn it into a quadratic approximation, we need to add a term, that is quadratic in  $(x - x_0)$ . For reasons that will become clear soon this additional term is  $\frac{1}{2}v(x - x_0)^2$  yielding  $g(x) = \frac{1}{2}v(x - x_0)^2 + w(x - x_0) + b$ . The additional work we have to do in comparison to a simple linear approximation is to determine the value of  $v$ . Analogously to determining the value of  $w$ , the value of  $v$  is simply  $v = f''(x_0)$ . By making this choice, we enforce that *the value, the slope and the curvature* of  $g$  and  $f$  match at  $x_0$ , i.e.

$$\begin{aligned} f(x_0) &= g(x_0) \\ f'(x_0) &= g'(x_0) \\ f''(x_0) &= g''(x_0). \end{aligned}$$

The reason why the quadratic term is multiplied with  $\frac{1}{2}$  is to cancel the exponent 2 when taking the derivatives. It is an instructive exercise to check that the function value and the first two derivatives really match at  $x_0$  for  $w = f'(x_0)$  and  $v = f''(x_0)$ .

### Example

For our example  $f(x) = e^x$ , the best quadratic approximation is  $g(x) = 1 + x + \frac{x^2}{2}$ .

◁

This example can be generalized to polynomials of arbitrary order. We can find better local approximations by matching the (higher order) derivatives. If we want to match the first three derivatives, then we have to use approximations by a polynomial of degree 3, 4 and so on. It is a remarkable fact that any function that is sufficiently differentiable can (locally) be approximated arbitrarily well by polynomials:

**Theorem (Taylor/MacLaurin)**

For a differentiable function  $f$  and a  $x$  near  $x_0$ ,  $f(x)$  can be approximated by

$$f(x) \approx f(x_0) + (x - x_0)f'(x_0) + \frac{1}{2}(x - x_0)^2 f''(x_0) + \frac{1}{6}(x - x_0)^3 f'''(x_0) + \dots$$

or in general

$$f(x) \approx f_{(n)}(x) = \sum_{k=0}^n \frac{1}{k!} (x - x_0)^k f^{(k)}(x_0),$$

where  $f^{(k)}(x_0)$  denotes the  $k$ -th order derivative of  $f$ , and  $f_{(n)}$  is referred to as the *Taylor series* approximation to  $f(x)$  at  $x_0$ , or simply the *Taylor approximation of order  $k$*  at  $x_0$ .

◇

The theorem states, that every function can locally be approximated by a polynomial. The higher the order of the polynomial, the more precise is the approximation. The linear approximation we had above is therefore really a special case for a polynomial of degree 1.

**Examples**

1. The cubic approximation to  $f(x) = e^x$  at  $x_0 = 0$  is given by

$$\begin{aligned} g(x) &= f(x_0) + f'(x_0)(x - x_0) + f''(x_0)\frac{1}{2}(x - x_0)^2 + f'''(x_0)\frac{1}{6}(x - x_0)^3 \\ &= e^0 + e^0 x + \frac{1}{2}e^0 x^2 + \frac{1}{6}e^0 x^3 \\ &= 1 + x + \frac{1}{2}x^2 + \frac{1}{6}x^3. \end{aligned}$$

We can generalize that to polynomials of any given order  $n$ : The  $n$ -th order approximation to  $f(x) = e^x$  is given by

$$g_n(x) = \sum_{k=0}^n \frac{1}{k!} (x - x_0)^k.$$

2. The linear and the quadratic approximation to  $\cos(x)$  for  $x_0 = 0$  are given by

$$\begin{aligned} g_{lin}(x) &= \cos(x_0) - \sin(x_0) \cdot x \\ &= 1 \end{aligned}$$

and

$$\begin{aligned} g_{quadr}(x) &= \cos(x_0) - \sin(x_0) \cdot x - \frac{1}{2} \cos(x_0) x^2 \\ &= 1 - x^2. \end{aligned}$$

3. The quadratic approximation to  $\sin(x)$  at  $x_0 = 0$  is given by

$$\begin{aligned} g(x) &= \sin(x_0) + \cos(x_0) \cdot x - \frac{1}{2} \sin(x_0) \cdot x^2 \\ &= x. \end{aligned}$$

4. Here is a small piece of matlab code that plots the Taylor approximation of  $e^x$  at  $x_0 = 1$  up to the order of  $n = 10$ .

```
>> t = [-1:0.01:2]; % define the range where we want to plot exp
>> x0 = 1; % set x0 = 1; you can change that
>> plot(t,exp(t),'k-', 'LineWidth',3), hold on % plot exp
>> fk = exp(x0); % compute constant part
>> n = 1; % set the polynomial order to 1
>> for k = 1:10 % loop over
    plot(t,fk,'-r', 'LineWidth',3); % plot the best approx. in red
    pause % wait until someone presses a key
    plot(t,fk,'-b', 'LineWidth',3); % plot the same curve in blue
    n = n*k; % update the factorial function
    fk = fk + 1/n*(t-x0).^k.*exp(x0); % get the next polynomial
end
>> hold off
```

Note that this code only works for the function  $f(x) = \exp(x)$  since  $\exp$  is its own derivative.

<

### Exercise

E

Find the quadratic approximation to  $f(x)=x^4 + 3x^3 + x^2 + x$  at  $x_0 = 1$ .

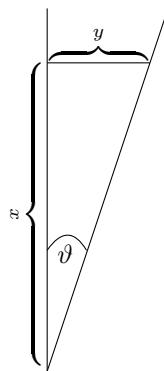
<

Being able to approximate functions by polynomials is extremely useful, and done very often. Because many functions are hard to analyze, people work with their linear or quadratic approximations in many cases. By using these approximations, one can often get away by just using polynomials. Furthermore, finding Taylor approximations is often not hard, you just have to be able to differentiate. You should keep in mind though that the approximation is only local, i.e. for  $x$  near  $x_0$ , and you should be careful that the approximation is really appropriate for your particular application.

Here are two applications where the Taylor expansion is used.

**Example: Estimating small angles**

A calculation that has to be done for many psychophysical experiments it to compute the degree of visual angle. A part of the problem (it is still part of the exercises which) is to figure out the angle  $\vartheta$  in the following setup:



The natural way to do so is simply by  $\tan \vartheta = \frac{y}{x}$  and, therefore,  $\vartheta = \arctan \frac{y}{x}$ . However, if you do not have a computer at your disposal (e.g. when reading the methods section of a paper on much more comfortable couch) you would still like to have an idea what  $\vartheta$  is, even if you cannot invert the  $\tan$  function in your head. Since the usual setup involves large  $x$  and small  $y$ , we expect  $\vartheta$  to be very small, i.e. close to zero. Therefore, we can simply start by computing the first order Taylor approximation of  $\tan$  around 0. First, we need the derivative of

course:

$$\begin{aligned}
 (\tan \vartheta)' &= \left( \frac{\sin \vartheta}{\cos \vartheta} \right)' \\
 &= \frac{\cos^2 \vartheta + \sin^2 \vartheta}{\cos^2 \vartheta} \\
 &= \frac{1}{\cos^2 \vartheta}.
 \end{aligned}$$

At zero, the first derivative is simply one. Since  $\tan 0 = 0$ , the first order Taylor approximation has the simple form

$$\tan \vartheta \approx f_{(1)}(\vartheta) = \vartheta.$$

Fortunately, this function is really easy to invert: It is simply the function itself  $f_{(1)}^{-1}(\vartheta) = \vartheta$ . This shows that we can use  $\frac{y}{x}$  for small  $y$  and large  $x$  itself as good estimate for the angle  $\vartheta$ .

<

### Example: An error analysis of depth perception

Very often a physical system shows a very specific error behavior, i.e. how the error in the inputs shows up in the outputs of a system. This error transformation is important to know, since it tells us how an error in the input will affect the quality of the system's output. A good example for that is depth perception from disparity. In this case the error in the estimated depth will grow quadratically with depth. This is what we will show in this example.

In the simplest case of two parallel pinhole cameras with a focal length  $f$  a distance of  $b$  between them, the distance of a point  $\mathbf{x}$  space to the view planes of the cameras can be estimated by  $d(\varrho) = \frac{fb}{\varrho}$  where  $\varrho = x_l - x_r$ , the disparity of the two images of  $\mathbf{x}$ , i.e. the distance between the  $x$ -coordinates of the image of  $\mathbf{x}$  in the two view planes of the cameras. Consider we have measured a certain disparity  $\hat{\varrho}$  which is  $\delta$  away from the true disparity  $\varrho$ , i.e. our measurement error was  $\delta$ , or  $\hat{\varrho} = \varrho + \delta$ . To see how this error affects the depth estimation  $d$  we make a first order Taylor expansion around the true value  $\varrho$ :

$$\begin{aligned}
 d(\hat{\varrho}) &\approx d_{(1)}(\hat{\varrho}) \\
 &= \frac{fb}{\varrho} - \frac{fb}{\varrho^2}(\hat{\varrho} - \varrho) \\
 &= \frac{fb}{\varrho} - \frac{fb}{\varrho^2}(\varrho + \delta - \varrho) \\
 &= \frac{fb}{\varrho} - \frac{fb}{\varrho^2}\delta.
 \end{aligned}$$

The true depth error is  $d(\hat{\varrho}) - d(\varrho)$ . Using the Taylor expansion instead of  $d$ , we get  $d(\hat{\varrho}) - d(\varrho) \approx \frac{fb}{\varrho^2}\delta$ . From  $d(\varrho) = \frac{fb}{\varrho}$  we can read off that the depth  $d$  is

inversely proportional to the true disparity, i.e.  $d \sim \frac{1}{\rho}$ . Plugging that into our expression for the depth estimation error, we obtain  $d(\hat{\varrho}) - d(\varrho) \sim fbd^2\rho$ , which shows that the influence of an error in the inputs grows quadratically with the depth that we want to estimate.

◁

Chapter 2

Appendix



## 2.1 Notation and Symbols

- The capital Greek letter sigma " $\sum$ " (sigma like sum) denotes a sum over several elements. Usually the components of the sum are indexed with lower case roman letters starting from " $i$ ". The starting index is indicated below the " $\sum$ " and final index is indicated on top. For example, the sum over  $n$  real numbers  $x_1, \dots, x_n \in \mathbb{R}$  is denoted by

$$x_1 + \dots + x_n = \sum_{i=1}^n x_i.$$

Sometimes, when summing over all elements of a set, the set is indicated below the sigma. For example, summing all elements of the set  $A = \{1, 2, 3, 4, \dots, 15\}$ , could be written as  $\sum_{x \in A} x$  as well as  $\sum_{n=1}^{15} n$ .

- The capital Greek letter pi " $\prod$ " is used in an analogous manner for products, i.e.

$$x_1 \cdot \dots \cdot x_n = \prod_{k=1}^n x_k.$$

# Bibliography

- [1] Eric C. Kandel and James H. Schwartz. *Principles of Neural Science*. Elsevier Science Publishing Co., Inc., 2 edition, 1985.