
CONVEX OPTIMIZATION

by

MARTIN LOTZ
School of Mathematics
The University of Manchester

MATH36061
Semester 1, 2016/17
December 2016

Contents

Contents	ii
I Introduction	1
Lecture 1	3
1.1 What is an optimization problem?	3
1.2 Examples of optimization problem	4
II Unconstrained Optimization	15
Lecture 2	17
2.1 Unconstrained optimization	17
2.2 Convex functions	18
Lecture 3	23
3.1 Gradient descent	23
Lecture 4	27
4.1 Step length selection	27
4.2 Convergence of iterative methods	29
4.3 Convergence of gradient descent	29
Lecture 5	31
5.1 Newton's Method	31
5.2 Quasinewton methods	35
Lecture 6	37
6.1 Supervised Learning	37
III Linear Programming	45
Lecture 7	47

7.1 Convex sets	47
Lecture 8	53
8.1 Linear Programming Duality: a first glance	53
8.2 Polyhedra	54
Lecture 9	59
9.1 Farkas' Lemma	59
9.2 Linear programming duality	61
Lecture 10	63
10.1 A first algorithm for linear programming	63
Lecture 11	67
1.1 An optimality condition	67
1.2 Newton's method for solving equations	69
Lecture 12	71
12.1 Towards an efficient interior-point method	71
Lecture 13	75
13.1 Path-following methods	76
13.2 Analysis of Path-following	78
IV Non-linear Convex Optimization	81
Lecture 14	83
14.1 Quadratic Programming and Portfolio Optimization	83
Lecture 15	89
15.1 A first-order optimality condition	89
15.2 Lagrangian duality	90
Lecture 16	93
16.1 Constraint qualification	94
16.2 Karush-Kuhn-Tucker optimality conditions	96
Lecture 17	99
17.1 The logarithmic barrier	100
17.2 The central path	102
Lecture 18	105
18.1 Linear Support Vector Machines	105
18.2 Extensions	108

V Semidefinite Programming	113
Lecture 19	115
19.1 Semidefinite programming	115
19.2 Semidefinite programming duality	117
Lecture 20	119
20.1 Semidefinite relaxation	119
20.2 The Goemans-Williamson Algorithm	121
Appendix	123
1 Asymptotic notation	123
2 Linear Algebra	124
3 Calculus	135
4 Finite precision arithmetic	143

Part I

Introduction

Lecture 1

“[N]othing at all takes place in the universe in which some rule of maximum or minimum does not appear.”

— Leonhard Euler

Mathematical optimization, traditionally also known as mathematical programming, is the theory of optimal decision making. Optimization problems arise in a large variety of contexts, including scheduling and logistics problems, finance, optimal control, signal processing, and machine learning. The underlying mathematical problem always amounts to finding parameters that minimize¹ (cost) or maximize (utility) an objective function in the presence or absence of a set of constraints. An important special case is the class of *convex optimization* problems. Such problems will be the main focus of this course.

1.1 What is an optimization problem?

A general mathematical optimization problem is a problem of the form

$$\begin{aligned} & \text{minimize} && f(\boldsymbol{x}) \\ & \text{subject to} && \boldsymbol{x} \in \Omega \end{aligned} \tag{1.1}$$

where $f: \mathbb{R}^n \rightarrow \mathbb{R}$ are real-valued **objective function** and $\Omega \subseteq \mathbb{R}^n$ is a set defining the **constraints**. Among all $\boldsymbol{x} \in \Omega$, we seek one with smallest f -value. Typically, the constraint set Ω will consist of such $\boldsymbol{x} \in \mathbb{R}^n$ that satisfy certain equations and inequalities,

$$f_1(\boldsymbol{x}) \leq 0, \dots, f_m(\boldsymbol{x}) \leq 0, g_1(\boldsymbol{x}) = 0, \dots, g_p(\boldsymbol{x}) = 0.$$

A vector \boldsymbol{x}^* satisfying the constraints is called an *optimum*, a *solution*, or a *minimizer* of the problem (1.1), if $f(\boldsymbol{x}^*) \leq f(\boldsymbol{x})$ for all other \boldsymbol{x} that satisfy the constraints. Note that replacing f by $-f$, we could equivalently state the problem as a maximization problem. In this course we are mostly concerned with functions and constraint sets that are **convex**.

- A set $C \subseteq \mathbb{R}^n$ is **convex**, if for all $\boldsymbol{x}, \boldsymbol{y} \in C$ and $\lambda \in [0, 1]$, $\lambda\boldsymbol{x} + (1 - \lambda)\boldsymbol{y} \in C$. That is, for any two points in C , the line segment connecting them is also in C .

¹For the sake of consistency with most of the literature, throughout these notes we use the American spelling of minimizer and maximizer with “z” instead of “s”.

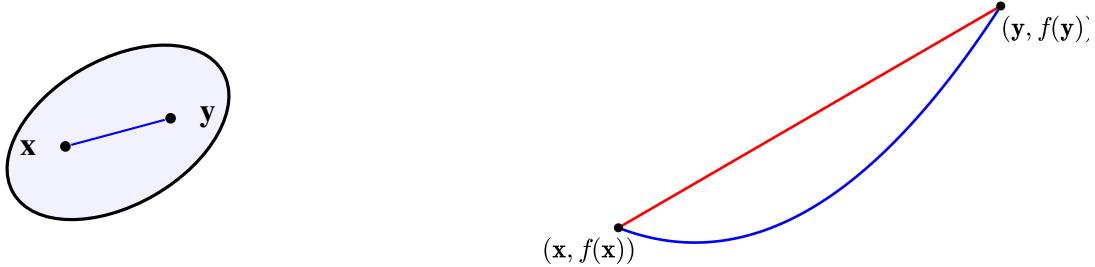


Figure 1.1: A convex set and a convex function

- A function $f: C \rightarrow \mathbb{R}$ is convex, if C is convex and for all $\mathbf{x}, \mathbf{y} \in C$ and $\lambda \in [0, 1]$, $f(\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y})$.

A **convex optimization** problem is one where the set of constraints Ω and the function f are convex. While most general optimization problems are practically intractable, convex optimization problems can be solved efficiently, and still cover a surprisingly large range of applications!

1.2 Examples of optimization problem

Countless problems from science and engineering can be cast as optimization problems. We present a few first examples, many more will follow in the course of this lecture. The examples below come with associated Python code. At this moment it is not expected that you understand them in detail, they are merely intended to illustrate some of the problems that convex optimization deals with, and how they can be solved.

Example 1.1. Suppose we want to understand the relationship of a quantity Y (for example, sales data) to a series of *predictors* X_1, \dots, X_p (for example, advertising budget in different media). We can often assume the relationship to be *approximately linear*,

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \varepsilon, \quad (1.2)$$

where ε is some error or noise term. The goal is to determine the *model parameters* β_0, \dots, β_p . To determine these, we can collect $n \geq p$ sample realizations (from observations or experiments),

$$Y = y_i, \quad X_1 = x_{i1}, \dots, X_p = x_{ip}, \quad 1 \leq i \leq n,$$

and assume that the data is related according to (1.2),

$$y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \varepsilon_i, \quad 1 \leq i \leq n.$$

Collecting the data in matrices and vectors,

$$\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix},$$

we can write the relationship concisely as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

We would then like to find $\boldsymbol{\beta}$ in such a way that the difference $\boldsymbol{\varepsilon} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$ is as *small* as possible. One way of measuring the size of a vector $\mathbf{x} \in \mathbb{R}^n$ is the square of its 2-norm, or Euclidean norm,

$$\|\mathbf{x}\|_2^2 = \mathbf{x}^\top \mathbf{x} = \sum_{i=1}^n x_i^2.$$

The best $\boldsymbol{\beta}$ is then the vector that solves the unconstrained optimization problem

$$\text{minimize } \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2.$$

This is an example of an optimization problem, with variables $\boldsymbol{\beta}$, no constraints (*all* $\boldsymbol{\beta}$ are valid candidates and the constraint set is $\Omega = \mathbb{R}^{p+1}$), and a *quadratic* objective function

$$f(\boldsymbol{\beta}) = \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 = (\mathbf{X}\boldsymbol{\beta} - \mathbf{y})^\top (\mathbf{X}\boldsymbol{\beta} - \mathbf{y}) = \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} + \mathbf{y}^\top \mathbf{y},$$

where \mathbf{X}^\top is the matrix transpose. As we will see later, quadratic functions are convex, so this is a convex optimization problem. This simple optimization problem has a *unique* closed form solution**,

$$\boldsymbol{\beta}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (1.3)$$

In practice one wouldn't compute $\boldsymbol{\beta}^*$ by evaluating [1], as there are more efficient methods available (see Lecture 2).

To illustrate the least squares setting using a concrete example, assume that we have data relating the basal metabolic rate (energy expenditure per time unit) in mammals to their mass.² The model we use is $Y = \beta_0 + \beta_1 X$, with Y the basal metabolic rate and X the mass. Using



²This example is from the episode “Size Matters” of the BBC series Wonders of Life.

data for 573 mammals from the PanTHERIA database³, we can assemble the vector \mathbf{y} and the matrix $\mathbf{X} \in \mathbb{R}^{n \times (p+1)}$ in order to compute the $\boldsymbol{\beta} = (\beta_0, \beta_1)^\top$. Here, $p = 1$ and $n = 573$.

We next illustrate how to solve this problem in Python. As usual, we first have to import some relevant libraries: **numpy** for numerical computation, **pandas** for loading and transforming datasets, **cvxpy** for convex optimization, and **matplotlib** for plotting.

```
In [1]: # Import some important Python modules
import numpy as np
import pandas as pd
from cvxpy import *
import matplotlib.pyplot as plt
```

We next have to load the data. The data is saved in a table with 573 rows and 2 columns, where the first column list the mass and the second the basal metabolic rate.

```
In [2]: # Load data into numpy array
bmr = pd.read_csv('../..../data/bmr.csv', header=None).as_matrix()
# We can find out the dimension of the data
bmr.shape
```

Out [2]: (573, 2)

To see the first three and the last three rows of the dataset, we can use the "print" command.

```
In [3]: print(bmr[0:3,:])
```

```
[[ 13.108   10.604 ]
 [  9.3918   8.2158]
 [ 10.366   9.3285]]
```

To visualise the whole dataset, we can make a scatterplot by interpreting each row as a coordinate on the plane, and marking it with a dot.

```
In [4]: # Display scatterplot of data (plot all the rows as points)
%matplotlib inline
bmrl = plt.plot(bmr[:,0],bmr[:,1], 'o')
plt.xlabel("Mass")
plt.ylabel("Basal metabolic rate")
plt.show()
```

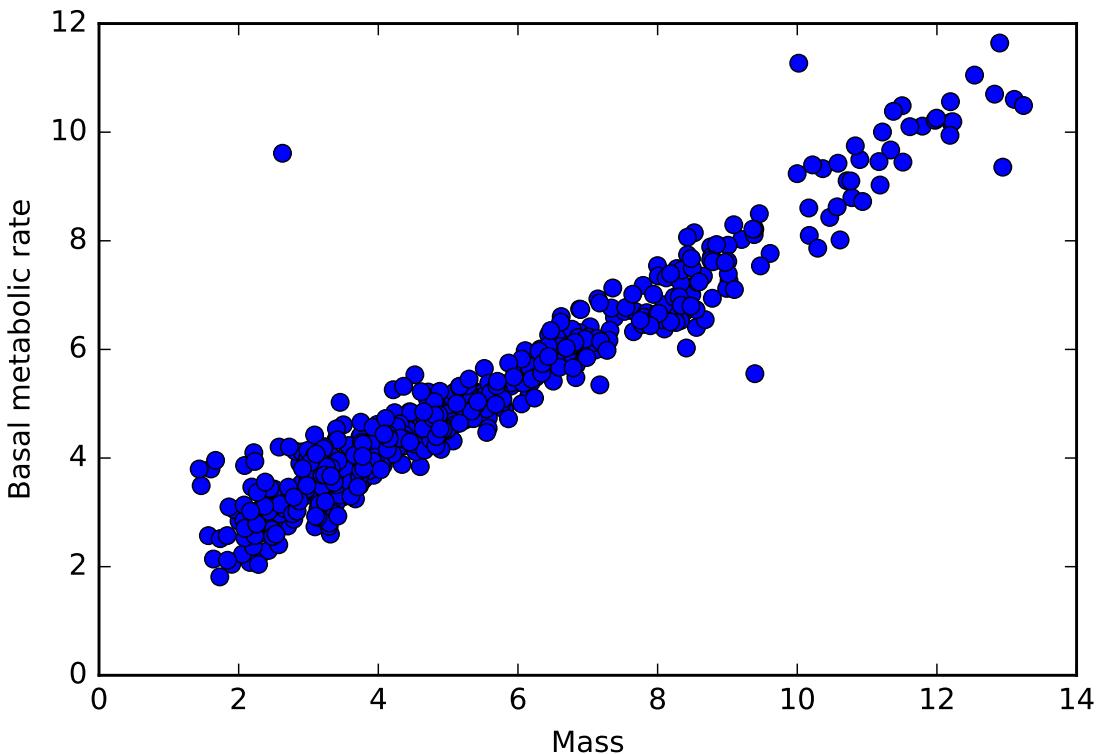
The plot above suggests that the relation of the basal metabolic rate to the mass is linear, i.e., of the form

$$Y = \beta_0 + \beta_1 X,$$

where X is the mass and Y the BMR. We can find β_0 and β_1 by solving an optimization problem as described above. We first have to assemble the matrix \mathbf{X} and the vector \mathbf{y} .

```
In [5]: n = bmr.shape[0]
p = 1
X = np.concatenate((np.ones((n,1)),bmr[:,0:p]),axis=1)
y = bmr[:, -1]
```

³<http://esapubs.org/archive/ecol/E090/184/#data>



```
In [6]: # Create a (p+1) vector of variables
Beta = Variable(p+1)

# Create sum-of-squares objective function
objective = Minimize(sum_entries(square(X*Beta - y)))

# Create problem and solve it
prob = Problem(objective)
prob.solve()

print("status: ", prob.status)
print("optimal value: ", prob.value)
print("optimal variables: ", Beta[0].value, Beta[1].value)
```

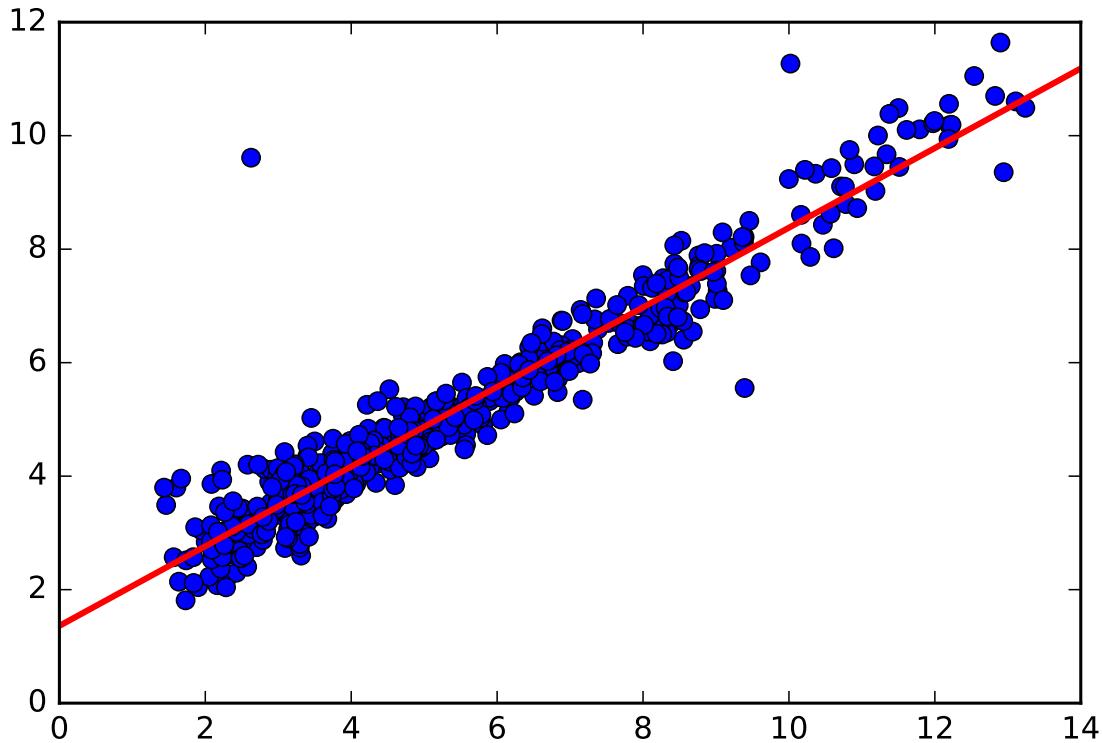
```
status:  optimal
optimal value:  152.736200529558
optimal variables:  1.3620698558275837  0.7016170245505547
```

Now that we solved the problem and have the values $\beta_0 = 1.362$ and $\beta_1 = 0.702$, we can plot the line and see how it fits the data.

```
In [6]: plt.plot(bmr[:,0],bmr[:,1],'o')

xx = np.linspace(0,14,100)
bmr = plt.plot(xx, Beta[0].value+Beta[1].value*xx, color='red', \
    linewidth=2)
plt.show()
```

Even though for illustration purposes we used the CVXPY package, this particular problem can be solved directly using the least squares solver in numpy.



```
In [7]: import numpy.linalg as la
beta = la.lstsq(X,y)
print(beta[0])
```

```
[ 1.36206997  0.70161692]
```

The above example is an example of a **machine learning** problem. In machine learning, one seeks to *learn* a function F mapping some inputs X to outputs Y , $Y = F(X)$. A few examples:

- X : economic data, Y : value of a stock;
- X : physiological data, Y : medical diagnosis;
- X : email, Y : 1 if email is spam, 0 otherwise;
- X : scanned image, Y : a letter represented by that image.

In *supervised learning* we have a set of sample input pairs, (y_i, x_i) , $1 \leq i \leq m$, and we typically try to find a function F that minimizes the **least squared error**,

$$\text{minimize} \quad \sum_{i=1}^m (\mathbf{y}_i - F(\mathbf{x}_i))^2,$$

where one minimizes over all functions F from some class. In the above example, we assumed our functions to be linear, in which case the can be parametrized by the coefficients β_0, \dots, β_p . As the course progresses, we will see examples of more sophisticated machine

learning problems, often with nonlinear objective function and other *loss functions* instead of the least square error.

Example 1.2. (Linear programming) Suppose a plane has two cargo compartments with weight capacities $C_1 = 35$ and $C_2 = 40$ tonnes, and volumes (space capacities) $V_1 = 250$ and $V_2 = 400$ cubic metres. Assume we have three types of cargo to transport, specified as follows.

	Volume (m^3 per tonne)	Weight (tonnes)	Profit (£/ tonne)
Cargo 1	8	25	£300
Cargo 2	10	32	£350
Cargo 3	7	28	£270

The problem is now to decide how much of each cargo to take on board, and how to distribute it in an optimal way among the two compartments.

1. The *decision variables* x_{ij} specify the amount, in tonnes, of cargo i to go into compartment j . We collect them in a vector \mathbf{x} .
2. The *objective function* is the total profit,

$$f(\mathbf{x}) = 300 \cdot (x_{11} + x_{12}) + 350 \cdot (x_{21} + x_{22}) + 270 \cdot (x_{31} + x_{32}).$$

3. The *constraints* are given by the space and weight limitations of the compartments, and the amount of cargo available.

$$\begin{aligned} x_{11} + x_{12} &\leq 25 && (\text{total amount of cargo 1}) \\ x_{21} + x_{22} &\leq 32 && (\text{total amount of cargo 2}) \\ x_{31} + x_{32} &\leq 28 && (\text{total amount of cargo 3}) \\ x_{11} + x_{21} + x_{31} &\leq 35 && (\text{weight constraint on compartment 1}) \\ x_{12} + x_{22} + x_{32} &\leq 40 && (\text{weight constraint on compartment 2}) \\ 8x_{11} + 10x_{21} + 7x_{31} &\leq 250 && (\text{volume constraint on compartment 1}) \\ 8x_{12} + 10x_{22} + 7x_{32} &\leq 400 && (\text{volume constraint on compartment 2}) \\ (x_{11} + x_{21} + x_{31})/35 - (x_{12} + x_{22} + x_{32})/40 &= 0 && (\text{maintain balance of weight ratio}) \\ x_{ij} &\geq 0 && (\text{cargo can't have negative weight}) \end{aligned}$$

It is customary to write the objective function as a scalar product, $f(\mathbf{x}) = \langle \mathbf{c}, \mathbf{x} \rangle := \mathbf{c}^\top \mathbf{x}$, and to express the constraints as systems of linear equations and inequalities using matrix-vector products,

$$\begin{aligned} &\text{maximize} && \langle \mathbf{c}, \mathbf{x} \rangle \\ &\text{subject to} && A\mathbf{x} \leq \mathbf{b} \\ & && B\mathbf{x} = \mathbf{d} \\ & && \mathbf{x} \geq 0 \end{aligned}$$

where the inequalities \geq and \leq are to be understood componentwise. This problem has a unique solution that can be found using CVXPY in Python,

```
In [8]: \\# Define all the matrices and vectors involved
c = np.array([300,300,350,350,270,270])
A = np.array([[1, 1, 0, 0, 0, 0],
              [0, 0, 1, 1, 0, 0],
              [0, 0, 0, 0, 1, 1],
              [1, 0, 1, 0, 1, 0],
              [0, 1, 0, 1, 0, 1],
              [8, 0, 10, 0, 7, 0],
              [0, 8, 0, 10, 0, 7]])
b = np.array([25,32,28,35,40,250,400]);
B = np.array([1/35, -1/40, 1/35, -1/40, 1/35, -1/40]);
d = np.zeros(1)

# Create variables, objective and constraints
x = Variable(6)
constraints = [A*x <= b, B*x == d, x >= 0]
objective = Maximize(c*x)

# Create a problem using the objective and constraints and solve it
prob = Problem(objective, constraints)
prob.solve()

print("Solution found: \n", np.round(np.abs(x.value), \
    decimals=2).transpose())
```

```
Solution found:
[[ 6.75  7.71   0.    32.    28.    0. ]]
```

The solution found is

$$x_{11} = 6.7500, x_{12} = 7.7143, x_{21} = 0, x_{22} = 32, x_{31} = 28, x_{32} = 0.$$

We made some simplifying assumptions, for example that the cargo can be split up into arbitrary fractions. Additional work is required to resolve these issues. Problems of this kind are known as **linear programming**, because the objective function and the constraints are given by linear functions. Such problems can be solved efficiently using the simplex algorithm or interior point methods. The highly developed theory of linear programming acts as a template for more general convex optimization that is developed in this course.

Example 1.3. (Image inpainting) Optimization methods play an increasingly important role in image and signal processing. An image can be viewed as an $m \times n$ matrix \mathbf{U} , with each entry u_{ij} corresponding to a light intensity (for greyscale images), or a colour vector, represented by a triple of red, green and blue intensities (usually with values between 0 and 255 each). For simplicity the following discussion assumes a greyscale image. For computational purposes, the matrix of an image is often viewed as an mn -dimensional vector \mathbf{u} , with the columns of the matrix stacked on top of each other.

In the *image inpainting* problem, one aims to *guess* the true value of missing or corrupted entries of an image. There are different approaches to this problem. A conceptually simple approach is to replace the image with the *closest* image among a set of images satisfying typical properties. But what are typical properties of a typical image? Some properties that come to mind are:

- Images tend to have large homogeneous areas in which the colour doesn't change much;

- Images have approximately low rank, when interpreted as matrices.

Total variation image analysis takes advantage of the first property. The **total variation** or TV-norm is the sum of the norm of the horizontal and vertical differences,

$$\|\mathbf{U}\|_{\text{TV}} = \sum_{i=1}^m \sum_{j=1}^n \sqrt{(u_{i+1,j} - u_{i,j})^2 + (u_{i,j+1} - u_{i,j})^2},$$

where we set entries with out-of-bounds indices to 0. The TV-norm naturally increases with increased variation or sharp edges in an image. Consider for example the two following matrices (imagine that they represent a 3×3 pixel block taken from an image).

$$\mathbf{U}_1 = \begin{pmatrix} 0 & 17 & 3 \\ 7 & 32 & 0 \\ 2 & 9 & 27 \end{pmatrix}, \quad \mathbf{U}_2 = \begin{pmatrix} 1 & 1 & 3 \\ 1 & 0 & 0 \\ 0 & 0 & 2 \end{pmatrix}$$

The left matrix has TV-norm $\|\mathbf{U}_1\|_{\text{TV}} = 200.637$, while the right one has TV-norm $\|\mathbf{U}_2\|_{\text{TV}} = 14.721$ (verify this!) Intuitively, we would expect a natural image with artifacts added to it to have a higher TV norm.

Now let \mathbf{U} be an image with entries u_{ij} , and let $\Omega \subset [m] \times [n] = \{(i, j) : 1 \leq i \leq m, 1 \leq j \leq n\}$ be the set of indices where the original image and the corrupted image coincide (all the other entries are missing). One could attempt to find the image with the *smallest* TV-norm that coincides with the known pixels u_{ij} for $(i, j) \in \Omega$. This is an optimization problem of the form

$$\text{minimize } \|\mathbf{X}\|_{\text{TV}} \quad \text{subject to } x_{ij} = u_{ij} \text{ for } (i, j) \in \Omega.$$

The TV-norm is an example of a convex function and the constraints are linear conditions which define a convex set. This is again an example of a **convex optimization problem** and can be solved efficiently by a range of algorithms. For the time being we will not go into the algorithms but solve it using CVXPY. The example below is based on an example from the CVXPY Tutorial⁴, and it is recommended to look at this tutorial for other interesting examples!

Warning: the example below uses some more advanced Python programming, it is not necessary to understand the details at this point.

In our first piece of code below, we load the image and a version of the image with text written on it, and display the images. The **Python Image Library (PIL)** is used for this purpose.

⁴<http://www=cvxpy.org/en/latest/tutorial/index.html>

```
In [9]: from PIL import Image

# Load the images and convert to numpy arrays for processing.
U = np.array(Image.open("../images/alanturing.png"))
Ucorr = np.array(Image.open("../images/alanturing-corr.png"))

# Display the images
%matplotlib inline
fig, ax = plt.subplots(1, 2, figsize=(10, 5))

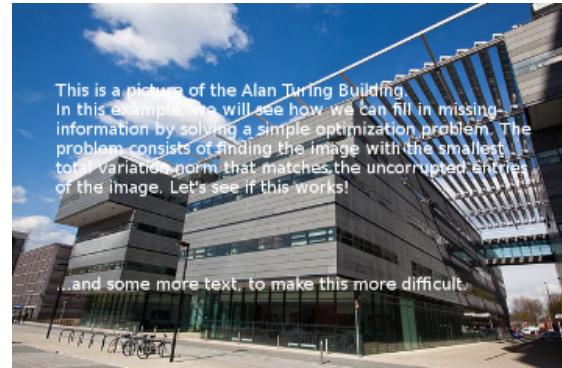
ax[0].imshow(U);
ax[0].set_title("Original Image")
ax[0].axis('off')

ax[1].imshow(Ucorr);
ax[1].set_title("Corrupted Image")
ax[1].axis('off');
```

Original Image



Corrupted Image



After having the images at our disposal, we determine which entries of the corrupted image are known. We store these in a *mask* M , with entries $m_{ijk} = 1$ if the colour k of the (i, j) -th pixel is known, and 0 otherwise.

```
In [10]: # Each image is now an m x n x 3 array, with each pixel
# represented by three numbers between 0 and 255,
# corresponding to red, green and blue
rows, cols, colours = U.shape

# Create a mask: this is a matrix with a 1 if the corresponding
# pixel is known, and zero else
M = np.zeros((rows, cols, colours))
for i in range(rows):
    for j in range(cols):
        for k in range(colours):
            if U[i, j, k] == Ucorr[i, j, k]:
                M[i, j, k] = 1
```

We are now ready to solve the optimization problem using CVXPY. As the problem is rather big ($400 \times 600 \times 3 = 720000$ variables), it is important to choose a good solver that will solve the problem to sufficient accuracy in an acceptable amount of time. For the example at hand, we choose the SCS solver, which can be specified when calling the **solve** function.

```
In [11]: # Determine the variables and constraints
variables = []
constraints = []
for k in range(colours):
    X = Variable(rows, cols)
    # Add variables
    variables.append(X)
    # Add constraints by multiplying the relevant variable matrix
    # elementwise with the mask
    constraints.append(mul_elemwise(M[:, :, k], X) ==
        \ (M[:, :, k], Ucorr[:, :, k]))

# Create a problem instance with
objective = Minimize(tv(variables[0], variables[1], variables[2]))

# Create a problem instance and solve it using the SCS solver
prob = Problem(objective, constraints)
prob.solve(verbose=True, solver=SCS)
```

Out 8263910.812250629

[11]:

Now that we solved the optimization problem, we have a solution stored in 'variables'. We have to transform this back into an image and display the result.

```
In [12]: %matplotlib inline

# Load variable values into a single array.
Urec = np.zeros((rows, cols, colours), dtype=np.uint8)
for i in range(colours):
    Urec[:, :, i] = variables[i].value

fig, ax = plt.subplots(1, 2, figsize=(10, 5))

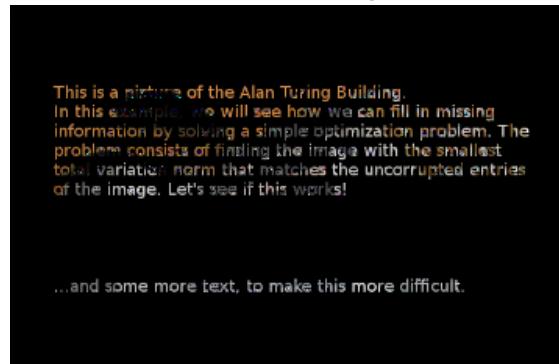
# Display the inpainted image.
ax[0].imshow(Urec);
ax[0].set_title("Inpainted Image")
ax[0].axis('off')

ax[1].imshow(np.abs(Ucorr[:, :, 0:3] - Urec));
ax[1].set_title("Difference Image")
ax[1].axis('off');
```

Inpainted Image



Difference Image



Another typical structure of images is that the **singular values** of the image, considered as matrix, decay quickly. The **singular value decomposition** (SVD) of a matrix $A \in \mathbb{R}^{m \times n}$ is the

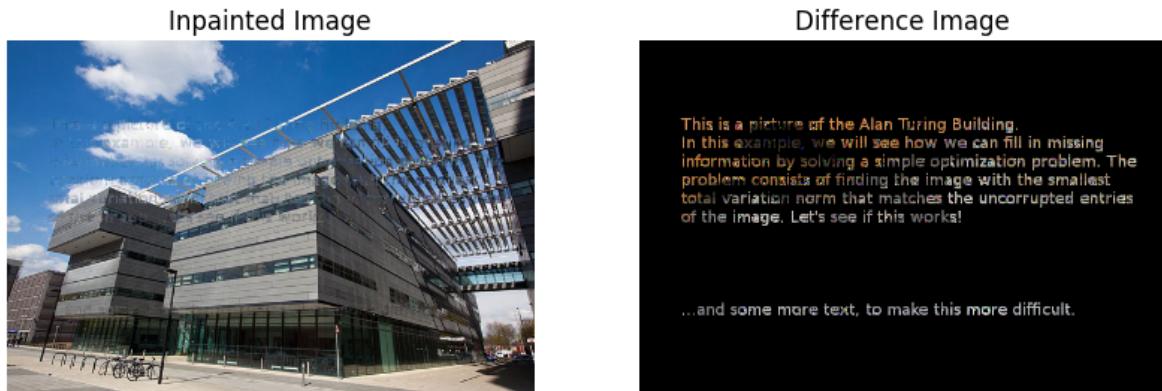
matrix product

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T,$$

where $\mathbf{U} \in \mathbb{R}^{m \times m}$ and $\mathbf{V} \in \mathbb{R}^{n \times n}$ are orthogonal matrices, and $\Sigma \in \mathbb{R}^{m \times n}$ is a diagonal matrix with entries $\sigma_1, \dots, \sigma_{\min\{m,n\}}$ on the diagonal. Instead of minimizing the TV-norm of an image \mathbf{X} , one may instead try to minimize the **Schatten 1-norm**, defined as the sum of the singular values, $\|\mathbf{U}\|_{S_1} = \sigma_1 + \dots + \sigma_{\min\{m,n\}}$. The problem is then

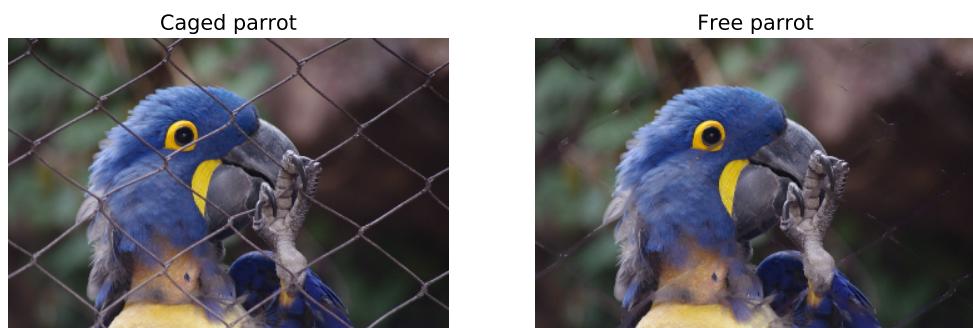
$$\text{minimize } \|\mathbf{X}\|_{S_1} \quad \text{subject to } x_{ij} = u_{ij} \text{ for } (i, j) \in \Omega.$$

As we will see towards the end of the course, this is an instance of a type of convex optimization problem known as **semidefinite programming**. Luckily, CVXPY includes the Schatten 1-norm (also known as nuclear norm) as valid objective function, so we don't have to deal with the details of this problem. As the problem is computationally intensive, we just reproduce the result.



In this example, the result appears worse as in the problem involving the TV-norm. Alternatively, one may also use the 1-norm of the image applied to a discrete cosine transform (DCT) or a discrete wavelet transform (DWT).

Of course, one could run the above examples for fun with different types of images in an attempt to get rid of certain parts. The image below show the result of applying the total variation inpainting procedure to set a parrot free.



Part II

Unconstrained Optimization

Lecture 2

In this lecture we will study the unconstrained problem

$$\text{minimize } f(\mathbf{x}), \quad (2.1)$$

where $\mathbf{x} \in \mathbb{R}^n$. Optimality conditions aim to identify properties that potential minimizers need to satisfy in relation to $f(\mathbf{x})$. We will review the well known local optimality conditions for differentiable functions from calculus. We then introduce convex functions and discuss some of their properties.

2.1 Unconstrained optimization

Solutions to (2.1) come in different flavours, as in the following definition.

Definition 2.1. A point $\mathbf{x}^* \in \mathbb{R}^n$ is a

- *global minimizer* of (2.1) if for all $\mathbf{x} \in \mathbb{R}^n$, $f(\mathbf{x}^*) \leq f(\mathbf{x})$;
- a *local minimizer*, if there is an open neighbourhood U of \mathbf{x}^* such that $f(\mathbf{x}^*) \leq f(\mathbf{x})$ for all $\mathbf{x} \in U$;
- a *strict local minimizer*, if there is an open neighbourhood U of \mathbf{x}^* such that $f(\mathbf{x}^*) < f(\mathbf{x})$ for all $\mathbf{x} \in U$;
- an *isolated minimizer* if there is an open neighbourhood U of \mathbf{x}^* such that \mathbf{x}^* is the only local minimizer in U .

Without any further assumptions on f , finding a minimizer is a hopeless task: we simply can't examine the function at *all* points in \mathbb{R}^n . The situation becomes more tractable if we assume some *smoothness* conditions. Recall that $C^k(U)$ denotes the set of functions that are k times continuously differentiable on some set U . The following *first-order* necessary condition for optimality is well known. We write $\nabla f(\mathbf{x})$ for the gradient of f at \mathbf{x} , i.e., the vector

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f}{\partial x_1}(\mathbf{x}), \dots, \frac{\partial f}{\partial x_n}(\mathbf{x}) \right)^\top$$

Theorem 2.2. Let \mathbf{x}^* be a local minimizer of f and assume that $f \in C^1(U)$ for a neighbourhood of U of \mathbf{x}^* . Then $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

There are simple examples that show that this is not a sufficient condition: maxima and saddle points will also have a vanishing gradient. If we have access to *second-order information*, in form of the second derivative, or Hessian, of f , then we can say more. Recall that the Hessian of f at \mathbf{x} , $\nabla^2 f(\mathbf{x})$, is the $d \times d$ symmetric matrix given by the second derivatives,

$$\nabla^2 f(\mathbf{x}) = \left(\frac{\partial^2 f}{\partial x_i \partial x_j} \right)_{1 \leq i, j \leq n}.$$

In the one-variable case we have learned that if x^* is a local minimizer of $f \in C^2([a, b])$, then $f'(x^*) = 0$ and $f''(x^*) \geq 0$. Moreover, the conditions $f'(x^*) = 0$ and $f''(x^*) > 0$ guarantee that we have a local minimizer. These conditions generalise to higher dimension, but first we need to know what $f''(x) > 0$ when we have more than one variable.

Recall also that a matrix A is **positive semidefinite**, written $A \succeq 0$, if for every $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{x}^\top A \mathbf{x} \geq 0$, and positive definite, written $A \succ 0$, if $\mathbf{x}^\top A \mathbf{x} > 0$. The property that the Hessian matrix is positive semidefinite is a multivariate generalization of the property that the second derivative is nonnegative. The known conditions for a minimizer involving the second derivative generalize accordingly.

Theorem 2.3. *Let $f \in C^2(U)$ for some open set U and $\mathbf{x}^* \in U$. If \mathbf{x}^* is a local minimizer, then $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}^*)$ is positive semidefinite. Conversely, if $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}^*)$ is positive definite, then \mathbf{x}^* is a strict local minimizer.*

Unfortunately, the above criteria are not able to identify global minimizers, as differentiability is a local property.

2.2 Convex functions

We now come to the central notion of this course.

Definition 2.4. A set $C \subseteq \mathbb{R}^n$ is **convex** if for all $\mathbf{x}, \mathbf{y} \in C$ and $\lambda \in [0, 1]$, the line $\lambda\mathbf{x} + (1 - \lambda)\mathbf{y} \in C$. A **convex body** is a convex set that is closed and bounded.

Definition 2.5. Let $S \subseteq \mathbb{R}^n$. A function $f: S \rightarrow \mathbb{R}$ is called *convex* if S is convex and for all $\mathbf{x}, \mathbf{y} \in S$ and $\lambda \in [0, 1]$,

$$f(\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}).$$

The function f is called *strictly convex* if

$$f(\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}) < \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}).$$

A function f is called *concave*, if $-f$ is convex.

Figure 2.2 illustrates how a convex function of one variable looks like. The graph of the function lies below any line connecting two points on it.

Convex functions have pleasant properties, while at the same time covering many of the functions that arise in applications. Perhaps the most important property is that local minima are global minima.

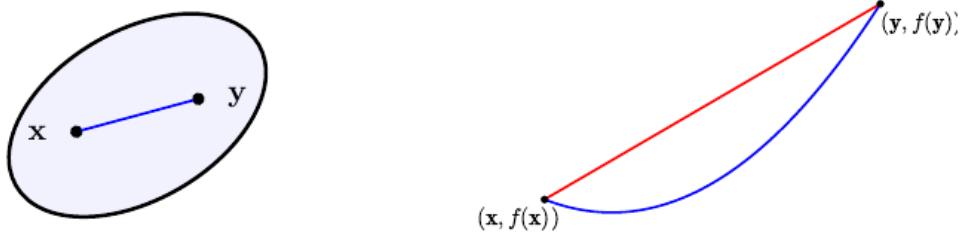


Figure 2.2: A convex set and a convex function

Theorem 2.6. Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function. Then any local minimizer of f is a global minimizer.

Proof. Let \mathbf{x}^* be a local minimizer and assume that it is not a global minimizer. Then there exists a vector $\mathbf{y} \in \mathbb{R}^d$ such that $f(\mathbf{y}) < f(\mathbf{x}^*)$. Since f is convex, for any $\lambda \in [0, 1]$ and $\mathbf{x} = \lambda\mathbf{y} + (1 - \lambda)\mathbf{x}^*$ we have

$$f(\mathbf{x}) \leq \lambda f(\mathbf{y}) + (1 - \lambda)f(\mathbf{x}^*) < \lambda f(\mathbf{x}^*) + (1 - \lambda)f(\mathbf{x}^*) = f(\mathbf{x}^*).$$

This holds for all \mathbf{x} on the line segment connecting \mathbf{y} and \mathbf{x}^* . Since every open neighbourhood U of \mathbf{x}^* contains a bit of this line segment, this means that every open neighbourhood U of \mathbf{x}^* contains an $\mathbf{x} \neq \mathbf{x}^*$ such that $f(\mathbf{x}) \leq f(\mathbf{x}^*)$, in contradiction to the assumption that \mathbf{x}^* is a local minimizer. It follows that \mathbf{x}^* has to be a global minimizer. \square

Remark 2.7. Note that in the above theorem we made no assumptions about the differentiability of the function f ! In fact, while a convex function is always *continuous*, it need not be differentiable. The function $f(x) = |x|$ is a typical example: it is convex, but not differentiable at $x = 0$.

Example 2.8. Affine functions $f(\mathbf{x}) = \langle \mathbf{x}, \mathbf{a} \rangle + b$ and the exponential function e^x are examples of convex functions.

Example 2.9. In optimization we will often work with functions of matrices, where an $m \times n$ matrix is considered as a vector in $\mathbb{R}^{m \times n} \cong \mathbb{R}^{mn}$. If the matrix is symmetric, that is, if $\mathbf{A}^\top = \mathbf{A}$, then we only care about the upper diagonal entries, and we consider the space \mathcal{S}^n of symmetric matrices as a vector space of dimension $n(n + 1)/2$ (the number of entries on and above the main diagonal). Important functions on symmetric matrices that are convex are the operator norm $\|\mathbf{A}\|_2$, defined as

$$\|\mathbf{A}\|_2 := \max_{\mathbf{x}: \|\mathbf{x}\|_2 \leq 1} \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{x}\|_2},$$

or the function $\log \det(\mathbf{X})$, defined on the set of *positive semidefinite* symmetric matrices \mathcal{S}_+^n .

There are useful ways of characterising convexity using differentiability.

Theorem 2.10. 1. Let $f \in C^1(\mathbb{R}^n)$. Then f is convex if and only if for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}).$$

2. Let $f \in C^2(\mathbb{R}^n)$. Then f is convex if and only if $\nabla^2 f(\mathbf{x})$ is positive semidefinite. If $\nabla^2 f(\mathbf{x})$ is positive definite, then f is strictly convex.

Example 2.11. Consider a quadratic function of the form

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} + \mathbf{b}^\top \mathbf{x} + c,$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric. Writing out the product, we get

$$\begin{aligned} \mathbf{x}^\top \mathbf{A} \mathbf{x} &= (x_1 \ \dots \ x_n) \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \\ &= (x_1 \ \dots \ x_n) \begin{pmatrix} a_{11}x_1 + \cdots + a_{1n}x_n \\ \vdots \\ a_{n1}x_1 + \cdots + a_{nn}x_n \end{pmatrix} = \sum_{i=1}^n \sum_{j=1}^n a_{ij}x_i x_j. \end{aligned}$$

Because \mathbf{A} is symmetric, we have $a_{ij} = a_{ji}$, and the above product simplifies to

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} = \sum_{i=1}^n a_{ii}x_i^2 + 2 \sum_{1 \leq i < j \leq n} a_{ij}x_i x_j.$$

This is a quadratic function, because it involves products of the x_i . The gradient and the Hessian of $f(\mathbf{x})$ are found by computing the partial derivatives of f :

$$\frac{\partial f}{\partial x_i} = \sum_{j=1}^n a_{ij}x_j + b_i, \quad \frac{\partial^2 f}{\partial x_i \partial x_j} = a_{ij}.$$

In summary, we have

$$\nabla f(\mathbf{x}) = \mathbf{A} \mathbf{x} + \mathbf{b}, \quad \nabla^2 f(\mathbf{x}) = \mathbf{A}.$$

Using the previous theorem, we see that f is convex **if and only if** \mathbf{A} is positive semidefinite. A typical example for such a function is

$$f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|^2 = (\mathbf{Ax} - \mathbf{b})^\top (\mathbf{Ax} - \mathbf{b}) = \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} - 2\mathbf{b}^\top \mathbf{A} \mathbf{x} + \mathbf{b}^\top \mathbf{b}.$$

The matrix $\mathbf{A}^\top \mathbf{A}$ is always symmetric and positive semidefinite (why?) so that the function f is convex.

A convenient way to visualise a function $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ is through **contour plots**. A **level set** of the function f is a set of the form

$$\{\mathbf{x} : f(\mathbf{x}) = c\},$$

where c is the **level**. Each such level set is a curve in \mathbb{R}^2 , and a contour plot is a plot of a collection of such curves for various c . If one colours the areas between adjacent curves, one

gets a plot as in the following figure. A *convex function* has the property that there is only one *sink* in the contour plot.

```
In [16]: # import numpy as np
import numpy.linalg as la
import matplotlib.pyplot as plt
%matplotlib inline

# Create random data: we use the randn function
X = np.random.randn(3,2)
y = np.random.randn(3)

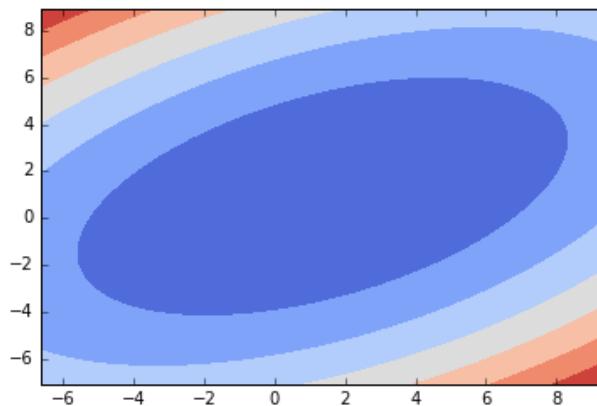
# Solve least squares problem minimize \|X\beta-y\|^2
# the index 0 says that we get the first component of the solution
# (the function lstsq give more output than just the beta vector)
beta = la.lstsq(X,y)[0]

# Create function and plot the contours
def f(a,b):
    return sum((a*X[:,0]+b*X[:,1]-y)**2)

# Find the "right" boundaries around the minimum
xx = np.linspace(beta[0]-8,beta[0]+8,100)
yy = np.linspace(beta[1]-8,beta[1]+8,100)
XX, YY = np.meshgrid(xx,yy)

Z = np.zeros(XX.shape)
for i in range(Z.shape[0]):
    for j in range(Z.shape[1]):
        Z[i,j] = f(XX[i,j],YY[i,j])

cmap = plt.cm.get_cmap("coolwarm")
plt.contourf(XX,YY,Z, cmap=cmap)
plt.show()
```



Lecture 3

Most modern optimization methods are iterative: they generate a sequence of points $\mathbf{x}_0, \mathbf{x}_1, \dots$ in \mathbb{R}^n in the hope that this sequences will converge to a local or global minimizer \mathbf{x}^* of a function $f(\mathbf{x})$. A typical rule for generating such a sequence would be to start with a vector \mathbf{x}_0 , chosen by an educated guess, and then for $k \geq 0$, move from step k to $k + 1$ by

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k,$$

in a way that ensures that $f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k)$. The parameter α_k is called the **step length**, while \mathbf{p}_k is the **search direction**. In this lecture we discuss one such method, the method of gradient descent, or steepest descent, and discuss how to select the right step length.

3.1 Gradient descent

In the method of gradient descent, the search direction is chosen as

$$\mathbf{p}_k = -\nabla f(\mathbf{x}_k). \quad (3.1)$$

To see why this makes sense, let \mathbf{p} be a direction and consider the Taylor expansion

$$f(\mathbf{x}_k + \alpha \mathbf{p}) = f(\mathbf{x}_k) + \alpha \langle \mathbf{p}, \nabla f(\mathbf{x}_k) \rangle + O(\alpha^2).$$

Considering this as a function of α , the rate of change in direction \mathbf{p} at \mathbf{x}_k is the derivative of this function at $\alpha = 0$,

$$\frac{df(\mathbf{x}_k + \alpha \mathbf{p})}{d\alpha}|_{\alpha=0} = \langle \mathbf{p}, \nabla f(\mathbf{x}_k) \rangle,$$

also known as the **directional derivative** of f at \mathbf{x}_k in the direction \mathbf{p} . This formula indicates that the rate of change is *negative*, and we have a **descent direction**, if $\langle \mathbf{p}, \nabla f(\mathbf{x}_k) \rangle < 0$.

The Cauchy-Schwarz inequality (see Preliminaries, Page 9) gives the bounds

$$-\|\mathbf{p}\|_2 \|\nabla f(\mathbf{x}_k)\|_2 \leq \langle \mathbf{p}, \nabla f(\mathbf{x}_k) \rangle \leq \|\mathbf{p}\|_2 \|\nabla f(\mathbf{x}_k)\|_2.$$

We see that the rate of change is the smallest when the first inequality is an equality, which happens if

$$\mathbf{p} = -\alpha \nabla f(\mathbf{x}_k)$$

for some $\alpha > 0$.

For a visual interpretation of what it means to be a descent direction, note that the **angle** θ between a vector p and the gradient $\nabla f(\mathbf{x})$ at a point \mathbf{x} is given by (see Preliminaries, Page 9)

$$\langle \mathbf{x}, \nabla f(\mathbf{x}) \rangle = \|p\|_2 \|\nabla f(\mathbf{x})\|_2 \cos(\theta).$$

This is negative if the vector p forms an angle greater than $\pi/2$ with the gradient. Recall that the gradient points in the direction of steepest ascent, and is orthogonal to the *level sets*. If you are standing on the slope of a mountain, walking along the level set lines will not change your elevation, the gradient points to the steepest upward direction, and the negative gradient to the steepest descent.

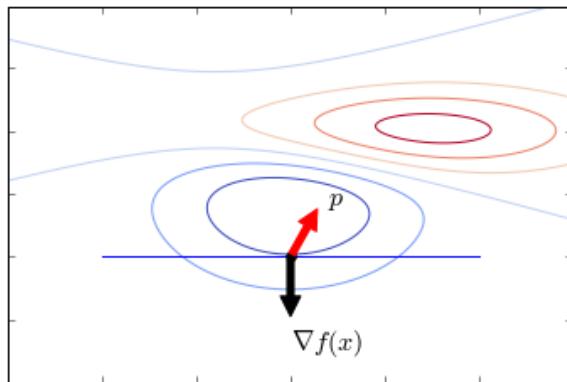


Figure 3.3: A descent direction

Any multiple $\alpha \nabla f(\mathbf{x}_k)$ points in the direction of steepest descent, but we have to choose a sensible parameter α to ensure that we make sufficient progress, but at the same time don't overshoot. Ideally, we would choose the value α_k that minimizes $f(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k))$. While finding such a minimizer is in general not easy (see Section Lecture 4 for alternatives), for quadratic functions it can be given in closed form.

Linear least squares

Consider a function of the form

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2.$$

In Problem Sheet 1 you will show that the Hessian is symmetric and positive semidefinite, with the gradient given by

$$\nabla f(\mathbf{x}) = \mathbf{A}^\top (\mathbf{Ax} - \mathbf{b}).$$

The method of gradient descent proceeds as

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{A}^\top (\mathbf{Ax}_k - \mathbf{b}).$$

To find the best α_k , we compute the minimum of the function

$$\alpha \mapsto f(\mathbf{x}_k - \alpha \mathbf{A}^\top (\mathbf{A}\mathbf{x}_k - \mathbf{b})). \quad (3.2)$$

If we set $\mathbf{r}_k := \mathbf{A}^\top (\mathbf{b} - \mathbf{A}\mathbf{x}_k) = -\nabla f(\mathbf{x}_k)$ and compute the minimum of (3.2) by differentiating, we get the step length

$$\alpha_k = \frac{\mathbf{r}_k^\top \mathbf{r}_k}{\mathbf{r}_k^\top \mathbf{A}^\top \mathbf{A} \mathbf{r}_k} = \frac{\|\mathbf{r}_k\|_2^2}{\|\mathbf{A} \mathbf{r}_k\|_2^2}.$$

(Verify this!) Note also that when we have \mathbf{r}_k and α_k , we can compute the next \mathbf{r}_k as

$$\begin{aligned} \mathbf{r}_{k+1} &= \mathbf{A}^\top (\mathbf{b} - \mathbf{A}\mathbf{x}_{k+1}) \\ &= \mathbf{A}^\top (\mathbf{b} - \mathbf{A}(\mathbf{x}_k + \alpha_k \mathbf{r}_k)) \\ &= \mathbf{A}^\top (\mathbf{b} - \mathbf{A}\mathbf{x}_k - \alpha_k \mathbf{A}^\top \mathbf{A} \mathbf{r}_k) = \mathbf{r}_k - \alpha_k \mathbf{A}^\top \mathbf{A} \mathbf{r}_k. \end{aligned}$$

The gradient descent algorithm for the linear least squares problem proceeds by first computing $\mathbf{r}_0 = \mathbf{A}^\top (\mathbf{b} - \mathbf{A}\mathbf{x}_0)$, and then at each step

$$\begin{aligned} \alpha_k &= \frac{\mathbf{r}_k^\top \mathbf{r}_k}{\mathbf{r}_k^\top \mathbf{A}^\top \mathbf{A} \mathbf{r}_k} \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \alpha_k \mathbf{r}_k \\ \mathbf{r}_{k+1} &= \mathbf{r}_k - \alpha_k \mathbf{A}^\top \mathbf{A} \mathbf{r}_k. \end{aligned}$$

Does this work? How do we know when to stop? It is worth noting that the residual satisfies $\mathbf{r} = 0$ if and only if \mathbf{x} is a stationary point, in our case, a minimizer. One criteria for stopping could then be to check whether $\|\mathbf{r}_k\|_2 \leq \varepsilon$ for some given tolerance $\varepsilon > 0$. One potential problem with this criterion is that the function can become *flat* long before reaching a minimum, so an alternative stopping method would be to stop when the difference between two successive points, $\|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2$, becomes smaller than some $\varepsilon > 0$.

Example 3.1. We plot the trajectory of gradient descent with the data

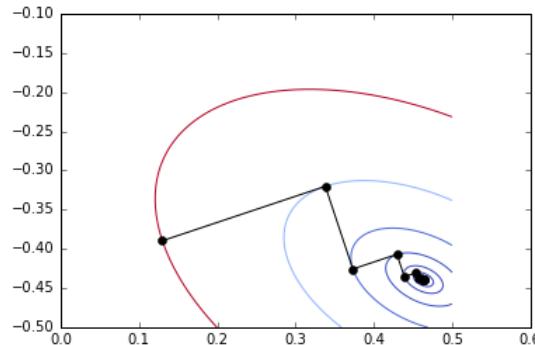


Figure 3.4: Trajectory of gradient descent

$$\mathbf{A} = \begin{pmatrix} 2 & 0 \\ 1 & 3 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix}.$$

As can be seen from the plot, we always move in the direction orthogonal to a level set, and stop at a point where we are tangent to a level set.

Lecture 4

Iterative algorithms for solving a problem of the form

$$\text{minimize } f(\mathbf{x}) \quad (4.1)$$

on \mathbb{R}^n generate a sequence of vectors $\mathbf{x}_0, \mathbf{x}_1, \dots$ in the hope that this sequence converges to a (local or global) minimizer \mathbf{x}^* of (4.1). In this lecture we first study step length selection procedures and then study what it means for a sequence to converge, and how to quantify the speed of convergence.

4.1 Step length selection

When moving in a descent direction \mathbf{p}_k (not necessarily the steepest descent direction),

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$$

we can only guarantee that the function value will decrease if the step length α_k is small enough. If we move too far in a descent direction, we might even land at a point where f is larger than where we started! It is therefore important to choose a step length that

- is not too small (so that the algorithm does not take too long);
- is not too large (so that we don't end up at a point with larger function value);
- is easy to compute.

An optimal step α for a descent method would be the minimizer of the function

$$\alpha \mapsto \varphi(\alpha) := f(\mathbf{x}_k + \alpha \mathbf{p}_k).$$

In practice, minimizing this function is not always the most efficient thing to do (or even possible). One would rather choose a step length that satisfies some criteria that ensure that the sequence \mathbf{x}_k converges to a minimizer \mathbf{x}^* under suitable conditions on a function f . One such condition is a **sufficient decrease condition**,

$$f(\mathbf{x}_k + \alpha \mathbf{p}_k) \leq f(\mathbf{x}_k) + c\alpha \langle \nabla f(\mathbf{x}_k), \mathbf{p}_k \rangle =: \ell(\alpha),$$

with $c \in (0, 1)$. Note that $\varphi'(0) = \langle \nabla f(\mathbf{x}_k), \mathbf{p}_k \rangle < 0$, because \mathbf{p}_k is a descent direction. The function $\ell(\alpha)$ is therefore a line through $f(\mathbf{x}_k)$ with a slope $c\langle \nabla f(\mathbf{x}_k), \mathbf{p}_k \rangle = c\varphi'(0) > \varphi'(0)$.

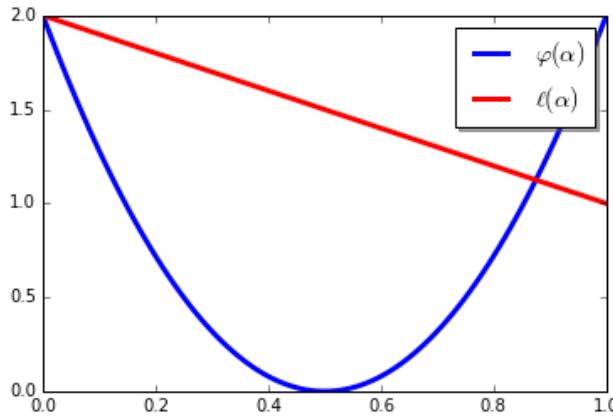


Figure 4.5: The sufficient decrease condition

To see why this condition is necessary, consider the function $f(x) = x^2 - 1$ and the sequence $x_k = \sqrt{1 + 1/k}$ for $k \geq 1$. Clearly, the sequence $f(x_k) = 1/k$ decreases, but fails to converge to the minimizer $f(0) = -1$.

The sufficient decrease condition (also called Armijo condition) can always be satisfied if α is chosen small enough, but the algorithm may become very slow. It is therefore common to supplement the sufficient descent condition with other criteria that guarantee that sufficient progress is made. Two of the commonly used criteria are:

- the Wolfe conditions, which add a *curvature condition*

$$\varphi'(\alpha) \geq \tilde{c}\varphi'(0)$$

for some $\tilde{c} \in (c, 1)$, which gives a lower bound on the slope of the new point;

- the Armijo-Goldstein conditions, which state that a step length α_k should additionally satisfy bound

$$f(\mathbf{x}_k) + (1 - c)\alpha_k \langle \nabla f(\mathbf{x}_k), \mathbf{p}_k \rangle \leq f(\mathbf{x}_k + \alpha_k \mathbf{p}_k), \quad (4.1)$$

which gives a lower bound on the step size.

Another common approach is **backtracking**: in this method one uses a high initial value of α (for example, $\alpha = 1$), and then decreases it until the sufficient descent condition is satisfied.

Example 4.1. Consider the function $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by $f(\mathbf{x}) = x_1^2 + x_2^2$. The gradient is $\nabla f(\mathbf{x}) = 2\mathbf{x}$, and the φ function at $\mathbf{x}_k = (1, 1)^\top$

$$\varphi(\alpha) = f(\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k)) = 2(1 - 2\alpha)^2, \quad \varphi'(\alpha) = -8(1 - 2\alpha).$$

The Armijo-Goldstein conditions (4.1) then state that we can choose α such that

$$2(1 - 4(1 - c)\alpha) \leq 2(1 - 2\alpha)^2 \leq 2(1 - 4c\alpha).$$

The optimal step length in this case would be $\alpha = 0.5$.

4.2 Convergence of iterative methods

A sequence of vectors $\{\mathbf{x}_k\}$ in \mathbb{R}^n , $k \geq 0$, converges to a vector \mathbf{x}^* with respect to a norm $\|\cdot\|$ as $k \rightarrow \infty$, written $\mathbf{x}_k \rightarrow \mathbf{x}$, if the sequence of numbers $\|\mathbf{x}_k - \mathbf{x}^*\|$ converges to zero. More formally, if for every $\varepsilon > 0$ there exists an index N such that for all $n \geq N$,

$$\|\mathbf{x}_n - \mathbf{x}^*\| < \varepsilon.$$

Iterative algorithms will rarely find the exact solution to a problem like (4.1), so we will usually be happy to find a solution that differs from the true one by at most some specified accuracy.

Definition 4.2. A sequence of vectors $\{\mathbf{x}_k\}$, $k \geq 0$, is said to converge to \mathbf{x}^*

- (a) linearly (or Q-linear, Q for quotient), if there exist an $r \in (0, 1)$ such that for sufficiently large k ,

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq r \|\mathbf{x}_k - \mathbf{x}^*\|.$$

- (b) superlinearly, if

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = 0,$$

- (c) with order p , if there exists a constant $M > 0$, such that for sufficiently large k ,

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq M \|\mathbf{x}_k - \mathbf{x}^*\|^p.$$

The case $p = 2$ is called *quadratic convergence*.

Of course, as mentioned earlier, these definitions depend on the choice of a norm. It can be shown that quadratic convergence implies superlinear convergence, and superlinear convergence implies linear convergence.

Example 4.3. Consider the sequence of numbers $x_k = 1/2^{r^k}$ for some $r > 1$. Clearly, $x_k \rightarrow x^* = 0$ as $k \rightarrow \infty$. Moreover,

$$x_{k+1} = \frac{1}{2^{r^{k+1}}} = \frac{1}{2^{r^k} \cdot 2^r} = \left(\frac{1}{2^{r^k}}\right)^r = x_k^r,$$

which shows that the sequence has rate of convergence r .

4.3 Convergence of gradient descent

In this section, a norm $\|\cdot\|$ will refer to the 2-norm, unless otherwise stated. We now study the convergence of gradient descent for the least squares problem

$$\text{minimize } f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2, \quad (4.1)$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m \geq n$ of full rank. As we have seen in Lecture 3, the gradient descent method is the procedure

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{r}_k,$$

where the step length and the *residual* are given by

$$\alpha_k = \frac{\|\mathbf{r}_k\|^2}{\|\mathbf{A}\mathbf{r}_k\|^2}, \quad \mathbf{r}_k = \mathbf{A}^\top(\mathbf{b} - \mathbf{A}\mathbf{x}_k) = -\nabla f(\mathbf{x}_k).$$

At the minimizer, the residual is

$$\mathbf{r} = -\nabla f(\mathbf{x}^*) = \mathbf{A}^\top(\mathbf{b} - \mathbf{A}\mathbf{x}^*) = 0, \quad (4.2)$$

and as the sequence \mathbf{x}_k converges to \mathbf{x}^* , the norms of the residuals converge to 0. Conversely, the residual is related to the difference $\mathbf{x}_k - \mathbf{x}^*$ by

$$\mathbf{r}_k = \mathbf{A}^\top(\mathbf{b} - \mathbf{A}\mathbf{x}_k) = \mathbf{A}^\top(\mathbf{b} - \mathbf{A}\mathbf{x}_k - (\mathbf{b} - \mathbf{A}\mathbf{x}^*)) = \mathbf{A}^\top \mathbf{A}(\mathbf{x}_k - \mathbf{x}^*), \quad (4.3)$$

where we used the “intelligent zero” (4.2). Therefore

$$\|\mathbf{x}_k - \mathbf{x}^*\| = \|(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{r}\| \leq \|(\mathbf{A}^\top \mathbf{A})^{-1}\| \|\mathbf{r}_k\|,$$

where $\|\mathbf{B}\| = \max_{\mathbf{x} \neq 0} \|\mathbf{B}\mathbf{x}\|/\|\mathbf{x}\|$ is the operator norm of a matrix \mathbf{B} with respect to the 2-norm. Consequently, if the sequence $\|\mathbf{r}_k\|$ converges to zero, so does the sequence $\|\mathbf{x}_k - \mathbf{x}^*\|$. A reasonable criterium is to stop the algorithm is therefore when the residual norm $\|\mathbf{r}_k\|$ is below a predefined tolerance ε .

The following theorem (whose proof we omit) shows that the gradient descent method for linear least squares converges linearly with respect to the \mathbf{A} norm. The statement involves the *condition number* of \mathbf{A} ⁵. This quantity is defined as

$$\kappa(\mathbf{A}) = \|\mathbf{A}\| \cdot \|\mathbf{A}^\dagger\|.$$

Theorem 4.4. *The error in the $k + 1$ -th iterate is bounded by*

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \left(\frac{\kappa^2(\mathbf{A}) - 1}{\kappa^2(\mathbf{A}) + 1} \right) \|\mathbf{x}_k - \mathbf{x}^*\|.$$

In particular, the gradient descent algorithm converges linearly.

⁵The concept of condition number, introduced by Alan Turing while in Manchester, is one of the most important ideas in numerical analysis, as it is indispensable in studying the performance of numerical algorithms.

Lecture 5

In Lecture 4 we found out that gradient descent works, and has linear convergence. In this lecture we introduce Newton's method, an algorithm that takes advantage of the second derivative and has quadratic convergence under certain circumstances. Throughout this lecture, $\|\cdot\|$ will refer to the 2-norm $\|\cdot\|_2$.

5.1 Newton's Method

Let $f \in C^2(\mathbb{R}^n)$ and let's look again at the unconstrained problem

$$\text{minimize } f(\mathbf{x}).$$

Newton's method starts with a guess \mathbf{x}_0 and then proceeds to compute a sequence of points $\{\mathbf{x}_k\}_{k \geq 0}$ in \mathbb{R}^n by the rule

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \nabla^2 f(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k), \quad k \geq 0. \quad (5.1)$$

The algorithm stops when $\|\mathbf{x}_{k+1} - \mathbf{x}_k\| < \varepsilon$ for some predefined tolerance $\varepsilon > 0$. In the context of the general scheme $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$, the step length is $\alpha_k = 1$, and the search direction is the inverse of the Hessian multiplied with the negative gradient.

Recall that the inner product $\langle \mathbf{p}, \nabla f(\mathbf{x}) \rangle$ is the directional derivative of f , and that a *descent direction* is a direction in which the rate of change (slope) is negative. The following gives a criterion for the search direction in Newton's method to be a descent direction.

Lemma 5.1. *Let $\mathbf{B} \in \mathbb{R}^{n \times n}$ be a positive definite symmetric matrix and $f \in C^1(\mathbf{R}^n)$. Then $\mathbf{p} = -\mathbf{B}^{-1} \nabla f(\mathbf{x})$ is a descent direction of f at \mathbf{x} .*

Proof. If $\mathbf{B} \in \mathbb{R}^{n \times n}$ is symmetric and positive definite, then \mathbf{B}^{-1} is also positive definite, since for all $\mathbf{v} \in \mathbb{R}^n$,

$$\mathbf{v}^\top \mathbf{B}^{-1} \mathbf{v} = (\mathbf{B} \mathbf{B}^{-1} \mathbf{v})^\top \mathbf{B}^{-1} \mathbf{v} = (\mathbf{B}^{-1} \mathbf{v})^\top \mathbf{B}^\top (\mathbf{B}^{-1} \mathbf{v}) = (\mathbf{B}^{-1} \mathbf{v})^\top \mathbf{B} (\mathbf{B}^{-1} \mathbf{v}) > 0.$$

(This can also be seen by noting that the eigenvalues of \mathbf{B}^{-1} are the inverses of the eigenvalues of \mathbf{B} .) For $\mathbf{p} = -\mathbf{B}^{-1} \nabla f(\mathbf{x})$ we then get

$$\langle \mathbf{p}, \nabla f(\mathbf{x}) \rangle = -\langle \mathbf{B}^{-1} \nabla f(\mathbf{x}), \nabla f(\mathbf{x}) \rangle = -\nabla f(\mathbf{x})^\top \mathbf{B}^{-1} \nabla f(\mathbf{x}) < 0,$$

which shows that \mathbf{p} is a descent direction. □

To better understand Newton's method, we first look at the one dimensional case.

Example 5.2. Let $f \in C^2(\mathbb{R})$. In this case Newton's method is described as

$$x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)}, \quad k \geq 0. \quad (5.2)$$

Newton's method looks for a local minimizer in the form of a point x^* such that $f'(x^*) = 0$ and $f''(x^*) > 0$. Setting $g(x) := f'(x)$, we are looking for a *root* x^* ,

$$g(x^*) = 0.$$

One approach to find such a root is to approximate the function $g(x)$ at a point x_k by its tangent line,

$$g(x) \approx g(x_k) + g'(x_k)(x - x_k),$$

and then identify the next iterate x_{k+1} as the root of this linear approximation:

$$g(x_k) + g'(x_k)(x_{k+1} - x_k) = 0 \iff x_{k+1} = x_k - \frac{g(x_k)}{g'(x_k)}$$

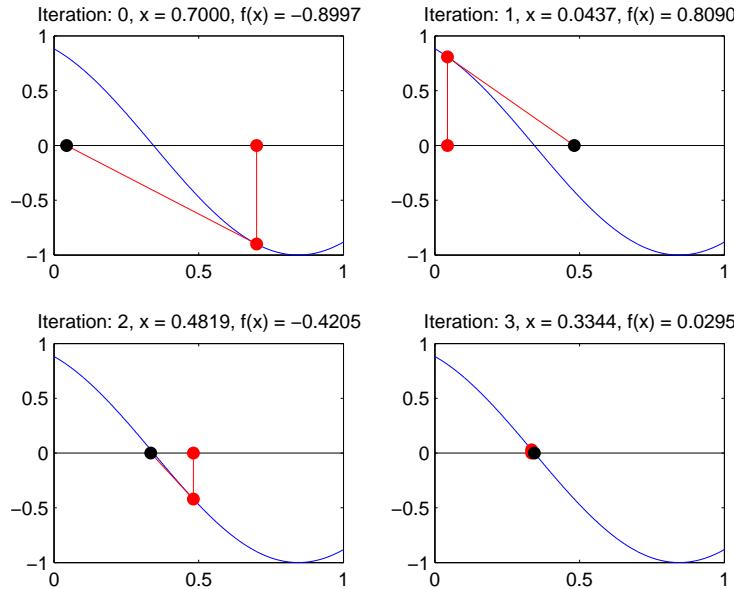


Figure 5.6: Newton's method

Geometrically this corresponds to taking the tangent to g at x_k and setting x_{k+1} to be the intersection of this tangent line with the x -axis, as shown in Figure 5.6. Replacing $g(x) = f'(x)$ gives precisely Newton's method (5.2).

Another way to understand Newton's method is to view it in contrast with gradient descent. While gradient descent corresponds to working with a *linear approximation*

$$f(\mathbf{x}_{k+1}) \approx f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle,$$

Newton's method is based on the *quadratic approximation*,

$$f(\mathbf{x}_{k+1}) \approx f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}_k)(\mathbf{x}_{k+1} - \mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle.$$

Example 5.3. Consider the function f on \mathbb{R}^2 ,

$$f(\mathbf{x}) = \frac{1}{2}(x_1^2 + 10x_2^2).$$

Starting with $\mathbf{x}_0 = (10, 1)^\top$, gradient descent takes 84 iterations to reach accuracy 10^{-6} , while Newton's method, unsurprisingly, takes only one.

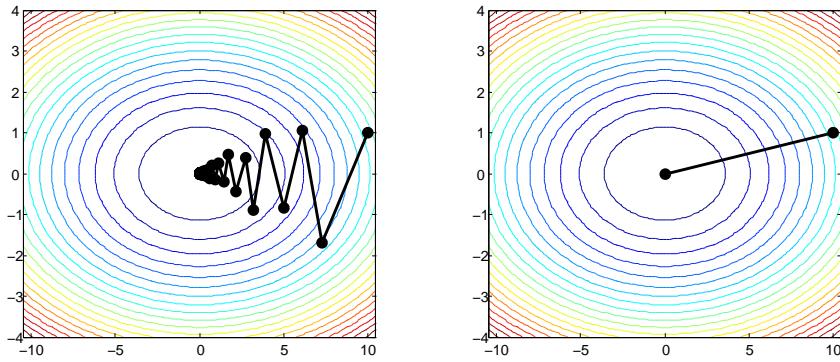


Figure 5.7: Gradient descent vs. Newton's method on a quadratic function.

In practice, when implementing Newton's method one does not explicitly compute the inverse of the Hessian. The reason is that one does not need the inverse itself, but only the product $\nabla^2 f(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k)$, which is the solution of a system of equations $\nabla^2 f(\mathbf{x}_k) \mathbf{y} = \nabla f(\mathbf{x}_k)$ that can be solved much more efficiently. One therefore replaces the update step (5.1) with the following two steps:

$$\begin{aligned} \text{Find } \mathbf{y} \text{ such that } \nabla^2 f(\mathbf{x}_k) \mathbf{y} &= -\nabla f(\mathbf{x}_k), \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \mathbf{y}. \end{aligned}$$

There is a lot that can go wrong with Newton's method. In particular, the matrix $\nabla^2 f(\mathbf{x})$ has to be non-singular, or invertible, at every step. If, however, we start at a point \mathbf{x}_0 that is not too far from a local minimizer \mathbf{x}^* with $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}^*)$ positive definite, then we are on the safe side.

Lemma 5.4. *Let $\mathbf{x}^* \in \mathbb{R}^n$ be such that $\nabla^2 f(\mathbf{x}^*)$ is positive definite. Then there exists an open neighbourhood U of \mathbf{x}^* such that for all $\mathbf{x} \in U$, $\nabla^2 f(\mathbf{x})$ is positive definite.*

For the main result of this lecture, namely the quadratic convergence of Newton's method, we make the additional assumption that the Hessian $\nabla^2 f(\mathbf{x})$ is Lipschitz continuous as a function of \mathbf{x} .

Definition 5.5. A function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is Lipschitz continuous on a domain $\Omega \subseteq \mathbb{R}^n$ with respect to a pair of norms on \mathbb{R}^n and \mathbb{R}^m if there is a constant $L > 0$ such that for all $\mathbf{x}, \mathbf{y} \in \Omega$,

$$\|f(\mathbf{x}) - f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|.$$

The constant L is called the *Lipschitz constant* of the map.

In particular, the Hessian of a function $f \in C^2(\mathbb{R}^n)$, considered as a map from \mathbb{R}^n to $\mathbb{R}^{n \times n}$, is Lipschitz continuous with respect to a norm on \mathbb{R}^n and the corresponding operator norm on $\mathbb{R}^{n \times n}$, if for any \mathbf{x}, \mathbf{y} we have

$$\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|.$$

Theorem 5.6. Let $f \in C^2(\mathbb{R}^n)$ and $\mathbf{x}^* \in \mathbb{R}^n$ a local minimizer with $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}^*) > 0$. Then for \mathbf{x}_0 sufficiently close to \mathbf{x}^* , Newton's method has quadratic convergence.

Proof. (Optional) Assume that $\nabla^2 f(\mathbf{x}_k)$ is positive definite. Consider the difference

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}^*\| &= \|\mathbf{x}_k - \mathbf{x}^* - \nabla^2 f(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k)\| \\ &\stackrel{(1)}{=} \|\mathbf{x}_k - \mathbf{x}^* - \nabla^2 f(\mathbf{x}_k)^{-1} (\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}^*))\| \\ &= \|\nabla^2 f(\mathbf{x}_k)^{-1} (\nabla^2 f(\mathbf{x}_k)(\mathbf{x}_k - \mathbf{x}^*) - (\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}^*))\| \end{aligned} \quad (5.3)$$

where (1) follows from $\nabla f(\mathbf{x}^*) = \mathbf{0}$. The Fundamental Theorem of Calculus tells us

$$\begin{aligned} \nabla f(\mathbf{x}^*) - \nabla f(\mathbf{x}_k) &= \int_0^1 \frac{d}{dt} \nabla f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k)) dt \\ &= \int_0^1 \nabla^2 f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k))(\mathbf{x}^* - \mathbf{x}_k) dt. \end{aligned}$$

Continuing from (5.3), by inserting this identity,

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}^*\| &= \left\| \nabla^2 f(\mathbf{x}_k)^{-1} \int_0^1 [\nabla^2 f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k))] (\mathbf{x}_k - \mathbf{x}^*) dt \right\| \\ &\leq \|\nabla^2 f(\mathbf{x}_k)^{-1}\| \cdot \\ &\quad \int_0^1 \|\nabla^2 f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k))\| dt \cdot \|(\mathbf{x}_k - \mathbf{x}^*)\|. \end{aligned}$$

Applying the Lipschitz bound to the term inside the integral gives

$$\|\nabla^2 f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k))\| \leq Lt\|\mathbf{x}_k - \mathbf{x}^*\|.$$

Integrating this out, we end up with the bound

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \frac{L}{2} \|\nabla^2 f(\mathbf{x}_k)^{-1}\| \cdot \|\mathbf{x}_k - \mathbf{x}^*\|^2.$$

The only remaining issue is that the “constant” on the right-hand side is not a constant. However, the Lipschitz continuity implies that $\nabla^2 f(\mathbf{x}_k)$ converges to $\nabla^2 f(\mathbf{x}^*)$ if \mathbf{x}_k converges to \mathbf{x}^* . Since

the inversion of a matrix is a continuous operation, also $\nabla^2 f(\mathbf{x}_k)^{-1}$ converges to $\nabla^2 f(\mathbf{x}^*)^{-1}$. In particular, if \mathbf{x}_k is sufficiently close to \mathbf{x}^* , we have $\|\nabla^2 f(\mathbf{x}_k)^{-1}\| \leq 2\|\nabla^2 f(\mathbf{x}^*)\|$. Setting $M := L\|\nabla^2 f(\mathbf{x}^*)^{-1}\|/2$, we end up with the bound

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq M \cdot \|\mathbf{x}_k - \mathbf{x}^*\|^2.$$

By Lemma 5.4 there exists an open neighbourhood around \mathbf{x}^* in which $\nabla^2 f(\mathbf{x})$ is positive definite, and within this neighbourhood there is an \mathbf{x}_0 such that $\|\mathbf{x}_0 - \mathbf{x}^*\|^2 < 1/M$, which ensures that all following iterates remain in U . This shows quadratic convergence. \square

Note that the conditions for quadratic convergence in an open neighbourhood U of \mathbf{x}^* are precisely that f is convex on U .

5.2 Quasinewton methods

One drawback of Newton's method is that it requires the computation of the Hessian matrix, which can be expensive. Quasinewton methods use an approximation of the Hessian, $\mathbf{B}_k \approx \nabla^2 f(\mathbf{x}_k)$, at each step of the algorithm. These are constructed in a way that \mathbf{B}_{k+1} can easily be computed from \mathbf{B}_k . A popular method, that is used often in practical applications because of its efficiency, is the **Broyden-Fletcher-Shanno-Goldfarb (BFGS)** method. The BFGS method may be described as follows.

- Start with $\mathbf{x}_0, \mathbf{B}_0$.
- For $k \geq 0$, compute

$$\begin{aligned} \mathbf{p}_k &= \mathbf{B}_k^{-1} \nabla f(\mathbf{x}_k) \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \alpha_k \mathbf{p}_k \text{ for a suitable step length } \alpha_k \\ \mathbf{s}_k &= \alpha_k \mathbf{p}_k \\ \mathbf{y}_k &= \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k) \\ \mathbf{B}_{k+1} &= \mathbf{B}_k + \frac{\mathbf{y}_k \mathbf{y}_k^\top}{\mathbf{y}_k^\top \mathbf{s}_k} - \frac{\mathbf{B}_k \mathbf{s}_k \mathbf{s}_k^\top \mathbf{B}_k}{\mathbf{s}_k^\top \mathbf{B}_k \mathbf{s}_k}. \end{aligned}$$

- Stop if $\|\nabla f(\mathbf{x}_k)\| < \varepsilon$ for some tolerance ε , or if $\|\mathbf{x}_{k+1} - \mathbf{x}_k\| < \varepsilon$.

Lecture 6

In this lecture we will embark on one of the most important modern applications of convex optimization: **machine learning**. The goal of machine learning is to develop methods to automatically *learn* a function

$$h: \mathcal{X} \rightarrow \mathcal{Y},$$

where \mathcal{X} is a space of *inputs* or *features*, and \mathcal{Y} consists of *outputs* or *responses*. The input space \mathcal{X} is usually an \mathbb{R}^n , though the inputs could represent anything such as images, texts, emails, genome sequences, or networks, financial time series, or demographic data. The output could be either **quantitative** values, such as a temperature or the amount of some substance in the body, or **qualitative** or **categorical**, such as YES, NO or $\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$ (for recognising spam or handwritten digits, respectively). The first type of problem is usually called **regression**, while the latter is called **classification**. Machine learning techniques underlie much of modern technology, from car electronics to online product recommendation systems and search engines.

6.1 Supervised Learning

In supervised learning, we have at our disposal a collection of input-output pairs

$$(\mathbf{x}_i, y_i) \in \mathcal{X} \times \mathcal{Y}, \quad 1 \leq i \leq m,$$

and the goal is to *learn* a function $h: \mathcal{X} \rightarrow \mathcal{Y}$ from this data. The given pairs $\{(\mathbf{x}_i, y_i)\}_{1 \leq i \leq m}$ are called the **training set**.

Example 6.1. (Handwriting recognition) Given a dataset of pixel matrices, each representing a grey-scale image, with associated **labels** telling us for each image the number it represents, the task is to use this to train a computer program to recognise new numbers. Such classification tasks are often carried out using **deep neural networks**.

Example 6.2. (Linear regression) Recall from Lecture 1 the problem of finding a function of the form

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p.$$

Given a set of input-output pairs (\mathbf{x}_i, y_i) , arranged in a matrix \mathbf{X} and a vector \mathbf{y} , we saw that we could *guess* the correct $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)^\top$ by solving the *least-squares* optimization problem

$$\underset{\boldsymbol{\beta}}{\text{minimize}} \quad \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2.$$

By now, we know how to solve this problem using gradient descent.

5	0	4	1	9	2	1	3	1	4
3	5	3	6	1	7	2	8	6	9
4	0	9	1	1	2	4	3	2	7
3	8	6	9	0	5	6	0	7	6
1	8	7	9	3	9	8	5	9	3
3	0	7	4	9	8	0	9	4	1
4	4	6	0	4	5	6	1	0	0
1	7	1	6	3	0	2	1	1	7
9	0	2	6	7	8	3	9	0	4
6	7	4	6	8	0	7	8	3	1

Example 6.3. In text classification, the task is to decide to which of a given set of categories a given text belongs. The training data consists of a *bag of words*: this is a large sparse matrix, whose columns represent words and the rows represent articles, with the (i, j) -th entry containing the number of times word j is contained in text i .

	Rooney	Boris
Article 1	5	0
Article 2	1	7

For example, in the above set we would classify the first article as "Sports" and the second one as "Politics". One such training dataset is the Reuters Corpus Volume I (RCV1), an archive of over 800,000 categorised newswire stories.

A typical binary classifier for such a problem would be a **linear classifier** of the form

$$h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} - \tau,$$

with $\mathbf{w} \in \mathbb{R}^n$ and $\tau \in \mathbb{R}$. Given a text, represented as a row of the dataset \mathbf{x} , it is classified into one of two classes $\{+1, -1\}$, depending on whether $h(\mathbf{x}) > 0$ or $h(\mathbf{x}) < 0$.

Suppose we have the training data $\{\mathbf{x}_i, y_i\}_{1 \leq i \leq m}$ and we found a candidate function $h: \mathcal{X} \rightarrow \mathcal{Y}$. How do we assess if this is a good fit? One usually assigns to the problem a **loss function** $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ that measures the mismatch of a prediction. One would then aim to find a function h among a set of candidates that minimizes the loss when applying the function to the training data:

$$\underset{h}{\text{minimize}} \quad \frac{1}{m} \sum_{i=1}^m \ell(h(\mathbf{x}_i, y_i)).$$

The form of the loss function depends on the problem at hand, but two typical candidates are the **square error** for regression problems,

$$\ell(h(\mathbf{x}), y) = (h(\mathbf{x}) - y)^2,$$

and the indicator loss function

$$\mathbf{1}\{h(\mathbf{x}) \neq y\},$$

where $\mathbf{1}\{A\}$ is the indication function, which takes the value 1 if A is true, and 0 else. As this function is not continuous, in practice one often encounters the *log loss* function,

$$\ell(h(\mathbf{x}), y) = \log(1 + e^{-h(\mathbf{x})y}).$$

Note that if $h(\mathbf{x})$ and y have the same sign (corresponding to a match), then the value of this function will be close to zero.

Suppose we have a binary classification task at hand, with $\mathcal{Y} = \{-1, 1\}$. We could *learn* the following function from our data:

$$h(\mathbf{x}) = \begin{cases} y_i & \text{if } \mathbf{x} = \mathbf{x}_i, \\ 1 & \text{otherwise.} \end{cases}$$

The **empirical misclassification risk** in for this problem is then 0,

$$R(h) = \frac{1}{m} \sum_{i=1}^m \mathbf{1}\{h(\mathbf{x}_i) \neq y_i\} = 0.$$

Nevertheless, this is not a good classifier: it will not perform very well outside of the training set. This is an example of **overfitting**: when the function is adapted too closely to the seen data. To remedy this, one often (randomly) splits the available data into a **training set** and a **test set**. The training set is used to find the function h , while the test set is for testing how good it is (in practise, one often adds a validation set, used to tune some parameters of the problem).

Example 6.4. In the linear regression problem with the quadratic loss function, we end up with the minimization problem

$$\underset{\boldsymbol{\beta}}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \boldsymbol{\beta} - y_i)^2 = \frac{1}{n} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2.$$

This is precisely the problem encountered in Lecture 1.

Example 6.5. In the example of classifying a text into two classes using a linear classifier, we can look at the problem

$$\underset{\mathbf{w}, \tau}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^m \mathbf{1}\{\mathbf{w}^\top \mathbf{x}_i - \tau \neq y_i\}.$$

Since this function is not smooth or even continuous, we can work with the log-loss function described above,

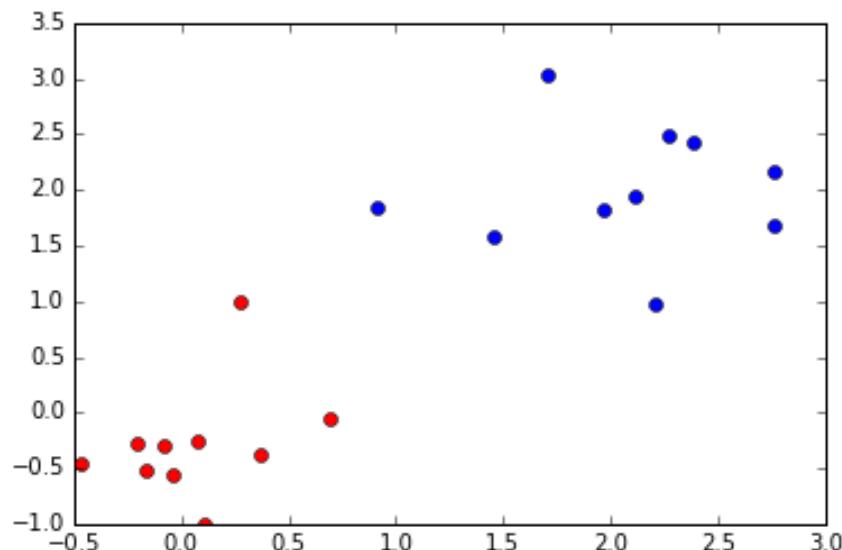
$$\underset{\mathbf{w}, \tau}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^m \ell(\mathbf{w}^\top \mathbf{x}_i - \tau, y_i).$$

This function can, in principle, be minimized using gradient descent or Newton's method.

In the following example, we use this classifier to separate an artificial data set of 20 points. 10 points are generated as normal (Gaussian) distributed vectors in \mathbb{R}^2 , with mean $(2, 2)^\top$ and variance 0.25, while the other 10 points are generated with mean $(0, 0)^\top$ and the same variance.

```
In [17]: X = np.zeros((2,20))
for i in range(10):
    X[:,i] = np.array([2,2]) + 0.5*rnd.randn(2)
for i in range(10):
    X[:,10+i] = np.array([0,0]) + 0.5*rnd.randn(2)

% matplotlib inline
plt.plot(X[0,0:10], X[1,0:10], 'o')
plt.plot(X[0,10:], X[1,10:], 'o', color='red')
plt.show()
```



We next minimize the *log-loss function*,

$$F(\mathbf{w}) = \sum_{i=1}^{20} \log(1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i - \tau)}),$$

where the $\mathbf{x}_i \in \mathbb{R}^2$ are the points, and the $y_i \in \{-1, 1\}$ classify these as either red or blue. As minimization algorithm we choose gradient descent with backtracking, though any other reasonable algorithm would do.

```
In [18]: def graddesc_bt(f, df, x0, tol, maxiter=100, rho=0.5, c=0.1):
    """
    Gradient descent with backtracking
    """
    x = np.vstack((x0+2*tol*np.ones(x0.shape),x0)).transpose()
    i = 1
    while (la.norm(x[:,i]-x[:,i-1]) > tol) and (i < maxiter):
        p = -df(x[:,i])
        # Start backtracking
        alpha = 1
        xnew = x[:,i] + alpha*p
        while (f(xnew) >= f(x[:,i]) + alpha*c*np.dot(p, df(x[:,i]))):
            alpha = alpha*rho
            xnew = x[:,i] + alpha*p
        x = np.concatenate((x,xnew.reshape((len(x0),1))), axis=1)
        i += 1
    return x[:,1:]
```

```
In [19]: y = np.concatenate((-np.ones(10), np.ones(10)))
X = np.concatenate((X, np.ones((1,20))), axis=0)

def f(w):
    return np.sum(np.log(1+np.exp(-y*(np.dot(w,X)))))

def df(w):
    return np.array([-np.sum(y*X[0,:]*np.exp(-y*np.dot(w,X))/(
        1+np.exp(-y*np.dot(w,X)))), 
        -np.sum(y*X[1,:]*np.exp(-y*np.dot(w,X))/(
        1+np.exp(-y*np.dot(w,X)))), 
        -np.sum(y*X[2,:]*np.exp(-y*np.dot(w,X))/(
        1+np.exp(-y*np.dot(w,X))))])

W = graddesc_bt(f, df, np.array([1.,1.,1.]), 1.e-8)
```

The algorithm gives as output a vector $w = (w_1, w_2, w_3) \in \mathbb{R}^3$, so that the line separating the two sets of points is given by

$$w_1x_1 + w_2x_2 + w_3 = 0.$$

Rearranging this as a function $\ell(x_1) = x_2 = -(w_1x_1 + w_3)/w_2$ and plotting the line into the pointcloud, we get the following graph.

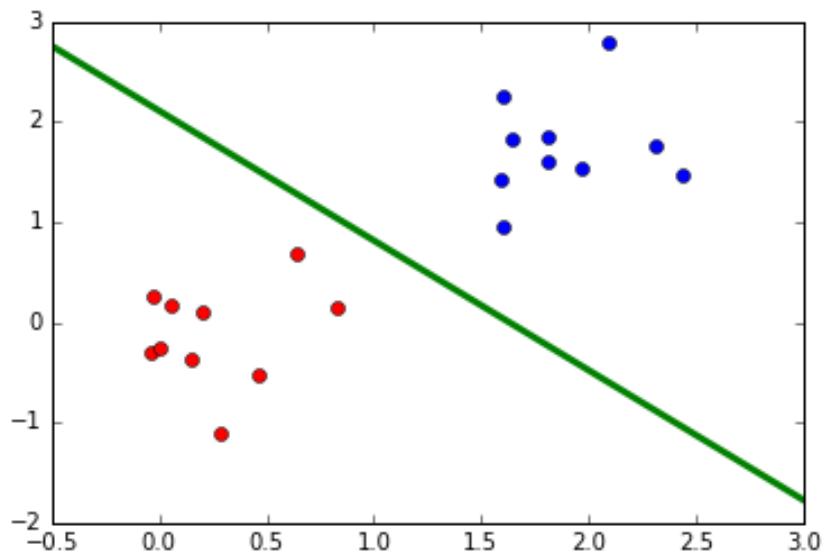
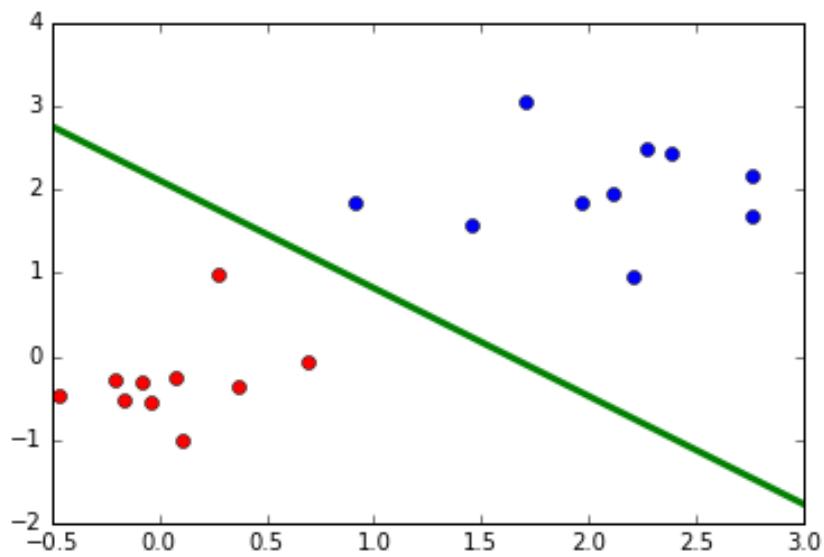
```
In [20]: def l(x):
    return (-W[0,-1]*x-W[2,-1])/W[1,-1]

xx = np.linspace(-0.5,3,100)
yy = l(xx)
plt.plot(X[0,0:10], X[1,0:10], 'o')
plt.plot(X[0,10:], X[1,10:], 'o', color='red')
plt.plot(xx,yy,linewidth=3)
plt.show()
```

Of course, we would like to know how our classifier does with a *test set* generated according to the same distribution. The result is seen the following plot.

The given example is unproblematic, since the two sets of points are generated in a way that makes it unlikely that they can't be separated using a linear classifier. With real data, the situation may not be as favourable.

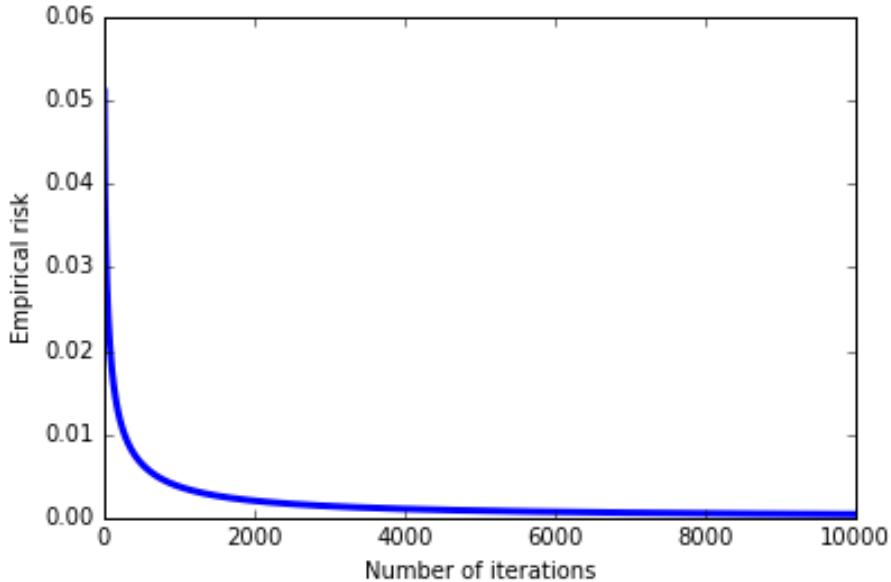
It is also interesting to see how fast the objective function (the *empirical risk*) approaches the minimum. While the value of the objective function decreases very quickly at the beginning, the



decrease slows down dramatically and the algorithm does not reach the desired accuracy within the limit we set for the number of iterations. Now, however, that for the classification task we are interested, we don't need the best possible precision! In fact, the noise in the data makes too much precision meaningless, and we can safely use the result obtained after 100 iterations. Other algorithms and approaches may perform better for this problem.

```
In [21]: risk = np.zeros(W.shape[1])
for i in range(W.shape[1]):
    risk[i] = f(W[:,i])

plt.plot(range(W.shape[1]), risk, linewidth=3)
plt.xlabel('Number of iterations')
plt.ylabel('Empirical risk')
plt.show()
```



We have seen how a machine learning problem can be turned into an optimization problem. Unfortunately, when the amount of data is large, computing gradients and Hessians of the resulting functions can be very expensive. Suppose the function h depends on parameters \mathbf{w} which we want to optimize over, $h = h(\mathbf{x}; \mathbf{w})$, and denote the composition of the loss function ℓ with h on the data \mathbf{x}_i by

$$f_i(\mathbf{w}) := \ell(h(\mathbf{x}_i; \mathbf{w}), y_i).$$

We can then write our optimization problem as

$$\underset{\mathbf{w}}{\text{minimize}} \quad F(\mathbf{w}), \quad F(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{w}).$$

The gradient is then

$$\nabla F(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^n \nabla f_i(\mathbf{w}).$$

For large datasets (for example, $m > 10^6$ is not uncommon), computing all the gradients in each step is highly ineffective.

An old algorithm that has recently been revived in light of big data is **stochastic gradient descent**. The idea is to compute, at each step, only the gradient of *one* (or sometimes a few) of the summands of the loss function:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f_{i_k}(\mathbf{w}_k),$$

where the index i_k is *chosen at random* from $\{1, \dots, m\}$ at each step.

Each step of this algorithm is clearly faster than gradient descent, but does it even work? It turns out that this algorithm converges quickly and that it is surprisingly effective on large datasets.

Part III

Linear Programming

Lecture 7

“[A]lmost every “convex” idea can be explained by a two-dimensional picture.”

— Alexander Barvinok

In this lecture we begin our study of the theory underlying constrained convex optimization. One way to define a *convex optimization problem* is

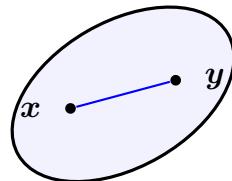
$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && f_1(\mathbf{x}) \leq 0 \\ & && \dots \\ & && f_m(\mathbf{x}) \leq 0 \\ & && \mathbf{x} \in \Omega \end{aligned}$$

where $f, f_1, \dots, f_m: \mathbb{R}^n \rightarrow \mathbb{R}$ are *convex* functions and $\Omega \subseteq \mathbb{R}^n$ is a *convex set*. The special case where the f and the f_i are linear functions and $\Omega = \mathbb{R}^n$ is known as linear programming, and is studied first. Before embarking on the study of models and algorithms for convex optimization, we need to study convex sets in more depth.

7.1 Convex sets

We recall the definition of a convex set.

Definition 7.1. A set $C \subseteq \mathbb{R}^n$ is a *convex set*, if for all $\mathbf{x}, \mathbf{y} \in C$ and $\lambda \in [0, 1]$, $\lambda\mathbf{x} + (1 - \lambda)\mathbf{y} \in C$. In words, for every two points in C , the line joining them is also in C . A compact (closed and bounded) convex set is called a *convex body*.



We will denote by $\mathcal{C}(\mathbb{R}^n)$ the collection of convex sets and by $\mathcal{K}(\mathbb{R}^n)$ the collection of convex bodies. The following Lemma is left as an exercise.

Lemma 7.2. Let $C, D \in \mathcal{C}(\mathbb{R}^n)$ be convex sets. Then the following are also convex.

- $C \cap D$;
- $C + D = \{\mathbf{x} + \mathbf{y} : \mathbf{x} \in C, \mathbf{y} \in D\}$;
- $AC = \{A\mathbf{x} : \mathbf{x} \in C\}$, where $A \in \mathbb{R}^{m \times n}$.

The *convex hull* $\text{conv } S$ of a set S is the intersection of all convex sets containing S . Clearly, if S is convex, then $S = \text{conv } S$.

Example 7.3. Let $S = \{(1, 1)^\top, (1, -1)^\top, (-1, 1)^\top, (-1, -1)^\top, (0, 0)^\top\}$. The convex hull of this set is the square.

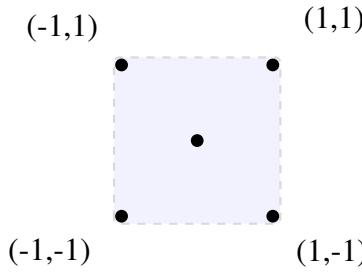


Figure 7.8: A convex hull of five points.

A *convex combination* of points $\mathbf{x}_1, \dots, \mathbf{x}_k$ is a linear combination

$$\sum_{i=1}^k \lambda_i \mathbf{x}_i$$

such that $\lambda_i \geq 0$ and $\sum_{i=1}^k \lambda_i = 1$. It can be shown inductively that convex sets are closed under convex combinations: any convex combination of points in $C \in \mathcal{C}(\mathbb{R}^n)$ is still in C . In fact, the set of all convex combinations of points in a set S is the convex hull of S .

Lemma 7.4. Let S be a set. Then

$$\text{conv } S = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} = \sum_{i=1}^k \lambda_i \mathbf{x}_i, \mathbf{x}_i \in S, \sum_{i=1}^k \lambda_i = 1, \lambda_i \geq 0\}.$$

Example 7.5. A hyperplane, defined as the solution set of one linear equation,

$$H = \{\mathbf{x} : \langle \mathbf{a}, \mathbf{x} \rangle = b\},$$

is a convex set. Define the halfspaces H_+ and H_- as the two sides that H divides \mathbb{R}^n into:

$$H_- = \{\mathbf{x} : \langle \mathbf{a}, \mathbf{x} \rangle \leq b\}, \quad H_+ = \{\mathbf{x} : \langle \mathbf{a}, \mathbf{x} \rangle \geq b\}$$

These are also convex sets.

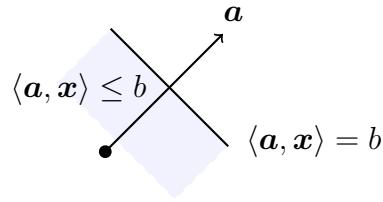


Figure 7.9: Hyperplane and halfspace

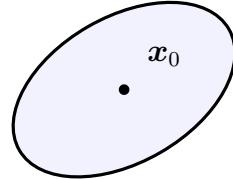


Figure 7.10: An ellipse

Example 7.6. Euclidean balls and ellipsoids are common examples of convex sets. Let \mathbf{P} be a positive semidefinite symmetric matrix. Then an ellipsoid with center \mathbf{x}_0 is a set of the form

$$\mathcal{E} = \{\mathbf{x} : \langle \mathbf{x} - \mathbf{x}_0, \mathbf{P}^{-1}(\mathbf{x} - \mathbf{x}_0) \rangle \leq 1\}.$$

A Euclidean unit ball is the special case $\mathbf{P} = \mathbf{I}$.

Example 7.7. A *convex cone* is a set C such that for all \mathbf{x}, \mathbf{y} and $\lambda \geq 0, \mu \geq 0$, $\lambda\mathbf{x} + \mu\mathbf{y} \in C$. It is easily verified that such a set is convex. Three important cones are the following:

1. The non-negative orthant $\mathbb{R}_+^n = \{\mathbf{x} \in \mathbb{R}^n : x_i \geq 0, 1 \leq i \leq n\}$,
2. The second order (ice cream) cone (or Lorentz cone)

$$C_\alpha = \{\mathbf{x} : \sum_{i=1}^{n-1} x_i^2 \leq x_n^2\},$$

3. The cone \mathcal{S}_+^n of positive semidefinite symmetric matrices.

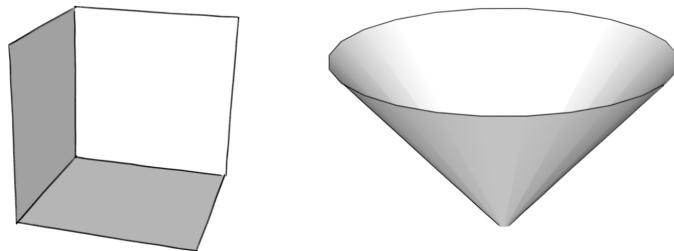


Figure 7.11: The orthant and the second order cone

Possibly the most important result in convex geometry is the *hyperplane separation theorem*. We first need the following.

Lemma 7.8. *Let C be a non-empty convex set and $\mathbf{x} \notin C$. Then there exists a point $\mathbf{y} \in C$ that minimizes the distance $\|\mathbf{x} - \mathbf{y}\|$. Moreover, for all $\mathbf{z} \in C$ we have*

$$\langle \mathbf{z} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle \leq 0.$$

In words, the vectors $\mathbf{z} - \mathbf{y}$ and $\mathbf{x} - \mathbf{y}$ form an obtuse angle.

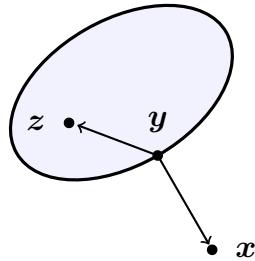


Figure 7.12: Internal and external directions

Proof. Since $C \neq \emptyset$, there exists $r > 0$ such that the ball $B(\mathbf{x}, r) := \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\| \leq r\}$ intersected with C is not empty. Since $K := C \cap B(\mathbf{x}, r)$ is compact (closed and bounded) and the function $\|\mathbf{y} - \mathbf{x}\|$ is continuous on K , it has a minimizer $\mathbf{y} \in K$. For the second claim, note that since C is convex, for every $\lambda \in [0, 1]$,

$$\mathbf{w} = \lambda \mathbf{z} + (1 - \lambda) \mathbf{y} \in C.$$

For the distance between \mathbf{z} and \mathbf{x} we then get

$$\begin{aligned} \|\mathbf{w} - \mathbf{x}\|^2 &= \|\lambda \mathbf{z} + (1 - \lambda) \mathbf{y} - \mathbf{x}\|^2 = \|\lambda(\mathbf{z} - \mathbf{y}) - (\mathbf{x} - \mathbf{y})\|^2 \\ &= \lambda^2 \|\mathbf{z} - \mathbf{y}\|^2 - 2\lambda \langle \mathbf{z} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle + \|\mathbf{x} - \mathbf{y}\|^2. \end{aligned}$$

We now prove the claim by contradiction. Assume $\langle \mathbf{z} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle > 0$. Then we can choose λ such that

$$0 < \lambda < \min \left\{ \frac{2\langle \mathbf{x} - \mathbf{y}, \mathbf{z} - \mathbf{y} \rangle}{\|\mathbf{z} - \mathbf{y}\|^2}, 1 \right\}.$$

With such a λ we get

$$\|\mathbf{w} - \mathbf{x}\|^2 = \lambda^2 \|\mathbf{z} - \mathbf{y}\|^2 - 2\lambda \langle \mathbf{z} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle + \|\mathbf{x} - \mathbf{y}\|^2 < \|\mathbf{x} - \mathbf{y}\|^2.$$

This inequality, however, contradicts the assumption that \mathbf{y} is a closest point, so that $\langle \mathbf{z} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle \leq 0$ has to hold. \square

In what follows write $\text{int}S$ for the *interior* of a set S .

Theorem 7.9. *Let C be a closed convex set and $\mathbf{x} \notin C$. Then there exists a hyperplane H such that $C \subset \text{int}H_-$ and $\mathbf{x} \in \text{int}H_+$.*

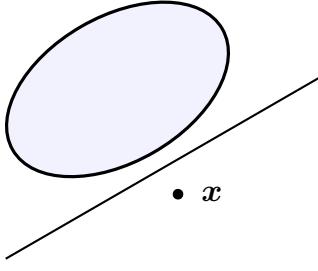


Figure 7.13: A separating hyperplane

Proof. Let $\mathbf{y} \in C$ be a nearest point to \mathbf{x} in C , i.e., a point such that for all other $\mathbf{z} \in C$, $\|\mathbf{x} - \mathbf{y}\| \leq \|\mathbf{x} - \mathbf{z}\|$. Define

$$\mathbf{a} = \mathbf{x} - \mathbf{y}, \quad b = (\|\mathbf{x}\|^2 - \|\mathbf{y}\|^2)/2.$$

We aim to show that $\langle \mathbf{a}, \mathbf{x} \rangle = b$ defines a separating hyperplane.

For this we have to show that

1. $\langle \mathbf{a}, \mathbf{x} \rangle > b$;
2. For all $\mathbf{z} \in C$, $\langle \mathbf{a}, \mathbf{z} \rangle < b$.

For (1), note that

$$\langle \mathbf{a}, \mathbf{x} \rangle = \langle \mathbf{x} - \mathbf{y}, \mathbf{x} \rangle > \langle \mathbf{x} - \mathbf{y}, \mathbf{x} \rangle - \frac{1}{2}\|\mathbf{x} - \mathbf{y}\|^2 = \frac{1}{2}(\|\mathbf{x}\|^2 - \|\mathbf{y}\|^2) = b.$$

To prove (2), assume on the contrary that there exists a $\mathbf{z} \in C$ such that $\langle \mathbf{a}, \mathbf{z} \rangle \geq b$. We know that the point $\mathbf{y} \in C$ satisfies the inequality (2), since

$$\langle \mathbf{a}, \mathbf{y} \rangle < \langle \mathbf{a}, \mathbf{y} \rangle + \frac{1}{2}\|\mathbf{a}\|^2 = \langle \mathbf{a}, \mathbf{y} \rangle + \frac{1}{2}\|\mathbf{x} - \mathbf{y}\|^2 = b.$$

Therefore,

$$\langle \mathbf{a}, \mathbf{z} - \mathbf{y} \rangle = \langle \mathbf{a}, \mathbf{z} \rangle - \langle \mathbf{a}, \mathbf{y} \rangle > b - b = 0,$$

but this contradicts Lemma 7.8. We therefore conclude $\langle \mathbf{a}, \mathbf{z} \rangle < b$. The separating hyperplane H is thus defined by the equation $\langle \mathbf{a}, \mathbf{x} \rangle = b$. \square

Lecture 8

Linear programming is about problems of the form

$$\begin{aligned} & \text{maximize} && \langle \mathbf{c}, \mathbf{x} \rangle \\ & \text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{b}, \end{aligned}$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{c} \in \mathbb{R}^n$, and $\mathbf{b} \in \mathbb{R}^m$, and the inequality sign means inequality in each row. The *feasible set* is the set of all possible candidates,

$$\mathcal{F} = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} \leq \mathbf{b}\}.$$

This set can be empty (example: $x \leq 1$, $-x \leq -2$), unbounded (example: $x \leq 1$) or bounded (example: $x \leq 1$, $-x \leq 0$). In any case, it is a convex set. To understand linear programming it is of paramount importance to understand the geometry of the feasible sets of linear programming, also called *polyhedra*.

8.1 Linear Programming Duality: a first glance

Suppose we are faced with a linear programming problem and would like to know if the feasible set \mathcal{F} is empty or not, i.e., if $\mathbf{A}\mathbf{x} \leq \mathbf{b}$ has a solution. If it is not empty, we can certify that by producing a vector from \mathcal{F} . To verify that the feasible set is empty is more tricky: we are asked to show that *no* vector lives in \mathcal{F} . What we can try to do, however, is to show that the assumption of a solution to $\mathbf{A}\mathbf{x} \leq \mathbf{b}$ would lead to a contradiction. Denote by \mathbf{a}_i^\top the rows of \mathbf{A} . Assuming $\mathbf{x} \in \mathcal{F}$, then given a vector $\mathbf{0} \neq \boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_m)^\top$ with $\lambda_i \geq 0$, the linear combination satisfies

$$\sum_{i=1}^m \lambda_i \mathbf{a}_i^\top \mathbf{x} \leq \sum_{i=1}^m \lambda_i b_i = \langle \boldsymbol{\lambda}, \mathbf{b} \rangle. \quad (8.1)$$

If we can find parameters $\boldsymbol{\lambda}$ such that the left-hand side of (8.1) is identically 0 and the right-hand side is strictly negative, then we have found a contradiction and can conclude that no \mathbf{x} satisfies $\mathbf{A}\mathbf{x} \leq \mathbf{b}$. A condition that ensures this is

$$\sum_{i=1}^m \lambda_i \mathbf{a}_i = \mathbf{0}, \quad \langle \boldsymbol{\lambda}, \mathbf{b} \rangle < 0. \quad (8.2)$$

In matrix form,

$$\exists \boldsymbol{\lambda} \geq \mathbf{0}, \mathbf{A}^\top \boldsymbol{\lambda} = \mathbf{0}, \langle \boldsymbol{\lambda}, \mathbf{b} \rangle < 0.$$

This condition will still be satisfied if we normalise the vector λ such that $\sum_{i=1}^m \lambda_i = 1$, so the statement says that $\mathbf{0}$ is a convex combination of the vectors defining the equations.

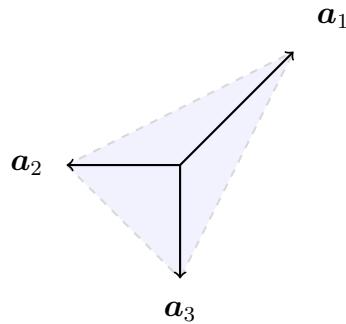
Example 8.1. Consider the system

$$\begin{aligned} x_1 + x_2 &\leq 2 \\ -x_1 &\leq -1 \\ -x_2 &\leq -1.5. \end{aligned} \tag{8.3}$$

The transpose matrix \mathbf{A}^\top and the vector \mathbf{b} are

$$\mathbf{A}^\top = \begin{pmatrix} 1 & -1 & 0 \\ 1 & 0 & -1 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 2 \\ -1 \\ -1.5 \end{pmatrix}$$

Drawing the columns of \mathbf{A}^\top we get We can get the origin as a convex combination of the



vectors \mathbf{a}_i (drawing a rope around them encloses the origin), and such a combination is given by $\lambda = (1/3, 1/3, 1/3)$. Taking the scalar product with the vector \mathbf{b} we get

$$\langle \lambda, \mathbf{b} \rangle = \frac{1}{3}(2 - 1 - 1.5) = -\frac{1}{6} < 0.$$

This shows that the system (8.3) does not have a solution (a fact that in this simple example can also be seen by drawing a picture).

It turns out that Condition (D) is not only sufficient but also necessary, and the separating hyperplane theorem is an essential part of this. We first make a detour in order to better understand the feasible sets, the polyhedra.

8.2 Polyhedra

Definition 8.2. A polyhedron (plural: polyhedra) is a set defined as the solution of linear equalities and inequalities,

$$P = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{Ax} \leq \mathbf{b}\}, \tag{8.1}$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$.

More classically, we can write out the equations.

$$\begin{aligned} a_{11}x_1 + \cdots + a_{1n}x_d &\leq b_1, \\ &\dots \\ a_{m1}x_1 + \cdots + a_{mn}x_d &\leq b_m. \end{aligned} \tag{8.2}$$

We now introduce some useful terminology and concepts associated to polyhedra, and illustrate them with a few examples. A supporting hyperplane H of a polyhedron P is a hyperplane such that $P \subseteq H_-$, where H_- is a halfspace associated to H . If H is a supporting hyperplane, then a set of the form $F = H \cap P$ is called a *face* of P . In particular, the polyhedron P is a face. Each of the inequalities $\langle \mathbf{a}_i, \mathbf{x} \rangle \leq b_i$ in (8.1) defines a supporting hyperplane, and therefore a face. The *dimension* of a face F , $\dim F$, is the smallest dimension of an affine space containing F . Faces of dimension $\dim F = \dim P - 1$ are called *facets*, faces of dimension 0 are vertices, and of dimension 1 edges. A vertex can equivalently be characterised as a point $\mathbf{x} \in P$ that can not be written as a convex combination of two other points in P .

Example 8.3. Polyhedra in one dimension are the sets $[a, b]$, $[a, \infty)$, $(-\infty, b]$, \mathbb{R} or \emptyset , where $a \leq b$. Each of them is clearly convex.

Example 8.4. The set

$$P = \{\mathbf{x} \in \mathbb{R}^2 : x_1 + x_2 \leq 1, x_1 \geq 0, x_2 \geq 0\}.$$

is the polyhedron shown in Figure 8.14. We can write the defining inequalities in standard form $\mathbf{Ax} \leq \mathbf{b}$ by setting

$$\mathbf{A} = \begin{pmatrix} 1 & 1 \\ -1 & 0 \\ 0 & -1 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}.$$

This polyhedron has one face of dimension 2 (itself), three facets of dimension 1 (the sides, cor-

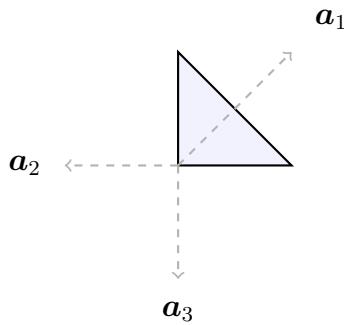


Figure 8.14: A two-dimensional polyhedron and defining equations.

responding to the three equations), and three vertices of dimension 0 (the corners, corresponding any two of the defining equations).

Example 8.5. The polyhedron in Albrecht Dürer’s famous “Melencolia I”, aka the “truncated triangular trapezohedron”, can be described using eight inequalities, see Figure 8.15.



$$\begin{aligned}
0.7071x_1 - 0.4082x_2 + 0.3773x_3 &\leq 1 \\
-0.7071x_1 + 0.4082x_2 - 0.3773x_3 &\leq 1 \\
0.7071x_1 + 0.4082x_2 - 0.3773x_3 &\leq 1 \\
-0.7071x_1 - 0.4082x_2 + 0.3773x_3 &\leq 1 \\
0.8165x_2 + 0.3773x_3 &\leq 1 \\
-0.8165x_2 - 0.3773x_3 &\leq 1 \\
0.6313x_3 &\leq 1 \\
-0.6313x_3 &\leq 1
\end{aligned}$$

Figure 8.15: Dürer's Melencolia and equations defining the polyhedron

We now move to a different characterization of bounded polyhedra. The main result of this lecture is that bounded polytopes can be described completely from knowing their vertices. A polyhedron P is called bounded, if there exists a ball $B(\mathbf{0}, r)$ with $r > 0$ such that $P \subset B(\mathbf{0}, r)$. For example, halfspaces are not bounded, but the polytope from Examples 8.4 and 8.5 are.

We first observe the nontrivial fact that a polyhedron has only finitely many vertices.

Definition 8.6. A *polytope* is the convex hull of finitely many points,

$$P = \text{conv}(\{x_1, \dots, x_k\}) = \left\{ \sum_{i=1}^k \lambda_i x_i, \lambda_i \geq 0, \sum_i \lambda_i = 1 \right\}.$$

Theorem 8.7. A bounded polyhedron P is the convex hull of its vertices.

Example 8.8. The triangle in Example 8.4 is the convex hull of the points $(0, 0)^\top$, $(0, 1)^\top$, and $(1, 0)^\top$.

Example 8.9. The Dürer polytope is the convex hull of the following 12 vertices:

$$\begin{aligned}
\mathbf{v}_1 &= \begin{pmatrix} -1.4142 \\ -0.8165 \\ -0.8835 \end{pmatrix}, \mathbf{v}_2 = \begin{pmatrix} -1.4142 \\ 0.8165 \\ 0.8835 \end{pmatrix}, \mathbf{v}_3 = \begin{pmatrix} -0.8536 \\ -0.4928 \\ -1.5840 \end{pmatrix}, \mathbf{v}_4 = \begin{pmatrix} -0.8536 \\ 0.4928 \\ 1.5840 \end{pmatrix}, \\
\mathbf{v}_5 &= \begin{pmatrix} -0.0000 \\ -1.6330 \\ 0.8835 \end{pmatrix}, \mathbf{v}_6 = \begin{pmatrix} 0.0000 \\ -0.9856 \\ 1.5840 \end{pmatrix}, \mathbf{v}_7 = \begin{pmatrix} -0.0000 \\ 0.9856 \\ -1.5840 \end{pmatrix}, \mathbf{v}_8 = \begin{pmatrix} 0.0000 \\ 1.6330 \\ -0.8835 \end{pmatrix}, \\
\mathbf{v}_9 &= \begin{pmatrix} 0.8536 \\ -0.4928 \\ -1.5840 \end{pmatrix}, \mathbf{v}_{10} = \begin{pmatrix} 0.8536 \\ 0.4928 \\ 1.5840 \end{pmatrix}, \mathbf{v}_{11} = \begin{pmatrix} 1.4142 \\ -0.8165 \\ -0.8835 \end{pmatrix}, \mathbf{v}_{12} = \begin{pmatrix} 1.4142 \\ 0.8165 \\ 0.8835 \end{pmatrix}.
\end{aligned}$$

The converse of Theorem 8.7 is also true.

Theorem 8.10. *A polytope is a bounded polyhedron.*

The equivalence between polytopes and bounded polyhedra gives a first glimpse into linear programming duality theory, a topic of central importance in both modeling and algorithm design.

Lecture 9

Linear programming duality associates to a linear programming problem a **dual problem**, with the property that the optimal values of the original and of the dual problem coincide. Duality is an important tool in applications and in the design of algorithms. Linear programming duality rests upon an important family of results in convex geometry, known collectively as Farkas' Lemma.

9.1 Farkas' Lemma

Recall the definition of a convex cone. This is a set C such that for all $\mathbf{x}, \mathbf{y} \in C$ and $\lambda_1, \lambda_2 \geq 0$ we have $\lambda_1\mathbf{x} + \lambda_2\mathbf{y} \in C$.

Lemma 9.1. (*Hyperplane separation for cones*) Let $C \neq \mathbb{R}^n$ be a closed convex cone and $\mathbf{z} \notin C$. Then there exists a linear hyperplane such that $C \subseteq H_-$ and $\mathbf{z} \in \text{int}H_+$.

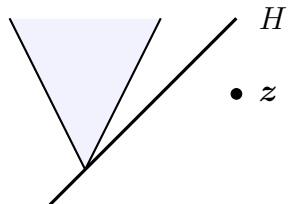


Figure 9.16: Separating hyperplane for a cone

Proof. From Lecture 7 we know that there exists an affine hyperplane H^a separating C and \mathbf{z} . Let this affine hyperplane be given by $\langle \mathbf{a}, \mathbf{x} \rangle = b$. We would like to show that $H = \{\mathbf{x} : \langle \mathbf{a}, \mathbf{x} \rangle = 0\}$ is a linear hyperplane separating C and \mathbf{z} . From $\langle \mathbf{a}, \mathbf{z} \rangle > b$ we clearly get $\langle \mathbf{a}, \mathbf{z} \rangle > 0$. Also, $\mathbf{0} \in H_-$, since $\langle \mathbf{a}, \mathbf{0} \rangle = 0$. Assume now that there exists a point $\mathbf{x} \in C$ such that $\langle \mathbf{a}, \mathbf{x} \rangle = c > 0$. Since C is a cone, for all $\lambda > 0$ we have that $\lambda\mathbf{x} \in C$. Choosing λ so that $\lambda > b/c$ we get

$$\langle \mathbf{a}, \lambda\mathbf{x} \rangle = \lambda c > b,$$

in contradiction to $C \subset H_-$. This shows that H is a linear separating hyperplane. \square

Theorem 9.2. (Farkas' Lemma) Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$, there exists a vector \mathbf{x} such that

$$\mathbf{A}\mathbf{x} = \mathbf{b}, \quad \mathbf{x} \geq 0$$

if and only if there is no $\mathbf{y} \in \mathbb{R}^m$ such that

$$\mathbf{A}^\top \mathbf{y} \geq 0, \quad \langle \mathbf{y}, \mathbf{b} \rangle < 0.$$

Proof. Assume $\mathbf{A}\mathbf{x} = \mathbf{b}$ has a solution $\mathbf{x} \geq 0$. Then for any $\mathbf{y} \neq 0$ such that $\mathbf{A}^\top \mathbf{y} \geq 0$,

$$0 \leq \langle \mathbf{A}^\top \mathbf{y}, \mathbf{x} \rangle = \langle \mathbf{y}, \mathbf{A}\mathbf{x} \rangle = \langle \mathbf{y}, \mathbf{b} \rangle,$$

which shows that $\mathbf{A}^\top \mathbf{y} \geq 0$ and $\langle \mathbf{y}, \mathbf{b} \rangle < 0$ are not simultaneously possible.

Assume now that $\mathbf{A}\mathbf{x} = \mathbf{b}$ has no solution that satisfies $\mathbf{x} \geq 0$. Let $\mathbf{a}_1, \dots, \mathbf{a}_n$ be the columns of \mathbf{A} . The set of $\mathbf{A}\mathbf{x}$ for $\mathbf{x} \geq 0$ is the set of all nonnegative linear combinations

$$C = \{\mathbf{z} \in \mathbb{R}^m : \mathbf{z} = x_1\mathbf{a}_1 + \dots + x_n\mathbf{a}_n, x_i \geq 0\},$$

and this set is a convex cone. The assumption that there is no nonnegative \mathbf{x} such that $\mathbf{A}\mathbf{x} = \mathbf{b}$ means that $\mathbf{b} \notin C$. By Lemma 9.1, there exists a linear hyperplane $H = \{\mathbf{x} : \langle \mathbf{y}, \mathbf{x} \rangle = 0\}$ such that $C \subset H_-$ and $\mathbf{b} \in \text{int}H_+$. Formulated differently, there exists a $\mathbf{y} \in \mathbb{R}^m$ such that

$$\forall \mathbf{z} \in C: \langle \mathbf{z}, \mathbf{y} \rangle \geq 0, \quad \langle \mathbf{b}, \mathbf{y} \rangle < 0.$$

Since every $\mathbf{z} \in C$ has the form $\mathbf{z} = \sum_{i=1}^n x_i \mathbf{a}_i$ with $x_i \geq 0$, the relation

$$\langle \mathbf{z}, \mathbf{y} \rangle = \sum_{i=1}^n x_i \langle \mathbf{a}_i, \mathbf{y} \rangle \geq 0$$

for all $\mathbf{x} \geq 0$ is equivalent to the condition that $\langle \mathbf{a}_i, \mathbf{y} \rangle \geq 0$ for $1 \leq i \leq n$, which again is equivalent to $\mathbf{A}^\top \mathbf{y} \geq 0$. This concludes the proof. \square

The following consequence is perhaps a more familiar form of Farkas' Lemma.

Corollary 9.3. Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$, there exists a vector $\mathbf{x} \neq 0$ such that

$$\mathbf{A}\mathbf{x} \leq \mathbf{b}$$

if and only if there is no \mathbf{y} such that

$$\mathbf{y} \geq 0, \quad \mathbf{A}^\top \mathbf{y} = 0, \quad \langle \mathbf{y}, \mathbf{b} \rangle < 0.$$

Proof. Consider the matrix

$$\mathbf{A}' := (\mathbf{A} \quad -\mathbf{A} \quad \mathbf{I}),$$

where \mathbf{I} is the $m \times m$ identity matrix. A nonnegative solution of $\mathbf{A}'\mathbf{x}' = \mathbf{b}$ has the form $\mathbf{x}' = (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)^\top$, and implies $\mathbf{A}(\mathbf{x}_1 - \mathbf{x}_2) + \mathbf{x}_3 = \mathbf{b}$. Therefore, such a solution \mathbf{x}' exists if and only if the system

$$\mathbf{A}\mathbf{x} \leq \mathbf{b}$$

has a solution. Applying Theorem 9.2, the complementary condition is

$$A'^\top \mathbf{y} \geq 0, \quad \langle \mathbf{b}, \mathbf{y} \rangle < 0,$$

which in terms of \mathbf{A} translates to

$$\mathbf{A}^\top \mathbf{y} = \mathbf{0}, \quad \mathbf{y} \geq 0, \quad \langle \mathbf{b}, \mathbf{y} \rangle < 0.$$

This concludes the proof. \square

One more important corollary will be given without proof.

Corollary 9.4. *Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. Then for $\delta > 0$ and every vector $\mathbf{x} \in \mathbb{R}^n$ with $\mathbf{A}\mathbf{x} \leq \mathbf{b}$, $\langle \mathbf{c}, \mathbf{x} \rangle \leq \delta$ holds if and only if there exists $\mathbf{y} \in \mathbb{R}^m$ such that*

$$\mathbf{y} \geq 0, \quad \mathbf{A}^\top \mathbf{y} = \mathbf{c}, \quad \langle \mathbf{y}, \mathbf{b} \rangle \leq \delta.$$

9.2 Linear programming duality

After studying the feasible sets of linear programming, the polyhedra, we now return to linear programming itself, in the form

$$\text{maximize } \langle \mathbf{c}, \mathbf{x} \rangle \quad \text{subject to } \mathbf{A}\mathbf{x} \leq \mathbf{b}. \quad (9.1)$$

Geometrically this amounts to moving the hyperplane orthogonal to \mathbf{c} to the highest level along \mathbf{c} , under the condition that it still intersects $P = \{\mathbf{x} : \mathbf{A}\mathbf{x} \leq \mathbf{b}\}$.

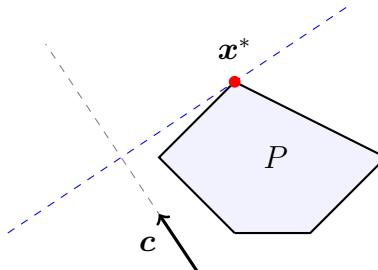


Figure 9.17: Geometry of linear programming

The famous duality theorem for linear programming states that if the maximum of (9.1) exists, then it coincides with the solution of a *dual* linear programming problem.

Theorem 9.5. *Let $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$ and $\mathbf{c} \in \mathbb{R}^n$. Then the optimal value of*

$$\text{maximize } \langle \mathbf{c}, \mathbf{x} \rangle \quad \text{subject to } \mathbf{A}\mathbf{x} \leq \mathbf{b} \quad (\text{P})$$

coincides with the optimal value of

$$\text{minimize } \langle \mathbf{b}, \mathbf{y} \rangle \quad \text{subject to } \mathbf{A}^\top \mathbf{y} = \mathbf{c}, \quad \mathbf{y} \geq 0, \quad (\text{D})$$

provided both (P) and (D) have a finite solution.

The problem (P) is called the *primal* problem, and (D) the *dual* problem.

Example 9.6. Consider the simple problem

$$\text{maximize } x_1 \quad \text{subject to } x_1 + x_2 \leq 1, x_1 \geq 0, x_2 \geq 0.$$

The dual problem is

$$\text{minimize } y_1 \quad \text{subject to } y_1 - y_2 = 1, y_1 - y_3 = 0, y_1 \geq 0, y_2 \geq 0, y_3 \geq 0.$$

For the proof we need the following observation.

Lemma 9.7. Let $P \subset \mathbb{R}^n$ be a polyhedron and $\mathbf{c} \in \mathbb{R}^n$ such that $\sup_{\mathbf{x} \in P} \langle \mathbf{c}, \mathbf{x} \rangle$ is finite. Then the supremum is attained, that is, it is a maximum.

Proof of Theorem 9.5. Let $P = \{\mathbf{x} : \mathbf{A}\mathbf{x} \in \mathbf{b}\}$ and $D = \{\mathbf{y} \in \mathbb{R}^m : \mathbf{A}^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq 0\}$. If $\mathbf{x} \in P$ and $\mathbf{y} \in Q$, then

$$\langle \mathbf{c}, \mathbf{x} \rangle = \langle \mathbf{A}^\top \mathbf{y}, \mathbf{x} \rangle = \langle \mathbf{y}, \mathbf{A}\mathbf{x} \rangle \leq \langle \mathbf{y}, \mathbf{b} \rangle,$$

so that in particular

$$\max_{\mathbf{x} \in P} \langle \mathbf{c}, \mathbf{x} \rangle \leq \min_{\mathbf{y} \in Q} \langle \mathbf{b}, \mathbf{y} \rangle,$$

which shows one inequality. To show the other inequality, set $\delta = \max_{\mathbf{x} \in P} \langle \mathbf{c}, \mathbf{x} \rangle$. By definition, if $\mathbf{A}\mathbf{x} \leq \mathbf{b}$, then $\langle \mathbf{c}, \mathbf{x} \rangle \leq \delta$. By Corollary 9.4, there exists a vector $\mathbf{y} \in \mathbb{R}^m$ such that

$$\mathbf{y} \geq 0, \quad \mathbf{A}^\top \mathbf{y} = \mathbf{c}, \quad \langle \mathbf{b}, \mathbf{y} \rangle \leq \delta.$$

In particular,

$$\min_{\mathbf{y} \in Q} \langle \mathbf{b}, \mathbf{y} \rangle \leq \delta = \max_{\mathbf{x} \in P} \langle \mathbf{c}, \mathbf{x} \rangle.$$

This finishes the proof. □

Lecture 10

So far we have seen aspects of the geometry of linear programming, discussed the feasible sets (polyhedra) and presented a useful duality theorem: the optimal value of

$$\text{maximize } \langle \mathbf{c}, \mathbf{x} \rangle \quad \text{subject to } \mathbf{A}\mathbf{x} \leq \mathbf{b} \tag{P}$$

coincides with the optimal value of

$$\text{minimize } \langle \mathbf{b}, \mathbf{y} \rangle \quad \text{subject to } \mathbf{A}^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq 0, \tag{D}$$

provided both (P) and (D) have a finite solution.

We now turn attention to algorithms for linear programming. We present the main idea behind the classical Simplex Algorithm and then discuss the more modern class of Interior Point Algorithms in more detail. The latter can be shown to be efficient in a very precise sense (polynomial time) and generalize to non-linear convex optimization.

10.1 A first algorithm for linear programming

For $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$ let $P = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} \leq \mathbf{b}\}$ be the polyhedron of feasible points for (P). Recall that a vertex of P is a zero-dimensional face of P , or equivalently, a point that cannot be written as convex combination of two distinct points of P .

If the optimization problem (P) has a solution, then it can be shown (Problem Set 4) that there exists a vertex \mathbf{x}^* such that

$$\max_{\mathbf{x} \in P} \langle \mathbf{c}, \mathbf{x} \rangle = \langle \mathbf{c}, \mathbf{x}^* \rangle.$$

Intuitively this is clear by looking at a picture, as in Figure 10.18: Move a hyperplane orthogonal to the objective \mathbf{c} along this direction to the highest level that still intersects a polyhedron P , and try to imagine that this intersection *does not* contain a vertex.

This crucial observation turns linear programming into a *finite* problem: we only have to test the objective function on the finite set of vertices to determine the maximizer. We have also seen (Problem Set 3) that the vertices can be identified as the solutions of the systems of equations

$$\mathbf{A}_I \mathbf{x} = \mathbf{b}_I,$$

for which \mathbf{A}_I has full rank and $\mathbf{x} \in P$, where $I \subset \{1, \dots, m\}$ is a subset of indices with $|I| = d$, and $\mathbf{A}_I, \mathbf{b}_I$ are the matrix and vector arising from \mathbf{A} and \mathbf{b} by taking the rows indexed by I . This gives a first primitive algorithm for solving the optimization problem (P).

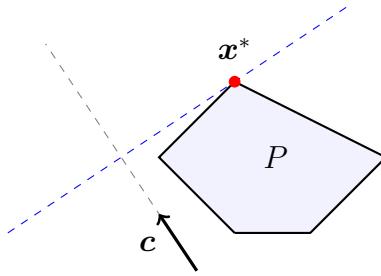
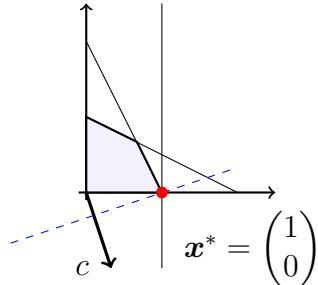


Figure 10.18: Geometry of linear programming

$$\begin{aligned} \text{maximize} \quad & x_1 - 3x_2 \\ \text{subject to} \quad & 0.5x_1 + x_2 \leq 1 \\ & x_1 + 0.5x_2 \leq 1 \\ & -x_1 \leq 0 \\ & -x_2 \leq 0 \\ & x_1 \leq 1 \end{aligned}$$



Example 10.1. Consider the linear programming problem

The matrix \mathbf{A} and vectors \mathbf{b} and \mathbf{c} are

$$\mathbf{A} = \begin{pmatrix} 0.5 & 1 \\ 1 & 0.5 \\ -1 & 0 \\ 0 & -1 \\ 1 & 0 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 1 \end{pmatrix}, \quad \mathbf{c} = \begin{pmatrix} 1 \\ -3 \end{pmatrix}.$$

We see that the minor $\mathbf{A}_{\{3,5\}} = \begin{pmatrix} -1 & 0 \\ 1 & 0 \end{pmatrix}$, corresponding to the two parallel hyperplanes $x_1 = 0$ and $x_1 = 1$, is singular and does not define a vertex, while all other minors are invertible. For the minor $I = \{1, 2\}$, $\mathbf{A}_I \mathbf{x} = \mathbf{b}_I$ has the form

$$\begin{aligned} 0.5x_1 + x_2 &= 1 \\ x_1 + 0.5x_2 &= 1, \end{aligned}$$

which has the solution $\mathbf{x} = (2/3, 2/3)^\top$. Since \mathbf{x} also satisfies all the other inequalities, it is a vertex. In a similar fashion we can find all the vertices as the points

$$\mathbf{x}_1 = \begin{pmatrix} 2/3 \\ 2/3 \end{pmatrix}, \quad \mathbf{x}_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad \mathbf{x}_3 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \quad \mathbf{x}_4 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$

The values of the objective function at these vertices are

$$\langle \mathbf{c}, \mathbf{x}_1 \rangle = -\frac{4}{3}, \quad \langle \mathbf{c}, \mathbf{x}_2 \rangle = 1, \quad \langle \mathbf{c}, \mathbf{x}_3 \rangle = -3, \quad \langle \mathbf{c}, \mathbf{x}_4 \rangle = 0.$$

It follows that the optimization problem has the optimal value 1 at $\mathbf{x}^* = (1, 0)^\top$.

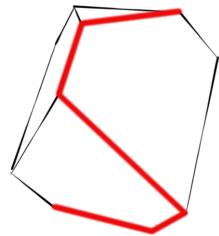
Unfortunately, this method requires solving up to $\binom{m}{n}$ systems of linear equations, which can be exponential in n in the worst case, making this method not practical.

Example 10.2. Consider the hypercube

$$\begin{aligned} -1 \leq x_1 &\leq 1 \\ \dots \\ -1 \leq x_n &\leq 1. \end{aligned}$$

It has 2^n vertices that need to be checked. Even though there are so many vertices, one can travel along edge between any two vertices in at most n steps.

The simplex algorithm is a way to search through the vertices in a more clever way than just listing all of them. The idea is to start with a vertex and see if there is a vertex connected to it by an edge that has a bigger objective value. If not, we are done. If there is a neighbouring vertex, we move there and continue the process. We keep walking along the vertices of the feasible polytope until we find an optimal one. The simplex algorithm is one of the most successful



algorithms around and usually very fast, even though there examples that show that its worst case running time can be exponential.

Lecture 11

We continue the study of algorithms for linear programming. Before doing so, we introduce a **standard form** for a linear programming problem, namely

$$\text{minimize } \langle \mathbf{c}, \mathbf{x} \rangle \quad \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}. \quad (\text{P})$$

for a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$ and $\mathbf{c} \in \mathbb{R}^n$. Any linear programming problem can be brought into this standard form. Start, for example, with the problem

$$\text{maximize } \langle \mathbf{c}, \mathbf{x} \rangle \quad \text{subject to } \mathbf{A}\mathbf{x} \leq \mathbf{b}, \quad (1.1)$$

in the form that we are used to. We first add a **slack variable** $\mathbf{s} \in \mathbb{R}^m$, so that we can reformulate the problem into an equivalent one of the form

$$\text{maximize } \langle \mathbf{c}, \mathbf{x} \rangle \quad \text{subject to } \mathbf{A}\mathbf{x} + \mathbf{s} = \mathbf{b}, \mathbf{s} \geq \mathbf{0}. \quad (1.2)$$

Define the big matrix $\mathbf{A}' \in \mathbb{R}^{m \times (n+n+m)}$ by

$$\mathbf{A}' := (\mathbf{A} \ -\mathbf{A} \ \mathbf{I}),$$

consider the big vector $\mathbf{x}' = (\mathbf{x}^+, \mathbf{x}^-, \mathbf{s})^\top$, where \mathbf{x}^+ is the *positive part*, with entries x_i if $x_i \geq 0$ and 0 else, and \mathbf{x}^- is the *negative part*, with entries $-x_i$ if $x_i < 0$ and 0 else (so that $\mathbf{x} = \mathbf{x}^+ - \mathbf{x}^-$), and set $\mathbf{c}' = (\mathbf{c}, -\mathbf{c}, \mathbf{0})^\top \in \mathbb{R}^{2n+m}$.

Then \mathbf{x} is a solution of (1.2) if and only if \mathbf{x}' is a solution of the problem

$$\text{minimize } \langle \mathbf{c}', \mathbf{x}' \rangle \quad \text{subject to } \mathbf{A}'\mathbf{x}' = \mathbf{b}, \mathbf{x}' \geq \mathbf{0},$$

which is in standard form (P). Note that since (P) has the form of a dual problem as derived in Lecture 9 and 10, its dual, in turn, has the form of a primal problem,

$$\text{maximize } \langle \mathbf{b}, \mathbf{y} \rangle \quad \text{subject to } \mathbf{A}^\top \mathbf{y} + \mathbf{s} = \mathbf{c}, \mathbf{s} \geq \mathbf{0}. \quad (\text{D})$$

For everything that follows, we assume $m \leq n$, as otherwise (D) is unbounded and (P) empty.

1.1 An optimality condition

We now work with a primal dual system *in standard form*,

$$\text{minimize } \langle \mathbf{c}, \mathbf{x} \rangle \quad \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}, \quad (\text{P})$$

and

$$\text{maximize } \langle \mathbf{b}, \mathbf{y} \rangle \quad \text{subject to } \mathbf{A}^\top \mathbf{y} + \mathbf{s} = \mathbf{c}, \mathbf{s} \geq \mathbf{0}. \quad (\text{D})$$

A useful consequence of the duality theorem, which states that a primal optimal solution is also dual optimal, is a characterisation of solution tuples $(\mathbf{x}^*, \mathbf{y}^*, \mathbf{s}^*)$. Note that for such an optimal pair,

$$\langle \mathbf{c}, \mathbf{x}^* \rangle = \langle \mathbf{y}^*, \mathbf{b} \rangle = \langle \mathbf{y}^*, \mathbf{A}\mathbf{x}^* \rangle = \langle \mathbf{A}^\top \mathbf{y}^*, \mathbf{x}^* \rangle,$$

where the first equality holds because of the duality theorem. By subtracting the last from the first expression above, we conclude

$$\langle \mathbf{x}^*, \mathbf{c} - \mathbf{A}^\top \mathbf{y}^* \rangle = \langle \mathbf{x}^*, \mathbf{s}^* \rangle = 0. \quad (\text{PD})$$

Since $\mathbf{x}^* \geq \mathbf{0}$ and $\mathbf{s}^* \geq \mathbf{0}$, each summand in (PD) is zero. This means that the individual components of \mathbf{s}^* and \mathbf{x}^* satisfy

$$x_i^* \cdot s_i^* = 0, \quad 1 \leq i \leq m.$$

Summarising, we have the following *optimality conditions* for linear programming: if vectors $(\mathbf{x}, \mathbf{y}, \mathbf{s})$ are primal/dual optimal solutions to a linear programming problem, then

$$\begin{aligned} \mathbf{Ax} + \mathbf{s} - \mathbf{b} &= \mathbf{0} \\ \mathbf{A}^\top \mathbf{y} - \mathbf{c} &= \mathbf{0} \\ y_i s_i &= 0, \quad 1 \leq i \leq m \\ \mathbf{y} &\geq \mathbf{0} \\ \mathbf{s} &\geq \mathbf{0}. \end{aligned} \quad (1.1)$$

We simplify this expression a bit further. Define the diagonal matrices

$$\mathbf{X} = \begin{pmatrix} x_1 & 0 & \cdots & 0 \\ 0 & x_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & x_n \end{pmatrix}, \quad \mathbf{S} = \begin{pmatrix} s_1 & 0 & \cdots & 0 \\ 0 & s_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & s_n \end{pmatrix}$$

and the vector $\mathbf{e} = (1, 1, \dots, 1)^\top$. Then the condition $x_i s_i = 0$ for $1 \leq i \leq n$ can be written concisely as $\mathbf{XSe} = \mathbf{0}$, and the whole system as

$$\begin{aligned} \mathbf{A}^\top \mathbf{y} + \mathbf{s} - \mathbf{c} &= \mathbf{0} \\ \mathbf{Ax} - \mathbf{b} &= \mathbf{0} \\ \mathbf{XSe} &= \mathbf{0} \\ \mathbf{x} &\geq \mathbf{0} \\ \mathbf{s} &\geq \mathbf{0}, \end{aligned} \quad (1.2)$$

Just as the optimality condition $\nabla f(\mathbf{x}) = \mathbf{0}$ serves as the basis for algorithms for unconstraint optimization, the optimality conditions for linear programming form the basis of the simplex method and of interior point methods. In the simplex method, the conditions are used to verify whether a candidate vertex is an optimal point, while primal/dual interior point methods view (15.1) as a multivariate function in $(\mathbf{x}, \mathbf{y}, \mathbf{s})$ for which a root satisfying the inequality constraints is sought using Newton's method or similar algorithms.

1.2 Newton's method for solving equations

Recall that we have seen Newton's method in two forms: as a minimization algorithm, and as a method for finding roots of a non-linear equation. In the latter, in the one-dimensional setting we want to solve an equation

$$f(x) = 0$$

and proceed by starting with an initial guess x_0 , and then computing successively

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}.$$

If f is the derivative of another function g , then this method is used to find a local minimizer of g . For a function $F = (f_1, \dots, f_n)^\top: \mathbb{R}^n \rightarrow \mathbb{R}^n$ we can use the same method, only that dividing by the derivative is replaced with multiplying with the inverse of the Jacobian matrix,

$$\mathbf{J}F(\mathbf{x}) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \cdots & \frac{\partial f_n}{\partial x_n} \end{pmatrix}.$$

Newton's method then starts with an initial guess \mathbf{x}_0 , and proceeds by

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{J}F(\mathbf{x}_k)^{-1} F(\mathbf{x}_k).$$

In practise, one does not compute the inverse, but solves a system of equations to get the update,

$$\mathbf{J}F(\mathbf{x}_k) \Delta \mathbf{x}_k = F(\mathbf{x}_k), \quad \mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \Delta \mathbf{x}_k. \quad (1.1)$$

Note that we allow for adjusting the step length in Newton's method. This can be rather convenient.

We can apply the multivariate Newton's method to the equalities in the optimality condition (15.1), by considering the function $F: \mathbb{R}^{2n+m} \rightarrow \mathbb{R}^{2n+m}$ defined by

$$F(\mathbf{x}, \mathbf{y}, \mathbf{s}) = \begin{pmatrix} \mathbf{A}^\top \mathbf{y} + \mathbf{s} - \mathbf{c} \\ \mathbf{A}\mathbf{x} - \mathbf{b} \\ \mathbf{X}\mathbf{S}\mathbf{e} \end{pmatrix}$$

We are then looking for a root $(\mathbf{x}^*, \mathbf{y}^*, \mathbf{s}^*)$ of this function that in addition satisfies the inequality constraints in (15.1). Before we get into the non-trivial issue enusing non-negativity, we first have a look at what the *update step* (1.1) looks like. Computing the Jacobian for F , the updates are computed by solving

$$\begin{pmatrix} \mathbf{0} & \mathbf{A}^\top & \mathbf{I} \\ \mathbf{A} & \mathbf{0} & \mathbf{0} \\ \mathbf{S} & \mathbf{0} & \mathbf{X} \end{pmatrix} \begin{pmatrix} \Delta \mathbf{x} \\ \Delta \mathbf{y} \\ \Delta \mathbf{s} \end{pmatrix} = \begin{pmatrix} \mathbf{A}^\top \mathbf{y} + \mathbf{s} - \mathbf{c} \\ \mathbf{A}\mathbf{x} - \mathbf{b} \\ \mathbf{X}\mathbf{S}\mathbf{e} \end{pmatrix}$$

Interior point method change the condition the condition $\mathbf{X}\mathbf{S}\mathbf{e} = \mathbf{0}$ to $\mathbf{X}\mathbf{S}\mathbf{e} = \tau \mathbf{e}$ for some $\tau > 0$, thus ensuring that we get non-negative solutions. Solving this problem for values $\tau \rightarrow 0$ then gives an approximation of the true solution. We discuss this in more detail in the next lecture.

Lecture 12

Recall the following primal dual standard forms of linear programming,

$$\text{minimize } \langle \mathbf{c}, \mathbf{x} \rangle \quad \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}. \quad (\text{P})$$

for a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$ and $\mathbf{c} \in \mathbb{R}^n$, and

$$\text{maximize } \langle \mathbf{b}, \mathbf{y} \rangle \quad \text{subject to } \mathbf{A}^\top \mathbf{y} + \mathbf{s} = \mathbf{c}, \mathbf{s} \geq \mathbf{0}. \quad (\text{D})$$

For what follows, we assume $m \leq n$, as otherwise (D) is unbounded and (P) empty.

Based on (P) and (D), we get the optimality conditions (see Lecture 11)

$$\begin{aligned} & \mathbf{A}^\top \mathbf{y} + \mathbf{s} - \mathbf{c} = \mathbf{0} \\ & \mathbf{A}\mathbf{x} - \mathbf{b} = \mathbf{0} \\ & \mathbf{XSe} = \mathbf{0} \\ & \mathbf{x} \geq \mathbf{0} \\ & \mathbf{s} \geq \mathbf{0}, \end{aligned} \quad (\text{O})$$

where \mathbf{X} is the diagonal matrix with the x_i in the diagonal, \mathbf{S} the diagonal matrix with the s_i on the diagonal, and \mathbf{e} the vector with all ones (\mathbf{XSe} is just a compact way of writing the vector with entries $x_i s_i$). If we define the function $\mathbb{R}^{2n+m} \rightarrow \mathbb{R}^{2n+m}$

$$F(\mathbf{x}, \mathbf{y}, \mathbf{s}) = \begin{pmatrix} \mathbf{A}^\top \mathbf{y} + \mathbf{s} - \mathbf{c} \\ \mathbf{A}\mathbf{x} - \mathbf{b} \\ \mathbf{XSe} \end{pmatrix} \quad (12.1)$$

then what we are looking for is a *root* $(\mathbf{x}^*, \mathbf{y}^*, \mathbf{s}^*)$ of this function, i.e., a point with $F(\mathbf{x}^*, \mathbf{y}^*, \mathbf{s}^*) = \mathbf{0}$, with the additional property that \mathbf{x}^* and \mathbf{s}^* be non-negative. For later reference, we record the form of the *Jacobian* matrix of F ,

$$\mathbf{JF}(\mathbf{x}, \mathbf{y}, \mathbf{s}) = \begin{pmatrix} \mathbf{0} & \mathbf{A}^\top & \mathbf{I} \\ \mathbf{A} & \mathbf{0} & \mathbf{0} \\ \mathbf{S} & \mathbf{0} & \mathbf{X} \end{pmatrix}.$$

12.1 Towards an efficient interior-point method

The challenge is to find an iterative method of solving this rootfinding problem while preserving the non-negativity constraints. We discuss three approaches, in increasing order of sophistication. The third of these approaches forms the basis of *primal dual interior point methods*.

A first attempt

A first attempt at solving the problem $F(\mathbf{x}, \mathbf{y}, \mathbf{s}) = \mathbf{0}$, with $\mathbf{x} \geq \mathbf{0}, \mathbf{s} \geq \mathbf{0}$, is to apply Newton's method with really small steps. For that, at each step we solve

$$\begin{pmatrix} \mathbf{0} & \mathbf{A}^\top & \mathbf{I} \\ \mathbf{A} & \mathbf{0} & \mathbf{0} \\ \mathbf{S}^{(k)} & \mathbf{0} & \mathbf{X}^{(k)} \end{pmatrix} \begin{pmatrix} \Delta \mathbf{x} \\ \Delta \mathbf{y} \\ \Delta \mathbf{s} \end{pmatrix} = \begin{pmatrix} \mathbf{c} - \mathbf{s}^{(k)} - \mathbf{A}^\top \mathbf{y}^{(k)} \\ \mathbf{b} - \mathbf{A} \mathbf{x}^{(k)} \\ -\mathbf{X}^{(k)} \mathbf{S}^{(k)} \mathbf{e} \end{pmatrix}.$$

and then compute

$$\begin{pmatrix} \mathbf{x}^{k+1} \\ \mathbf{y}^{k+1} \\ \mathbf{s}^{k+1} \end{pmatrix} = \begin{pmatrix} \mathbf{x}^k \\ \mathbf{y}^k \\ \mathbf{s}^k \end{pmatrix} + \alpha_k \begin{pmatrix} \Delta \mathbf{x} \\ \Delta \mathbf{y} \\ \Delta \mathbf{s} \end{pmatrix},$$

choosing the step length α_k so that the non-negativity of \mathbf{x} and \mathbf{s} remains. Note that above we have a + sign in front of alpha. That is because in the system of equations above, we took the negative $-F(\mathbf{x}_k, \mathbf{y}_k, \mathbf{s}_k)$.

Example 12.1. Consider the linear programming problem

$$\begin{aligned} \text{minimize} \quad & x_1 + 2x_2 - 2x_3 \\ & x_1 - 2x_3 \\ & x_2 - x_3 = -1 \\ & x_1 \geq 0, x_2 \geq 0, x_3 \geq 0. \end{aligned}$$

The data for this problem is given by

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & -2 \\ 0 & 1 & -1 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \quad \mathbf{c} = \begin{pmatrix} 1 \\ 2 \\ -2 \end{pmatrix}.$$

Based on this data, the function F given by

$$F(\mathbf{x}, \mathbf{y}, \mathbf{s}) = \begin{pmatrix} y_1 + s_1 - 1 \\ y_2 + s_2 - 2 \\ -2y_1 - y_2 + s_3 + 2 \\ x_1 - 2x_3 - 1 \\ x_2 - x_3 + 1 \\ x_1 s_1 \\ x_2 s_2 \\ x_3 s_3 \end{pmatrix}$$

and the Jacobian matrix

$$\mathbf{J}F(\mathbf{x}, \mathbf{y}, \mathbf{s}) = \begin{pmatrix} 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & -2 & -1 & 0 & 0 & 1 \\ 1 & 0 & -2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 & 0 & 0 & 0 & 0 \\ s_1 & 0 & 0 & 0 & 0 & x_1 & 0 & 0 \\ 0 & s_2 & 0 & 0 & 0 & 0 & x_2 & 0 \\ 0 & 0 & s_3 & 0 & 0 & 0 & 0 & x_3 \end{pmatrix}.$$

With this data at hand, we can easily use a Python program to solve the problem for us, making sure that the steplength is small enough. Starting with $\mathbf{x}^{(0)} = \mathbf{s}^{(0)} = (1, 1, 1)^\top$, we get the sequence of step lengths

$$\alpha_0 = 0.5455, \alpha_1 = 0.5455, \alpha_2 = 1,$$

with the corresponding sequence of \mathbf{x} vectors

$$\mathbf{x}^{(1)} = \begin{pmatrix} 1.1818 \\ 0 \\ 0.5455 \end{pmatrix}, \mathbf{x}^{(2)} = \begin{pmatrix} 2.1736 \\ 0 \\ 0.7934 \end{pmatrix}, \mathbf{x}^{(3)} = \begin{pmatrix} 3 \\ 0 \\ 1 \end{pmatrix}.$$

One verifies that $\mathbf{x}^{(3)}$ is in fact a vertex of the feasible set, and $\langle \mathbf{x}^{(3)}, \mathbf{c} \rangle = 1$ gives the optimal value.

A second attempt

In general, the method of choosing small step lengths can be slow. A variation would be to solve for

$$F(\mathbf{x}, \mathbf{y}, \mathbf{s}) = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ \tau \mathbf{e} \end{pmatrix}, \quad (12.1)$$

for a parameter $\tau > 0$. Given some fixed initial values $\mathbf{x} \geq 0$ and $\mathbf{s} \geq 0$, compute the *duality measure*

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i s_i$$

as the average of the products and using the *centering parameter* $\sigma \in (0, 1)$, set $\tau = \sigma\mu$. When aiming to solve (12.1) using Newton's method, we are aiming towards a solution $(\mathbf{x}^*, \mathbf{y}^*, \mathbf{s}^*)$ where instead of asking that $x_i^* s_i^* = 0$, we want that $x_i^* s_i^* = \sigma\mu$, a value that is strictly positive, but smaller than average of the initial value.

In an additional twist to this approach, after starting with an initial guess $(\mathbf{x}, \mathbf{y}, \mathbf{s})$, we perform only *one* Newton step, and then update the duality measure μ with the new values of \mathbf{x} and \mathbf{s} . This way we arrive at the following algorithm:

- Start with $(\mathbf{x}^{(0)}, \mathbf{y}^{(0)}, \mathbf{s}^{(0)})$;
- For each $k \geq 0$, compute

$$\mu^{(k)} = \frac{1}{n} \sum_{i=1}^n x_i s_i$$

and σ_k . Solve

$$\begin{pmatrix} \mathbf{0} & \mathbf{A}^\top & \mathbf{I} \\ \mathbf{A} & \mathbf{0} & \mathbf{0} \\ \mathbf{S}^{(k)} & \mathbf{0} & \mathbf{X}^{(k)} \end{pmatrix} \begin{pmatrix} \Delta \mathbf{x} \\ \Delta \mathbf{y} \\ \Delta \mathbf{s} \end{pmatrix} = \begin{pmatrix} \mathbf{c} - \mathbf{s}^{(k)} - \mathbf{A}^\top \mathbf{y}^{(k)} \\ \mathbf{b} - \mathbf{A} \mathbf{x}^{(k)} \\ -\mathbf{X}^{(k)} \mathbf{S}^{(k)} \mathbf{e} + \sigma \mu^{(k)} \mathbf{e} \end{pmatrix}$$

and compute

$$\begin{pmatrix} \mathbf{x}^{k+1} \\ \mathbf{y}^{k+1} \\ \mathbf{s}^{k+1} \end{pmatrix} = \begin{pmatrix} \mathbf{x}^k \\ \mathbf{y}^k \\ \mathbf{s}^k \end{pmatrix} + \alpha_k \begin{pmatrix} \Delta \mathbf{x} \\ \Delta \mathbf{y} \\ \Delta \mathbf{s} \end{pmatrix},$$

for a small enough $\alpha_k > 0$ to ensure non-negativity.

A third attempt

For the purpose of analysis, it is convenient to let an interior point method operate on vectors that satisfy the first two equalities in (O) exactly. Define the feasible and strictly feasible sets as

$$\begin{aligned} \mathcal{F} &= \{(\mathbf{x}, \mathbf{y}, \mathbf{s}) : \mathbf{A}^\top \mathbf{y} + \mathbf{s} = \mathbf{c}, \mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}, \mathbf{s} \geq \mathbf{0}\} \\ \mathcal{F}^\circ &= \{(\mathbf{x}, \mathbf{y}, \mathbf{s}) : \mathbf{A}^\top \mathbf{y} + \mathbf{s} = \mathbf{c}, \mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} > \mathbf{0}, \mathbf{s} > \mathbf{0}\} \end{aligned}$$

Restricting to points in \mathcal{F}° , the computation of the Newton update in the second approach would change to

$$\begin{pmatrix} \mathbf{0} & \mathbf{A}^\top & \mathbf{I} \\ \mathbf{A} & \mathbf{0} & \mathbf{0} \\ \mathbf{S}^{(k)} & \mathbf{0} & \mathbf{X}^{(k)} \end{pmatrix} \begin{pmatrix} \Delta \mathbf{x} \\ \Delta \mathbf{y} \\ \Delta \mathbf{s} \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ -\mathbf{X}^{(k)} \mathbf{S}^{(k)} \mathbf{e} + \sigma \mu^{(k)} \mathbf{e} \end{pmatrix}$$

In each iteration, a Newton step is taken in the direction of the *central path*. This is a curve in \mathcal{F}° defined as the set of solutions of

$$\begin{aligned} \mathbf{A}^\top \mathbf{y} + \mathbf{s} - \mathbf{c} &= \mathbf{0} \\ \mathbf{A}\mathbf{x} - \mathbf{b} &= \mathbf{0} \\ \mathbf{X}\mathbf{S}\mathbf{e} &= \tau \mathbf{e} \\ \mathbf{x} &> \mathbf{0} \\ \mathbf{s} &> \mathbf{0}, \end{aligned} \tag{12.2}$$

where $\tau > 0$. As $\tau \rightarrow 0$, any solution of this system will converge to an optimal primal-dual vector $(\mathbf{x}, \mathbf{y}, \mathbf{s})$ for the original linear programming problem. As we will see, practical primal-dual interior point methods will try to ensure that we always move within a neighbourhood of the central path. While this third approach lends itself well to analysis, one problem is that a starting point $(\mathbf{x}^{(0)}, \mathbf{y}^{(0)}, \mathbf{s}^{(0)}) \in \mathcal{F}^\circ$ may be hard to find.

Lecture 13

Primal-dual interior point methods aim to solve the problem

$$\text{minimize} \quad \langle \mathbf{c}, \mathbf{x} \rangle \quad \text{subject to} \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0} \quad (\text{P})$$

for a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$ and $\mathbf{c} \in \mathbb{R}^n$, by applying Newton-type iterations to the optimality conditions of linear programming. More precisely, we have seen the following algorithm. Recall the feasible sets

$$\begin{aligned} \mathcal{F} &= \{(\mathbf{x}, \mathbf{y}, \mathbf{s}) : \mathbf{A}^\top \mathbf{y} + \mathbf{s} = \mathbf{c}, \mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}, \mathbf{s} \geq \mathbf{0}\} \\ \mathcal{F}^\circ &= \{(\mathbf{x}, \mathbf{y}, \mathbf{s}) : \mathbf{A}^\top \mathbf{y} + \mathbf{s} = \mathbf{c}, \mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} > \mathbf{0}, \mathbf{s} > \mathbf{0}\} \end{aligned}$$

The a simple primal-dual interior point method can be described as follows.

- Start with $(\mathbf{x}^{(0)}, \mathbf{y}^{(0)}, \mathbf{s}^{(0)}) \in \mathcal{F}^\circ$;
- For each $k \geq 0$, compute the duality parameter

$$\mu^{(k)} = \frac{1}{n} \sum_{i=1}^n x_i s_i$$

and choose $\sigma_k \in [\sigma_{\min}, \sigma_{\max}]$. Solve

$$\begin{pmatrix} \mathbf{0} & \mathbf{A}^\top & \mathbf{I} \\ \mathbf{A} & \mathbf{0} & \mathbf{0} \\ \mathbf{S}^{(k)} & \mathbf{0} & \mathbf{X}^{(k)} \end{pmatrix} \begin{pmatrix} \Delta \mathbf{x} \\ \Delta \mathbf{y} \\ \Delta \mathbf{s} \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ -\mathbf{X}^{(k)} \mathbf{S}^{(k)} \mathbf{e} + \sigma \mu^{(k)} \mathbf{e} \end{pmatrix}$$

and compute

$$\begin{pmatrix} \mathbf{x}^{k+1} \\ \mathbf{y}^{k+1} \\ \mathbf{s}^{k+1} \end{pmatrix} = \begin{pmatrix} \mathbf{x}^k \\ \mathbf{y}^k \\ \mathbf{s}^k \end{pmatrix} + \alpha_k \begin{pmatrix} \Delta \mathbf{x} \\ \Delta \mathbf{y} \\ \Delta \mathbf{s} \end{pmatrix},$$

for a small enough $\alpha_k > 0$ to ensure non-negativity.

In each iteration, a Newton step is taken in the direction of the **path**. This is a curve in \mathcal{F}° defined as the set of solutions of

$$\begin{aligned} \mathbf{A}^\top \mathbf{y} + \mathbf{s} - \mathbf{c} &= \mathbf{0} \\ \mathbf{A}\mathbf{x} - \mathbf{b} &= \mathbf{0} \\ \mathbf{X}\mathbf{S}\mathbf{e} &= \tau\mathbf{e} \\ \mathbf{x} &> \mathbf{0} \\ \mathbf{s} &> \mathbf{0}, \end{aligned} \tag{13.1}$$

where $\tau > 0$.

13.1 Path-following methods

A path-following method tries to ensure that each iterate is *close* to the central path. What it means to be close to the central path depends on the neighbourhood we choose. Here, we will look at the (one-sided) ∞ -norm neighbourhood

$$\mathcal{N}_{-\infty}(\gamma) = \{(\mathbf{x}, \mathbf{y}, \mathbf{s}) \in \mathcal{F}^\circ : x_i s_i \geq \gamma\mu, 1 \leq i \leq n\}$$

for some $\gamma \in (0, 1]$ (say, $\gamma = 10^{-3}$). In words, each $x_i s_i$ has to be at least some small multiple of their average value. To see what this has to do with the ∞ -norm neighbourhood, consider the set of \mathbf{x} such that

$$\|\mathbf{X}\mathbf{S}\mathbf{e} - \mu\mathbf{e}\|_\infty \leq (1 - \gamma)\mu \iff \forall 1 \leq i \leq n, \gamma\mu \leq x_i s_i \leq 2 - \gamma,$$

and we are only interested in the lower inequality.

The so-called *long-step path-following* interior point method can then be described as follows.

- Start with $(\mathbf{x}^{(0)}, \mathbf{y}^{(0)}, \mathbf{s}^{(0)}) \in \mathcal{N}_{-\infty}(\gamma)$;
- For each $k \geq 0$, compute the duality parameter

$$\mu^{(k)} = \frac{1}{n} \sum_{i=1}^n x_i s_i$$

and choose $\sigma_k \in [\sigma_{\min}, \sigma_{\max}]$. Solve

$$\begin{pmatrix} \mathbf{0} & \mathbf{A}^\top & \mathbf{I} \\ \mathbf{A} & \mathbf{0} & \mathbf{0} \\ \mathbf{S}^{(k)} & \mathbf{0} & \mathbf{X}^{(k)} \end{pmatrix} \begin{pmatrix} \Delta\mathbf{x} \\ \Delta\mathbf{y} \\ \Delta\mathbf{s} \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ -\mathbf{X}^{(k)}\mathbf{S}^{(k)}\mathbf{e} + \sigma\mu^{(k)}\mathbf{e} \end{pmatrix}$$

and compute

$$\begin{pmatrix} \mathbf{x}^{k+1} \\ \mathbf{y}^{k+1} \\ \mathbf{s}^{k+1} \end{pmatrix} = \begin{pmatrix} \mathbf{x}^k \\ \mathbf{y}^k \\ \mathbf{s}^k \end{pmatrix} + \alpha_k \begin{pmatrix} \Delta\mathbf{x} \\ \Delta\mathbf{y} \\ \Delta\mathbf{s} \end{pmatrix},$$

for a small enough $\alpha_k \in [0, 1]$ is the largest value such that $(\mathbf{x}^{(k+1)}, \mathbf{y}^{(k+1)}, \mathbf{s}^{(k+1)}) \in \mathcal{N}_{-\infty}(\gamma)$.

Remark 13.1. As noted at the end of Lecture 12, to find an initial point in \mathcal{F}° might not be trivial. In practice one can therefore also use the algorithm described above using *infeasible* points, though in this case we have to make sure that the residual norms $\|\mathbf{b} - \mathbf{A}\mathbf{x}\|$ and $\|\mathbf{c} - \mathbf{s} - \mathbf{A}^\top \mathbf{y}\|$ remain bounded.

Visualising the algorithm

The feasible set \mathcal{F}° usually lives in a space that can't be easily visualised, but if the dual version is two-dimensional,

$$\{\mathbf{y} \in \mathbb{R}^2 : \mathbf{A}^\top \mathbf{y} + \mathbf{s} = \mathbf{c}, \mathbf{s} \geq 0\},$$

then we have a chance to see how the trajectories of the iterates in \mathbf{y} look like. Consider, for example, the linear programming problem whose dual is given by

$$\begin{aligned} & \text{maximize} && y_1 + y_2 \\ & \text{subject to} && 0.2py_1 + y_2 + s_p = 1 + 0.01p^2, \quad 0 \leq p \leq 10. \end{aligned}$$

The points $\mathbf{y} = (0, 0)^\top$, $\mathbf{x} = (1, \dots, 1)^\top/11$ and $\mathbf{s} = \mathbf{c} - \mathbf{A}^\top \mathbf{y}$ are strictly feasible starting points. The trajectory in the \mathbf{y} -plane of the algorithm (red path) and the constraint equations (blue lines) are given in the following diagram, where the parameter $\sigma = 0.5$ was used. It is

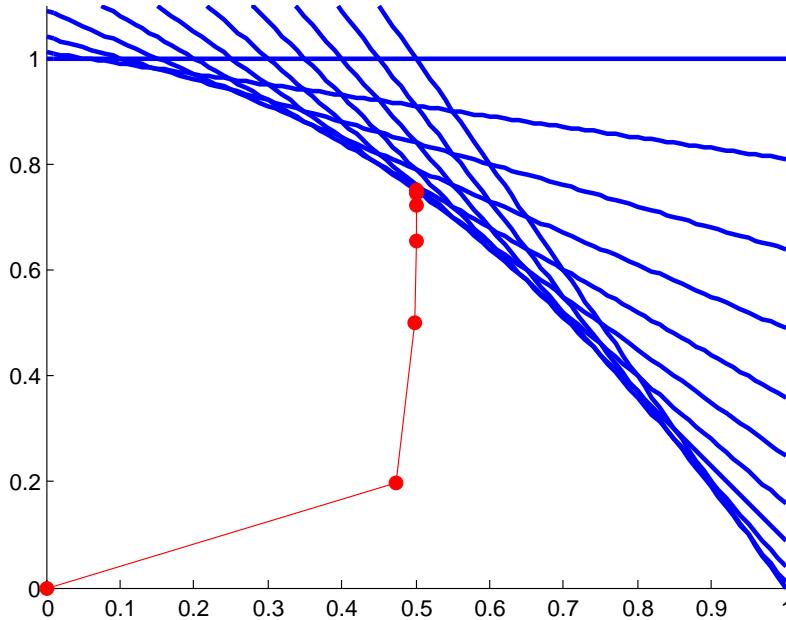


Figure 13.19: Trajectory of long-step path-following in the $y_1 - y_2$ plane.

instructive to play around with the parameter σ and to try to determine the form of the central path in this example.

Another way to visualise the trajectory is to plot the pairs $x_i s_i$ and $x_j y_j$ against each other. Figure 13.20 shows the trajectory of the above example in the $x_2 s_2 - x_5 s_5$ plane. Note that the central path, plotted in blue, is trivial in these coordinates, as it is defined by the property of the $x_i s_i = \tau$ being equal.

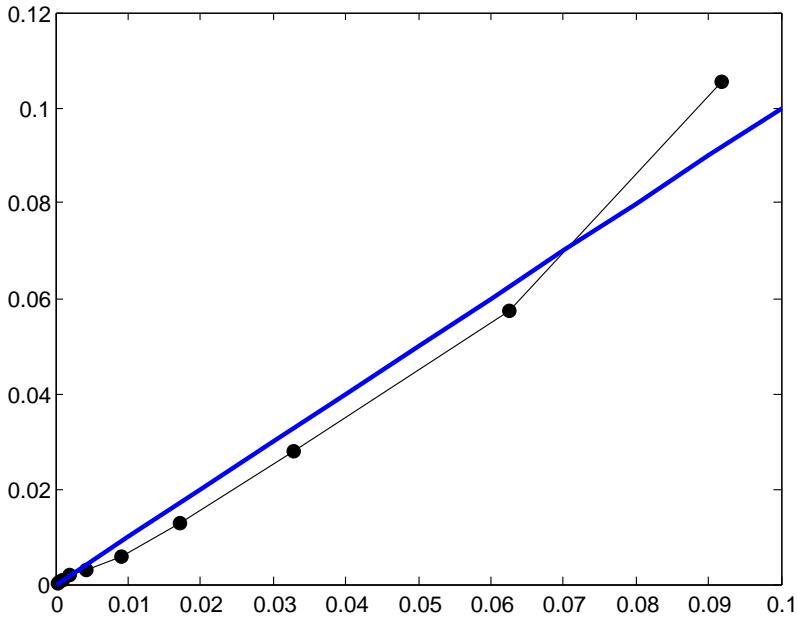


Figure 13.20: Trajectory and central path in $x_2s_2 - x_5s_5$ coordinates.

13.2 Analysis of Path-following

In the analysis of the long-step path-following algorithm, it is enough to establish that the duality measure $\mu^{(k)}$ converges to 0 as $k \rightarrow \infty$. The reason is that $\mu = 0$ forces all the products $x_i s_i = 0$, and since by design the other constraints are satisfied, this means that the sequence of points converges to a solution. The first theorem tells us that the μ_k decrease as k increases. An elementary proof is given in Theorem 14.3 in Nocedal and Wright. It depends crucially on the assumption that the iterates remain inside the neighbourhood $\mathcal{N}_\infty(\gamma)$ of the central path.

Theorem 13.2. *Given parameters γ , σ_{\min} and σ_{\max} , there is a constant $\delta > 0$, independent of n , such that*

$$\mu_{k+1} \leq \left(1 - \frac{\delta}{n}\right) \mu_k. \quad (13.1)$$

The next theorem gives a bound on the number of iterations needed to reduce the duality measure beyond any given ε .

Theorem 13.3. *Let $\varepsilon > 0$ and $\gamma \in (0, 1)$. Let $(\mathbf{x}^{(0)}, \mathbf{y}^{(0)}, \mathbf{s}^{(0)}) \in \mathcal{N}_\infty(\gamma)$ be a starting point such that the duality measure satisfies $\mu^{(0)} \leq \varepsilon^{-\kappa}$ for some constant κ . Then there is an index $K = O(n \log(1/\varepsilon))$ such that for all $k > K$,*

$$\mu_k \leq \varepsilon.$$

In particular, the long-step path-following algorithm converges.

Proof. Repeatedly applying (19.4), we get

$$\mu_k \leq \left(1 - \frac{\delta}{n}\right)^k \mu_0.$$

Taking logarithms on both sides,

$$\begin{aligned}\log \mu_k &\leq k \log \left(1 - \frac{\delta}{n}\right) + \log \mu_0 \leq k \log \left(1 - \frac{\delta}{n}\right) + \kappa \log \left(\frac{1}{\varepsilon}\right) \\ &\leq k \frac{-\delta}{n} + \kappa \log \left(\frac{1}{\varepsilon}\right).\end{aligned}$$

We have $\mu_k < \varepsilon$ if

$$-k \frac{\delta}{n} + \kappa \left(\frac{1}{\varepsilon}\right) \leq \log \varepsilon,$$

or equivalently, if

$$k \geq (1 + \kappa) \frac{n}{\delta} \log \left(\frac{1}{\varepsilon}\right) = K.$$

This was to be shown. □

Part IV

Non-linear Convex Optimization

Lecture 14

So far we have only dealt with constrained optimization problems where the objective and the constraints are linear. We now turn attention to general problems of the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{f}(\mathbf{x}) \leq \mathbf{0} \\ & && \mathbf{h}(\mathbf{x}) = \mathbf{0}, \end{aligned} \tag{1}$$

where $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{f} = (f_1, \dots, f_m)^\top$, $\mathbf{h} = (h_1, \dots, h_p)$, and the inequalities are componentwise. The problem (1) is *convex*, if f and the g_i are convex, and the h_j are linear. We also denote by $\text{dom}(f)$ the *domain* of f , which is the set of points \mathbf{x} where f takes a finite value. The feasible set

$$\mathcal{F} = \{\mathbf{x} : f_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0, 1 \leq i \leq m, 1 \leq j \leq \ell\}$$

for a convex constrained problem is a convex set.

14.1 Quadratic Programming and Portfolio Optimization

The simplest case of non-linear, constrained convex optimization is **quadratic programming**. Quadratic programming problems are of the form

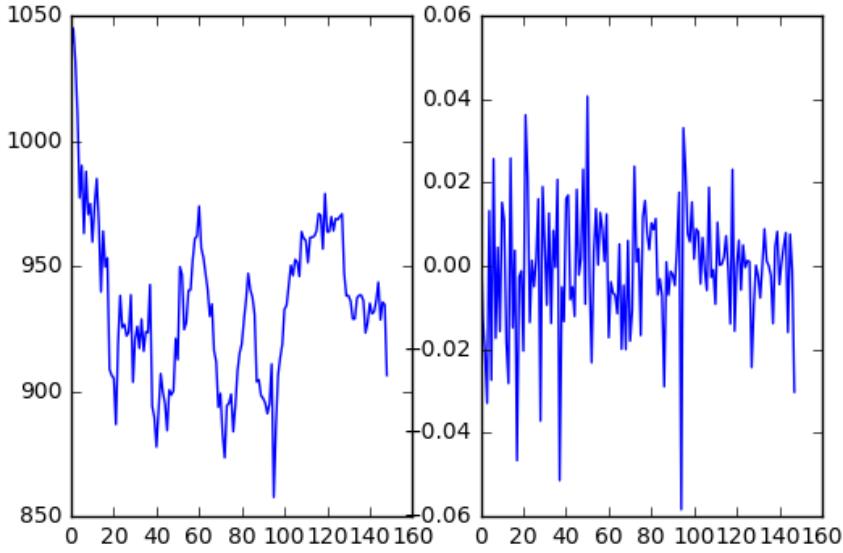
$$\begin{aligned} & \text{minimize} && \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A} \mathbf{x} \leq \mathbf{b}, \end{aligned}$$

with \mathbf{Q} symmetric and positive semidefinite. Note that this problem combines two problems we studied in some detail before: minimizing a quadratic function, and linear constraints. Just as we did for unconstrained minimization and for linear programming, we will derive optimality conditions for such problems. Before studying the theory, we first present an important application: portfolio optimization.

Mean-variance portfolio theory

If we invest an amount x^0 into a product at period 0, and at period 1 (for example, one day) the value is x^1 , then the **relative return** is defined as

$$r = \frac{x^1 - x^0}{x^0}.$$



The following figure shows the price movement and the returns of a stock from January 2016 until now.

As we can't predict the future, we usually work with the **expected return** $\mathbf{E}[r] = \mu$, which is a statistical estimate of the future return. One naive method of estimating the future return is by taking the average of past returns, but other more sophisticated methods are possible.

In **portfolio optimization**, we have a proportion x_i of our available funds that we want to invest in a stock i (in particular, as x_i measures a proportion, we have $\sum_{i=1}^n x_i = 1$). We may or may not allow $x_i < 0$, which would correspond to short-selling or borrowing. Given this allocation, the overall return is $\mathbf{x}^\top \mathbf{r}$, where $\mathbf{x} = (x_1, \dots, x_n)^\top$ and $\mathbf{r} = (r_1, \dots, r_n)^\top$ is the vector of (relative) returns. If $\boldsymbol{\mu} = \mathbf{E}[r]$ denotes the vector of expected returns, then the total expected return is

$$\mu = \mathbf{x}^\top \boldsymbol{\mu} = \sum_{i=1}^n x_i \mu_i.$$

The **risk** of an investment is measured in terms of the **variance** of the returns $r = \mathbf{x}^\top \mathbf{r}$. Let Σ denote the **covariance matrix**, where the (i, j) -th entry is $\text{Cov}(r_i, r_j) = \mathbf{E}[(r_i - \mu_i)(r_j - \mu_j)]$. The (i, j) -th entry measures how much products i and j are correlated. The **variance** of the returns r is then the quadratic function

$$\mathbf{x}^\top \Sigma \mathbf{x}.$$

A portfolio optimization problem either seeks to maximize the return while bounding the

risk,

$$\begin{aligned} & \text{maximize} && \mathbf{x}^\top \boldsymbol{\mu} \\ & \text{subject to} && \mathbf{x}^\top \Sigma \mathbf{x} \leq \sigma, \\ & && \sum_{i=1}^n x_i = 1 \\ & && x_i \geq 0, \end{aligned}$$

or minimize the risk given a certain target return μ ,

$$\begin{aligned} & \text{minimize} && \mathbf{x}^\top \Sigma \mathbf{x} \\ & \text{subject to} && \mathbf{x}^\top \boldsymbol{\mu} = \mu \\ & && \sum_{i=1}^n x_i = 1 \\ & && x_i \geq 0. \end{aligned}$$

Both of these problems are convex optimization problems. The constraints $x_i \geq 0$ mean that we are not allowed to short-sell; when dealing with futures or options, or if we are a large institutional investor, we may drop these constraints. As we will see later, dropping the inequality constraints from the second formulation allows the problem to be solved in closed form.

A third form of the portfolio optimization problem is to combine the expected return and the risk into a single objective function, as follows:

$$\begin{aligned} & \text{maximize} && \boldsymbol{\mu}^\top \mathbf{x} - \gamma \mathbf{x}^\top \Sigma \mathbf{x} \\ & \text{subject to} && \sum_{i=1}^n x_i = 1 \\ & && x_i \geq 0. \end{aligned}$$

Note that this is a convex quadratic problem: we can transform it into a minimization problem with the matrix Σ by changing the sign. The parameter γ is called a **risk aversion parameter**; it adjusts the level of risk we are willing to take. If $\gamma = 0$, then the quadratic term does not feature in the objective and we just aim to maximize the expected return. If γ is big, then the risk term weighs heavily on the objective function, and the optimal value will likely be one where $\mathbf{x}^\top \Sigma \mathbf{x}$ is small. By varying the value of γ , we get different expected return / risk (variance) trade-offs. The following code computes this trade-off curve for a portfolio of 10 stocks from the FT100 index. The mean and covariance were computed by averaging over the past 60 trading days (3 months). The x -axis is the **standard deviation**, which is given as the square root of the variance (risk).

```
In [22]: from datetime import datetime, date
import pandas as pd
from pandas_datareader import data, wb
import numpy as np
import matplotlib.pyplot as plt
from cvxpy import *
```

We first use the functionality of the Pandas module to load the price data from the Yahoo Finance web site.

```
In [23]: START = datetime(2016,1,1)
END = date.today()
TICKER = ['ADN', 'AZN', 'EZJ', 'GSK', 'ITV', 'LSE', 'TSCO', 'PSON', 'PRU', 'DGE']
mydata = data.DataReader('TSCO', "yahoo", START, END)
dates = mydata.index
df = pd.DataFrame(index=dates, columns=TICKER)
for x in TICKER:
    mydata = data.DataReader(x, "yahoo", START, END)
    df.loc[:,x] = mydata['Adj Close']
df = df.dropna()
df.head(2)
```

The following shows a small sample of the data loaded.

	ADN	AZN	EZJ	GSK	ITV	LSE	TSCO
Date							
2016-01-04	2.41197	31.926706	84.900002	37.813187	0.4	0.002	82.800593
2016-01-05	2.41197	32.404650	86.620003	38.028194	0.4	0.001	82.741275

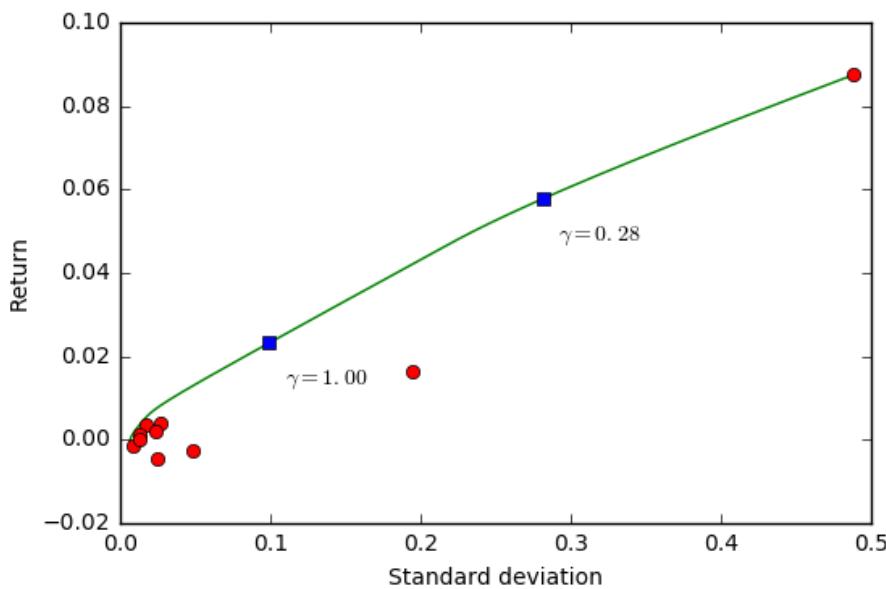
We next compute the returns and the mean and covariance.

```
In [24]: price_matrix = df.values
returns = (price_matrix[1:]-price_matrix[0:-1])/price_matrix[0:-1]
mu = np.mean(returns[-60:,:], axis=0)
Sigma = np.cov(returns[-60:,:].T)
n = 10
```

In the next step, we start the CVXPY engine and compute the mean-variance trade-off curve. We include the constraint $x \geq 0$ to disallow for short-selling.

```
In [25]: x = Variable(n)
gamma = Parameter(sign='positive')
ret = mu.T*x
risk = quad_form(x, Sigma)
prob = Problem(Maximize(ret - gamma*risk),
               [sum_entries(x) == 1,
                x >= 0])
```

```
In [26]: SAMPLES = 1000
risk_data = np.zeros(SAMPLES)
ret_data = np.zeros(SAMPLES)
gamma_vals = np.logspace(-2, 3, num=SAMPLES)
for i in range(SAMPLES):
    gamma.value = gamma_vals[i]
    prob.solve()
    risk_data[i] = sqrt(risk).value
    ret_data[i] = ret.value
```



```
In [27]: markers_on = [290, 400]
fig = plt.figure()
ax = fig.add_subplot(111)
plt.plot(risk_data, ret_data, 'g-')
for marker in markers_on:
    plt.plot(risk_data[marker], ret_data[marker], 'bs')
    ax.annotate(r"$\gamma = %.2f$" % gamma_vals[marker], xy=(risk_data[marker]+0.01, ret_data[marker]-0.01))
for i in range(n):
    plt.plot(sqrt(Sigma[i,i]).value, mu[i], 'ro')
plt.xlabel('Standard deviation')
plt.ylabel('Return')
plt.show()
```

The red dots represent the standard deviation and expected return of the individual stocks. The two blue dots indicate the standard deviation and expected return of the portfolio for two values of γ . Which value of γ we go for depends on how much risk we are willing to take: if we are risk-averse, we may prefer the $\gamma = 1$ value to the $\gamma = 0.028$ value, at the expense of smaller expected returns.

In practical applications there are a lot of other factors to be considered, such as whether the estimation procedure for the mean and covariance makes sense. In addition, it is often common to add terms that account for **transaction costs** into the objective function.

Lecture 15

In this lecture we study optimality conditions for convex problems of the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{f}(\mathbf{x}) \leq \mathbf{0} \\ & && \mathbf{h}(\mathbf{x}) = \mathbf{0}, \end{aligned} \tag{1}$$

where $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{f} = (f_1, \dots, f_m)^\top$, $\mathbf{h} = (h_1, \dots, h_p)^\top$, and the inequalities are componentwise. We assume that f and the f_i are convex, and the h_j are linear. It is also customary to write the conditions $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ as $\mathbf{Ax} = \mathbf{b}$, with $h_j(\mathbf{x}) = \mathbf{a}_j^\top \mathbf{x} - b_j$, \mathbf{a}_j being the j -th row of \mathbf{A} .

15.1 A first-order optimality condition

So far we have seen two examples of first order optimality conditions: for unconstraint optimization ($\nabla f(\mathbf{x}) = \mathbf{0}$) and for linear programming. We now generalize these to the setting of constrained convex optimization.

Theorem 15.1. *Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex, differentiable function, and*

$$\mathcal{F} = \{\mathbf{x} : f_i(\mathbf{x}) \leq 0, \mathbf{Ax} = \mathbf{b}\}$$

a feasible set, with f_i convex. Then \mathbf{x}^ is an optimal point of the optimization problem*

$$\text{minimize } f(\mathbf{x}) \quad \text{subject to } \mathbf{x} \in \mathcal{F}$$

if and only if for all $\mathbf{y} \in \mathcal{F}$,

$$\langle \nabla f(\mathbf{x}^*), \mathbf{y} - \mathbf{x}^* \rangle \geq 0. \tag{15.1}$$

Proof. Suppose \mathbf{x}^* is such that (1) holds. Then, since f is a convex function, for all $\mathbf{y} \in \mathcal{F}$ we have, by Theorem 2.10.1,

$$f(\mathbf{y}) \geq f(\mathbf{x}^*) + \langle \nabla f(\mathbf{x}^*), \mathbf{y} - \mathbf{x}^* \rangle \geq f(\mathbf{x}^*),$$

which shows that \mathbf{x}^* is a minimizer in \mathcal{F} . To show the opposite direction, assume that \mathbf{x}^* is a minimizer but that (1) does not hold. This means that there exists a $\mathbf{y} \in \mathcal{F}$ such that

$\langle \nabla f(\mathbf{x}^*), \mathbf{y} - \mathbf{x}^* \rangle < 0$. Since both \mathbf{x}^* and \mathbf{y} are in \mathcal{F} and \mathcal{F} is convex, any point $\mathbf{z}(\lambda) = (1 - \lambda)\mathbf{x}^* + \lambda\mathbf{y}$ with $\lambda \in [0, 1]$ is also in \mathcal{F} . At $\lambda = 0$ we have

$$\frac{df}{d\lambda} f(\mathbf{z}(\lambda))|_{\lambda=0} = \langle \nabla f(\mathbf{x}^*), \mathbf{y} - \mathbf{x}^* \rangle < 0.$$

Since the derivative at $\lambda = 0$ is negative, the function $f(\mathbf{z}(\lambda))$ is decreasing at $\lambda = 0$, and therefore, for small $\lambda > 0$, $f(\mathbf{z}(\lambda)) < f(\mathbf{z}(0)) = f(\mathbf{x}^*)$, in contradiction to the assumption that \mathbf{x}^* is a minimizer. \square

Example 15.2. In the absence of constraints, $\mathcal{F} = \mathbb{R}^n$, and the statement says that

$$\forall \mathbf{y} \in \mathbb{R}^n : \langle \nabla f(\mathbf{x}^*), \mathbf{y} - \mathbf{x}^* \rangle \geq 0.$$

Given \mathbf{y} such that $\langle \nabla f(\mathbf{x}^*), \mathbf{y} - \mathbf{x}^* \rangle \geq 0$, then replacing \mathbf{y} by $2\mathbf{x} - \mathbf{y}$ we also have the converse inequality, and therefore the optimality condition is equivalent to saying that $\nabla f(\mathbf{x}^*) = \mathbf{0}$. We therefore recover the well-known first order optimality condition from Lecture 2.

Geometrically, the first order optimality condition means that the set

$$\{\mathbf{x} : \langle \nabla f(\mathbf{x}^*), \mathbf{x} \rangle = \langle \nabla f(\mathbf{x}^*), \mathbf{x}^* \rangle\}$$

defines a supporting hyperplane to the set \mathcal{F} .

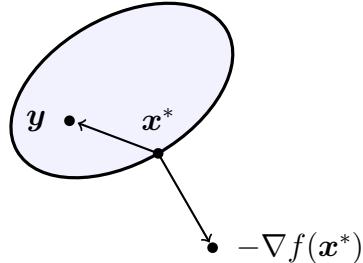


Figure 15.21: Optimality condition

15.2 Lagrangian duality

Recall the method of Lagrange multipliers. Given two functions $f(x, y)$ and $h(x, y)$, if the problem

$$\text{minimize } f(x, y) \quad \text{subject to } h(x, y) = 0$$

has a solution (x^*, y^*) , then there exists a parameter λ , the *Lagrange multiplier*, such that

$$\nabla f(x^*, y^*) = \lambda \nabla h(x^*, y^*). \tag{15.1}$$

In other words, if we define the *Lagrangian* as

$$\mathcal{L}(x, y, \lambda) = f(x, y) - \lambda h(x, y),$$

then (15.1) says that $\nabla \mathcal{L}(x^*, y^*, \lambda) = 0$ for some λ . The intuition is as follows. The set

$$M = \{(x, y) \in \mathbb{R}^2 : h(x, y) = 0\}$$

is a curve in \mathbb{R}^2 , and the gradient $\nabla h(x, y)$ is perpendicular to M at every point $(x, y) \in M$. For someone living inside M , a vector that is perpendicular to M is not visible, it is zero. Therefore the gradient $\nabla f(x, y)$ is zero as viewed from within M if it is perpendicular to M , or equivalently, a multiple of $\nabla h(x, y)$.

Alternatively, we can view the graph of $f(x, y)$ in three dimensions. A maximum or minimum of $f(x, y)$ along the curve defined by $h(x, y) = 0$ will be a point at which the direction of steepest ascent $\nabla f(x, y)$ is perpendicular to the curve $h(x, y) = 0$.

Example 15.3. Consider the function $f(x, y) = x^2y$ with the constraint $h(x, y) = x^2 + y^2 - 3$ (a circle of radius $\sqrt{3}$). The Lagrangian is the function

$$\mathcal{L}(x, y, \lambda) = x^2y - \lambda(x^2 + y^2 - 3).$$

Computing the partial derivatives gives the three equations

$$\begin{aligned}\frac{\partial}{\partial x} \mathcal{L} &= 2xy - 2\lambda x = 0 \\ \frac{\partial}{\partial y} \mathcal{L} &= x^2 - 2\lambda y = 0 \\ \frac{\partial}{\partial \lambda} \mathcal{L} &= x^2 + y^2 - 3 = 0.\end{aligned}$$

From the second equation we get $\lambda = \frac{x^2}{2y}$, and the first and third equations become

$$\begin{aligned}2xy - \frac{x^3}{y} &= 0 \\ x^2 + y^2 - 3 &= 0.\end{aligned}$$

Solving this system, we get six critical points $(\pm\sqrt{2}, \pm 1)$, $(0, \pm\sqrt{2})$. To find out which one of these is the minimizers, we just evaluate the function f on each of these.

We now turn to convex problems of the more general form

$$\begin{aligned}&\text{minimize} && f(\mathbf{x}) \\ &\text{subject to} && \mathbf{f}(\mathbf{x}) \leq \mathbf{0} \\ & && \mathbf{h}(\mathbf{x}) = \mathbf{0},\end{aligned}\tag{15.2}$$

Denote by \mathcal{D} the *domain* of all the functions f, f_i, h_j , i.e.,

$$\mathcal{D} = \text{dom}(f) \cap \text{dom}(f_1) \cap \cdots \cap \text{dom}(f_m) \cap \text{dom}(h_1) \cap \cdots \cap \text{dom}(h_p).$$

Assume that \mathcal{D} is not empty and let p^* be the optimal value of (15.2).

The *Lagrangian* of the system is defined as

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\lambda}^\top \mathbf{f}(\mathbf{x}) + \boldsymbol{\mu}^\top h(\mathbf{x}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{i=1}^p \mu_i h_i(\mathbf{x}).$$

The vectors $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$ are called the *dual variables* or *Lagrange multipliers* of the system. The domain of \mathcal{L} is $\mathcal{D} \times \mathbb{R}^m \times \mathbb{R}^p$.

Definition 15.4. The *Lagrange dual* of the problem (15.2) is the function

$$g(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\mathbf{x} \in \mathcal{D}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}).$$

If the Lagrangian \mathcal{L} is unbounded from below, then the value is $-\infty$.

The Lagrangian \mathcal{L} is linear in the λ_i and μ_j variables. The infimum of a family of linear functions is concave, so that the Lagrange dual is a concave function. Therefore the negative $-g(\boldsymbol{\lambda}, \boldsymbol{\mu})$ is a convex function.

Lemma 15.5. For any $\boldsymbol{\mu} \in \mathbb{R}^p$ and $\boldsymbol{\lambda} \geq \mathbf{0}$ we have

$$g(\boldsymbol{\lambda}, \boldsymbol{\mu}) \leq p^*.$$

Proof. Let \mathbf{x}^* be a feasible point for (15.2), that is,

$$f_i(\mathbf{x}^*) \leq 0, \quad h_j(\mathbf{x}^*) = 0, \quad 1 \leq i \leq m, \quad 1 \leq j \leq p.$$

Then for $\boldsymbol{\lambda} \geq \mathbf{0}$ and any $\boldsymbol{\mu}$, since each $h_j(\mathbf{x}^*) = 0$ and $\lambda_j f_j(\mathbf{x}^*) \leq 0$,

$$\mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}^*) + \sum_{j=1}^p \mu_j h_j(\mathbf{x}^*) \leq f(\mathbf{x}^*).$$

In particular,

$$g(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq f(\mathbf{x}^*).$$

Since this holds for all feasible \mathbf{x}^* , it holds in particular for the \mathbf{x}^* that minimizes (15.2), for which $f(\mathbf{x}^*) = p^*$. \square

A point $(\boldsymbol{\lambda}, \boldsymbol{\mu})$ with $\boldsymbol{\lambda} \geq \mathbf{0}$ and $(\boldsymbol{\lambda}, \boldsymbol{\mu}) \in \text{dom}(g)$ is called a *feasible point* of the dual problem.

The *Lagrange dual* of the optimization problem (15.2) is the problem

$$\text{maximize } g(\boldsymbol{\lambda}, \boldsymbol{\mu}) \quad \text{subject to } \boldsymbol{\lambda} \geq \mathbf{0}. \tag{15.3}$$

We have seen that if q^* is the optimal value of (15.3), then $q^* \leq p^*$, and the example above implies that in the special case of linear programming we actually have $q^* = p^*$. We will see that under certain conditions, we have $q^* = p^*$ for more general problems, but this is not always the case.

Lecture 16

For convex problems of the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{f}(\mathbf{x}) \leq \mathbf{0} \\ & && \mathbf{Ax} = \mathbf{b}, \end{aligned} \tag{1}$$

we introduced the Lagrangian $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ and defined the Lagrange dual as

$$g(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\mathbf{x} \in \mathcal{D}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}).$$

We saw that $g(\boldsymbol{\lambda}, \boldsymbol{\mu})$ is a lower bound on the optimal value of (1). Note that here we wrote the conditions $h_j(\mathbf{x}) = 0$ as system of linear equations $\mathbf{Ax} = \mathbf{b}$, since for the problem to be convex, we require that the h_j be linear functions.

Example 16.1. Consider a linear programming problem of the form

$$\begin{aligned} & \text{minimize} && \langle \mathbf{c}, \mathbf{x} \rangle \\ & \text{subject to} && \mathbf{Ax} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

The inequality constraints are $-x_i \leq 0$, while the equality constraints are $\mathbf{a}_i^\top \mathbf{x} - b_i$. The Lagrangian has the form

$$\begin{aligned} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) &= \langle \mathbf{c}, \mathbf{x} \rangle - \sum_{i=1}^n \lambda_i x_i + \sum_{j=1}^m \mu_j (\mathbf{a}_j^\top \mathbf{x} - b_j) \\ &= (\mathbf{c} - \boldsymbol{\lambda} + \mathbf{A}^\top \boldsymbol{\mu})^\top \mathbf{x} - \mathbf{b}^\top \boldsymbol{\mu}. \end{aligned}$$

The infimum over \mathbf{x} of this function is $-\infty$ unless $\mathbf{c} - \boldsymbol{\lambda} + \mathbf{A}^\top \boldsymbol{\mu} = \mathbf{0}$. The Lagrange dual is therefore

$$g(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \begin{cases} -\boldsymbol{\mu}^\top \mathbf{b} & \text{if } \mathbf{c} - \boldsymbol{\lambda} + \mathbf{A}^\top \boldsymbol{\mu} = \mathbf{0} \\ -\infty & \text{else.} \end{cases}$$

From Lemma 15.5 we conclude that

$$\max\{-\boldsymbol{\mu}^\top \mathbf{b} : \mathbf{c} - \boldsymbol{\lambda} + \mathbf{A}^\top \boldsymbol{\mu} = \mathbf{0}, \boldsymbol{\lambda} \geq \mathbf{0}\} \leq \min\{\mathbf{c}^\top \mathbf{x} : \mathbf{Ax} = \mathbf{0}, \mathbf{x} \geq \mathbf{0}\}.$$

Note that if we write $\mathbf{y} = -\boldsymbol{\mu}$ and $\mathbf{s} = \boldsymbol{\lambda}$, then we get the dual version of the linear programming problem we started out with, and in this case we know that

$$\max_{\boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\mu}} g(\boldsymbol{\lambda}, \boldsymbol{\mu}) = p^*.$$

16.1 Constraint qualification

In the example of linear programming, we have seen that the optimal value of the dual problem is equal to the optimal value of the primal problem. In general, we have

$$d^* = \sup_{\boldsymbol{\lambda} > \mathbf{0}, \boldsymbol{\mu}} g(\boldsymbol{\lambda}, \boldsymbol{\mu}) \leq \inf_{\mathbf{x} \in \mathcal{D}} \{f(\mathbf{x}) : f_i(\mathbf{x}) \leq 0, \mathbf{A}\mathbf{x} = \mathbf{b}\} = p^*.$$

Once certain conditions, called *constraint qualifications*, hold, we can ensure that *strong duality* holds, which means $d^* = p^*$. One particular such constraint qualification is Slater's Theorem.

Theorem 16.2. (*Slater conditions*) Assume that the interior of the domain \mathcal{D} of (1) is non-empty, that the problem (1) is convex, and that there exists a point $\mathbf{x} \in \mathcal{D}$ such that

$$f_i(\mathbf{x}) < 0, \quad 1 \leq i \leq m, \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \quad 1 \leq j \leq p.$$

Then $d^* = p^*$, the primal optimal value coincides with the dual optimal value.

Proof. (Optional) Assume \mathbf{A} has rank p (the number of rows). Assume moreover that p^* is finite, since if $p^* = -\infty$, then by weak duality we already have $d^* = p^*$.

Define the convex set

$$\mathcal{A} = \{(\mathbf{u}, \mathbf{v}, t) \in \mathbb{R}^m \times \mathbb{R}^p \times \mathbb{R} : \forall \mathbf{x} \in \mathcal{D}, f_i(\mathbf{x}) \leq u_i, \mathbf{A}\mathbf{x} - \mathbf{b} = \mathbf{v}, f(\mathbf{x}) \leq t\}.$$

Then the optimal value of (1) is

$$p^* = \inf \{t : (\mathbf{0}, \mathbf{0}, t) \in \mathcal{A}\}.$$

Define the convex set \mathcal{B} as

$$\mathcal{B} = \{(0, 0, s) \in \mathbb{R}^m \times \mathbb{R}^p \times \mathbb{R} : s < p^*\}.$$

The sets \mathcal{A} and \mathcal{B} are disjoint. To see this, assume to the contrary that there is a point $\mathbf{w} \in \mathcal{A} \cap \mathcal{B}$. Then since $\mathbf{w} \in \mathcal{B}$, $\mathbf{w} = (0, 0, s)$ with $s < p^*$, but also since $\mathbf{w} \in \mathcal{A}$, there exists an $\mathbf{x} \in \mathcal{D}$ with $f_i(\mathbf{x}) \leq 0$, $\mathbf{A}\mathbf{x} - \mathbf{b} = \mathbf{0}$, and $f(\mathbf{x}) \leq s < p^*$, in contradiction to the optimality of p^* as a value.

By the separating hyperplane theorem, there exist a hyperplane separating \mathcal{A} and \mathcal{B} (but not necessarily strictly!), defined by a vector $(\tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\mu}}, \nu) \neq \mathbf{0}$ and $\alpha \neq 0$ with the property that

$$(u, v, t) \in \mathcal{A} \implies \boldsymbol{\lambda}^\top \mathbf{u} + \boldsymbol{\mu}^\top \mathbf{v} + \nu t \geq \alpha \tag{16.1}$$

and

$$(u, v, t) \in \mathcal{B} \implies \boldsymbol{\lambda}^\top \mathbf{u} + \boldsymbol{\mu}^\top \mathbf{v} + \nu t \leq \alpha. \tag{16.2}$$

If $\tilde{\boldsymbol{\lambda}} < \mathbf{0}$ or $\nu < 0$, we could make the right-hand side of (16.1) arbitrary small, contradicting the bound by α . It follows that $\tilde{\boldsymbol{\lambda}} \geq \mathbf{0}$ and $\nu \geq 0$. Condition (16.2) simply means that $\nu t \leq \alpha$

for all $t < p^*$, so that $\nu p^* \leq \alpha$. Combining this bound with (16.1), we get the two inequalities, valid for any $\mathbf{x} \in \mathcal{D}$,

$$\sum_{i=1}^m \tilde{\lambda}_i f_i(\mathbf{x}) + \tilde{\boldsymbol{\mu}}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) + \nu f(\mathbf{x}) \geq \alpha \geq \nu p^*. \quad (16.3)$$

Note that the left-hand side has the form of a Lagrangian function scaled by ν . If $\nu > 0$ we can divide by μ and set $\lambda_i = \tilde{\lambda}_i/\nu$, $\mu_i = \tilde{\mu}_i/\nu$, to obtain

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \geq p^*.$$

By weak duality we have $p^* \geq g(\boldsymbol{\lambda}, \boldsymbol{\mu})$, so that we get strong duality if $\nu > 0$. If $\nu = 0$, then (16.3) implies

$$\sum_{i=1}^m \tilde{\lambda}_i f_i(\mathbf{x}) + \tilde{\boldsymbol{\mu}}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) \geq 0, \quad (16.4)$$

and for a point $\tilde{\mathbf{x}}$ satisfying the conditions $\mathbf{A}\mathbf{x} = \mathbf{b}$ and $f_i(\mathbf{x}) < 0$ for $1 \leq i \leq m$, this means that $\tilde{\boldsymbol{\lambda}} = 0$. As $(\tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\mu}}, \nu) \neq 0$, we have $\boldsymbol{\mu} \neq 0$. Since $\tilde{\boldsymbol{\mu}}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) \geq 0$ and $= 0$ for some $\tilde{\mathbf{x}}$ in the interior of \mathcal{D} , there must be a \mathbf{x} in the interior such that $\tilde{\boldsymbol{\mu}}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) < 0$, in contradiction to (16.4) (unless $\tilde{\boldsymbol{\mu}}^\top \mathbf{A} = 0$, which would contradict the condition that \mathbf{A} has maximal rank p). \square

Example 16.3. The problem minimize e^{-x} subject to $x^2/y \leq 0$, $y \geq 0$ is an example of a convex problem that does not satisfy strong duality.

Example 16.4. Consider the problem

$$\text{minimize } \mathbf{x}^\top \mathbf{x} \text{ subject to } \mathbf{A}\mathbf{x} = \mathbf{b}.$$

The Lagrangian is $\mathcal{L}(\mathbf{x}, \boldsymbol{\mu}) = \mathbf{x}^\top \mathbf{x} + \boldsymbol{\mu}^\top (\mathbf{A}\mathbf{x} - \mathbf{b})$. For any $\boldsymbol{\mu}$, we can find the infimum

$$g(\boldsymbol{\mu}) = \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\mu})$$

by setting the derivative of the Lagrangian to \mathbf{x} to zero:

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\mu}) = 2\mathbf{x} + \mathbf{A}^\top \boldsymbol{\mu} = 0,$$

which gives the solution

$$\mathbf{x} = -\frac{1}{2} \mathbf{A}^\top \boldsymbol{\mu}.$$

The dual function is therefore

$$g(\boldsymbol{\mu}) = -\frac{1}{4} \boldsymbol{\mu}^\top \mathbf{A}^\top \mathbf{A} \boldsymbol{\mu} - \mathbf{b}^\top \boldsymbol{\mu}.$$

As the negative of a positive semidefinite quadratic function, it is concave. Moreover, we get the lower bound

$$-\frac{1}{4} \boldsymbol{\mu}^\top \mathbf{A}^\top \mathbf{A} \boldsymbol{\mu} - \mathbf{b}^\top \boldsymbol{\mu} \leq \inf \{ \mathbf{x}^\top \mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b} \}.$$

The problem we started out with is convex, and if we assume that there exists a feasible primal point, then the above inequality is in fact an equality by Slater's conditions.

16.2 Karush-Kuhn-Tucker optimality conditions

Consider now a not necessarily convex problem of the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{f}(\mathbf{x}) \leq \mathbf{0} \\ & && \mathbf{h}(\mathbf{x}) = \mathbf{0}, \end{aligned} \tag{16.1}$$

If p^* is the optimal solution of (16.1) and $(\boldsymbol{\lambda}, \boldsymbol{\mu})$ dual variables, then we have seen that (this holds even in the non-convex case)

$$p^* \geq g(\boldsymbol{\lambda}, \boldsymbol{\mu}).$$

From this follows that for any primal feasible point \mathbf{x} ,

$$f(\mathbf{x}) - p^* \leq f(\mathbf{x}) - g(\boldsymbol{\lambda}, \boldsymbol{\mu}).$$

The difference $f(\mathbf{x}) - g(\boldsymbol{\lambda}, \boldsymbol{\mu})$ between the primal objective function at a primal feasible point and the dual objective function at a dual feasible point is called the *duality gap* at \mathbf{x} and $(\boldsymbol{\lambda}, \boldsymbol{\mu})$. For any such points we know that

$$p^*, q^* \in [g(\boldsymbol{\lambda}, \boldsymbol{\mu}), f(\mathbf{x})],$$

and if the gap is small we have a good approximation of the primal and dual optimal values. The duality gap can be used in iterative algorithms to define stopping criteria: if the algorithm generates a sequence of primal-dual variables $(\mathbf{x}^k, \boldsymbol{\lambda}^k, \boldsymbol{\mu}^k)$, then we can stop if the duality gap is less than, say, a predefined tolerance ε .

Now suppose that we have points $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ such that the duality gap is zero. Then

$$\begin{aligned} f(\mathbf{x}^*) &= g(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) \\ &= \inf_{\mathbf{x}} \left(f(\mathbf{x}) + \sum_{i=1}^m \lambda_i^* f_i(\mathbf{x}) + \sum_{j=1}^p \mu_j^* h_j(\mathbf{x}) \right) \\ &\leq f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* f_i(\mathbf{x}^*) + \sum_{j=1}^p \mu_j^* h_j(\mathbf{x}^*) \\ &\leq f(\mathbf{x}^*), \end{aligned}$$

where the last inequality follows from the fact that $h_j(\mathbf{x}^*) = 0$ and $\lambda_i^* f_i(\mathbf{x}^*) \leq 0$ for $1 \leq j \leq p$ and $1 \leq i \leq m$. It follows that the inequalities are in fact equalities. From the identity

$$f(\mathbf{x}^*) = f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* f_i(\mathbf{x}^*)$$

and $\lambda_i^* \geq 0$ and $f_i(\mathbf{x}^*) \leq 0$ we also conclude that at such optimal points,

$$\lambda_i^* f_i(\mathbf{x}^*) = 0, \quad 1 \leq i \leq m.$$

This condition is known as *complementary slackness*. From the above we also see that \mathbf{x}^* minimizes the Lagrangian $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$, so that the gradient of that function is zero:

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = \mathbf{0}.$$

Collecting these conditions (primal and dual feasibility, complementary slackness, vanishing gradient), we arrive at a set of optimality conditions known as the Karush-Kuhn-Tucker (KKT) conditions.

Theorem 16.5. (KKT conditions) Let \mathbf{x}^* and $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ be primal and dual optimal solutions of (16.1) with zero duality gap. Then the following conditions are satisfied:

$$\begin{aligned} \mathbf{f}(\mathbf{x}^*) &\leq \mathbf{0} \\ \mathbf{h}(\mathbf{x}^*) &= \mathbf{0} \\ \boldsymbol{\lambda}^* &\geq \mathbf{0} \\ \lambda_i^* f_i(\mathbf{x}^*) &= 0, \quad 1 \leq i \leq m \\ \nabla_{\mathbf{x}} f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* \nabla_{\mathbf{x}} f_i(\mathbf{x}^*) + \sum_{j=1}^p \mu_j^* \nabla_{\mathbf{x}} h_j(\mathbf{x}^*) &= \mathbf{0}. \end{aligned}$$

Moreover, if the problem is convex, then any points satisfying the KKT conditions have zero duality gap.

Lecture 17

In this lecture we introduce interior point methods for solving nonlinear convex optimization problems of the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{f}(\mathbf{x}) \leq \mathbf{0} \\ & && \mathbf{A}\mathbf{x} = \mathbf{b}. \end{aligned} \tag{1}$$

with $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{p \times n}$ and $\mathbf{f} = (f_1, \dots, f_m)^\top: \mathbb{R}^n \rightarrow \mathbb{R}^m$. The particular form of interior point methods we discuss is the *barrier method*, which differs slightly from the primal-dual method discussed for linear programming.

Recall the Karush-Kuhn-Tucker (KKT) conditions

$$\begin{aligned} & \mathbf{f}(\mathbf{x}^*) \leq \mathbf{0} \\ & \mathbf{A}\mathbf{x}^* = \mathbf{b} \\ & \boldsymbol{\lambda}^* \geq \mathbf{0} \\ & \lambda_i^* f_i(\mathbf{x}^*) = 0, \quad 1 \leq i \leq m \\ & \nabla_{\mathbf{x}} f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* \nabla_{\mathbf{x}} f_i(\mathbf{x}^*) + \mathbf{A}^\top \boldsymbol{\mu}^* = \mathbf{0}, \end{aligned} \tag{2}$$

where $\nabla_{\mathbf{x}}$ denotes the gradient of a function with respect to \mathbf{x} . As discussed in the previous lecture, for convex problems these are necessary and sufficient conditions for $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) \in \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^p$ being a set of optimal primal and dual solutions. By this we mean that \mathbf{x}^* is a minimizer of (1) and $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ is a maximizer of the Lagrange dual function $g(\boldsymbol{\lambda}, \boldsymbol{\mu})$ subject to $\boldsymbol{\lambda} \geq \mathbf{0}$, which in turn was defined as the infimum

$$g(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}),$$

where $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\lambda}^\top \mathbf{f}(\mathbf{x}) + \boldsymbol{\mu}^\top (\mathbf{A}\mathbf{x} - \mathbf{b})$ is the Lagrangian. In this section we make the additional assumption that the functions f and f_i , $1 \leq i \leq m$, are two time continuously differentiable. The reason is that we want to apply variants of Newton's method to the KKT conditions. For later reference we list some of the more popular forms of convex optimization problems that fall into our scope, with the associated KKT conditions.

Example 17.1. (LP) For a linear programming problem of the form

$$\begin{aligned} & \text{minimize} && \langle \mathbf{c}, \mathbf{x} \rangle \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

the KKT conditions have the form

$$\begin{aligned} & \mathbf{A}\mathbf{x} = \mathbf{b} \\ & \mathbf{A}^\top \mathbf{y} + \mathbf{s} = \mathbf{c} \\ & \mathbf{s} \geq \mathbf{0} \\ & \mathbf{x} \geq \mathbf{0} \\ & x_i s_i = 0, \quad 1 \leq i. \end{aligned}$$

Example 17.2. (QP) A *quadratic programming problem* is a problem with quadratic objective and linear constraints of the form

$$\begin{aligned} & \text{minimize} && \mathbf{x}^\top \mathbf{P}\mathbf{x} + \mathbf{q}^\top \mathbf{x} + r \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{B}\mathbf{x} \leq \mathbf{c}. \end{aligned}$$

One typical example is portfolio optimization.

Many other problems can be cast in the required form, and we will see later that the framework also encompasses semidefinite programming.

17.1 The logarithmic barrier

The idea of barrier methods is to replace the constraints of the optimization problem (1) with a parametrized system of equality constraints, such that the solutions of the equality constraint system can be found, and converge to the solution of the original system. The first observation is that we can replace the inequality constrained system

$$\text{minimize } f(\mathbf{x}) \quad \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{b}, \quad f_i(\mathbf{x}) \leq 0, \quad 1 \leq i \leq m,$$

with the equality constrained, but nonsmooth, problem

$$\text{minimize } f(\mathbf{x}) + \sum_{i=1}^m I_-(f_i(\mathbf{x})), \quad \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{b},$$

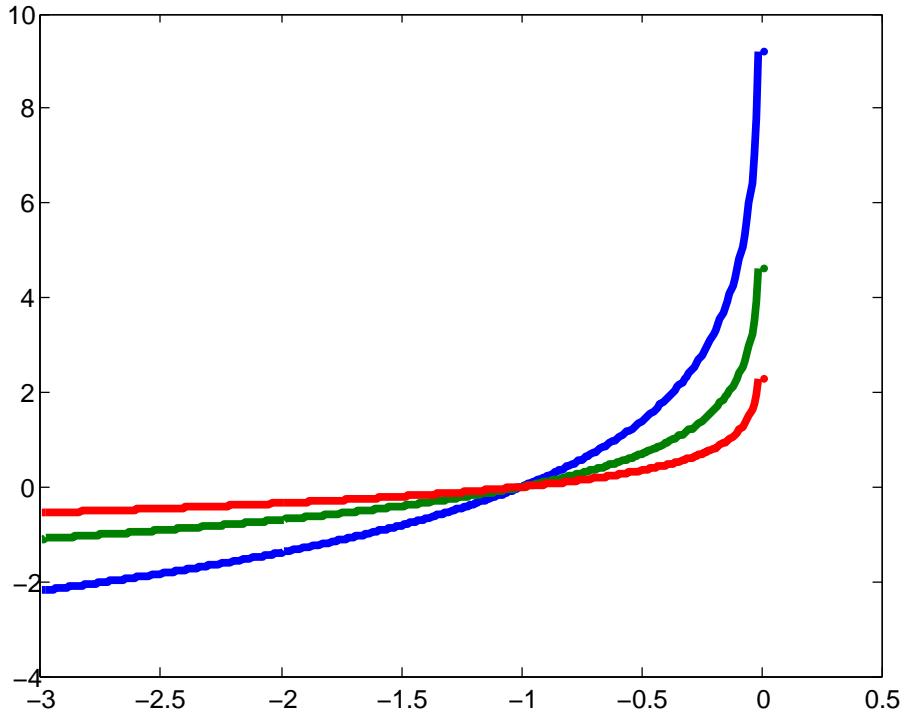
where

$$I_-(t) = \begin{cases} 0 & t \leq 0, \\ \infty & t > 0. \end{cases}$$

Clearly, the objective is $f(\mathbf{x})$ whenever the inequality constraints $f_i(\mathbf{x}) \leq 0$ are satisfied, and ∞ if not, so that points that do not satisfy the constraints are forbidden. Since this function is not easy to deal with analytically, we *approximate* it with the function

$$\hat{I}(u) := -\frac{1}{t} \ln(-u)$$

for suitable values t . The figure below shows the approximating functions for $t = 0.5, 1, 2$. As is easily seen, as t increases the approximation becomes better. An approximation to the



nonsmooth problem is given by the following equality constrained problem,

$$\begin{aligned} \text{minimize} \quad & f(\mathbf{x}) - \sum_{i=1}^m \frac{1}{t} \log(f_i(\mathbf{x})) \\ \text{subject to} \quad & \mathbf{A}\mathbf{x} = \mathbf{b}. \end{aligned}$$

The domain \mathcal{D} of this problem is the set of points such that the inequalities $f_i(\mathbf{x}) < 0$ are strictly satisfied. One can also verify that the objective function is convex, giving rise to a convex optimization problem. The function

$$\varphi(\mathbf{x}) = - \sum_{i=1}^m \log(f_i(\mathbf{x}))$$

is called the *logarithmic barrier function* of the system of inequalities $f_i(\mathbf{x}) < 0$: the reason is that it prevents candidates \mathbf{x} from becoming too close to the boundary of the feasible set by

becoming very large near it. With this notation, we can write the approximate problem as

$$\begin{aligned} & \text{minimize} && t f(\mathbf{x}) + \varphi(\mathbf{x}) \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b}. \end{aligned} \tag{17.1}$$

Note that here we multiplied the objective with t , which does not change the optimal points (but the values). For reference, we record the gradient and Hessian of φ ,

$$\begin{aligned} \nabla \varphi &= \sum_{i=1}^m \frac{1}{-f_i(\mathbf{x})} \nabla f_i(\mathbf{x}) \\ \nabla^2 \varphi &= \sum_{i=1}^m \frac{1}{f_i^2(\mathbf{x})} \nabla f_i(\mathbf{x}) \nabla f_i(\mathbf{x})^\top + \sum_{i=1}^m \frac{1}{-\nabla f_i(\mathbf{x})} \nabla^2 f_i(\mathbf{x}). \end{aligned}$$

Two questions are of importance:

1. How well does a solution of (17.1) approximate one of (1), and
2. How do we solve (17.1)?

17.2 The central path

The set of solutions $\mathbf{x}^*(t)$ of (17.1) for $t > 0$ is called the *central path*. Points on the central path are strictly feasible: they satisfy $\mathbf{A}\mathbf{x}^*(t) = \mathbf{b}$ and $f_i(\mathbf{x}^*(t)) < 0$. Moreover, by the Lagrange multiplier theorem, they satisfy the optimality condition

$$t \nabla f(\mathbf{x}) + \nabla \varphi(\mathbf{x}) + \mathbf{A}^\top \boldsymbol{\mu}^* = \mathbf{0}$$

for some suitable Lagrange multipliers $\boldsymbol{\mu} \in \mathbb{R}^p$. An important property of the central path is that we can derive dual points $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$, and therefore (by Lagrange duality) lower bounds on the optimal value p^* of (1), from any point \mathbf{x}^* on the central path. In fact, the central path can be derived from the solutions of a variant of the KKT conditions.

Lemma 17.3. *A point \mathbf{x} is equal to a point $\mathbf{x}^*(t)$ on the central path if and only if there exist dual multipliers λ and μ such that the following conditions are satisfied.*

$$\begin{aligned} \mathbf{f}(\mathbf{x}^*) &\leq \mathbf{0} \\ \mathbf{A}\mathbf{x}^* &= \mathbf{b} \\ \boldsymbol{\lambda}^* &\geq \mathbf{0} \\ -\lambda_i^* f_i(\mathbf{x}^*) &= \frac{1}{t}, \quad 1 \leq i \leq m \\ \nabla_{\mathbf{x}} f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* \nabla_{\mathbf{x}} f_i(\mathbf{x}^*) + \mathbf{A}^\top \boldsymbol{\mu}^* &= \mathbf{0}, \end{aligned}$$

The proof is left as an exercise. Note the analogy with the definition of the central path in the case of linear programming! Accordingly, the solution methods discussed in the context of linear programming carry over to the nonlinear case. We conclude by sketching the *barrier method*.

1. Start with a strictly feasible \mathbf{x} and $t := t^{(0)} > 0$, $\mu > 1$ and tolerance $\varepsilon > 0$;
2. Compute $\mathbf{x}^*(t)$ by minimizing $tf + \varphi$ subject to $\mathbf{A}\mathbf{x} = \mathbf{b}$, starting with \mathbf{x} ;
3. Update $\mathbf{x} = \mathbf{x}^*(t)$;
4. Stop if $m/t < \varepsilon$, otherwise set $t := \mu t$ (increase t).

The only unexplained part is how to solve the equality constrained minimization problem. One way is to use Newton's method on the optimality conditions of this problem.

Lecture 18

In this lecture we return to the task of classification. As seen earlier, examples include spam filters, letter recognition, or text classification. In this lecture we introduce a popular method for classification, **Support Vector Machines (SVMs)**, from the point of view of convex optimization.

18.1 Linear Support Vector Machines

In the simplest case there is a set of labels $\mathcal{Y} = \{-1, 1\}$ and the set of training points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is *linearly separable*: this means that there exists an affine hyperplane $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$ such that $h(\mathbf{x}_i) > 0$ if $y_i = 1$ and $h(\mathbf{x}_j) < 0$ if $y_j = -1$. We call the points for which $y_i = 1$ *positive*, and the ones for which $y_j = -1$ *negative*. The problem of finding such a hyperplane can be posed as a linear programming feasibility problem as follows: we look for a vector of *weights* \mathbf{w} and a *bias term* b (together a $(p+1)$ -dimensional vector) such that

$$\mathbf{w}^\top \mathbf{x}_i + b \geq 1, \text{ for } y_i = 1, \quad \mathbf{w}^\top \mathbf{x}_j + b \leq -1, \text{ for } y_j = -1.$$

Note that we can replace the $+1$ and -1 with any other positive or negative quantity by rescaling the \mathbf{w} and b , so this is just convention. We can also describe the two inequalities concisely as

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1 \geq 0. \tag{18.1}$$

A hyperplane separating the two point sets will in general not be unique. As we want to use the linear classifier on new, yet unknown data, we want to find a separating hyperplane with best possible **margin**. Let d_+ and d_- denote the distance of a separating hyperplane to the closest positive and closest negative point, respectively. The quantity $d = d_+ + d_-$ is then called the margin or the classifier, and we want to find a hyperplane with largest possible margin.

We next show that the margin for a separating hyperplane that satisfies (18.1) is $d = 2/\|\mathbf{w}\|_2$. Given a hyperplane H described in (18.1) and a point \mathbf{x} such that we have the equality $\mathbf{w}^\top \mathbf{x} + b = 1$ (the point is as close as possible to the hyperplane, also called a **support vector**), the distance of that point to the hyperplane can be computed by first taking the difference of \mathbf{x} with a point \mathbf{p} on H (an *anchor*), and then computing the dot product of $\mathbf{x} - \mathbf{p}$ with the unit vector $\mathbf{w}/\|\mathbf{w}\|$ orthogonal to H (see Calculus and Vectors A, math10121).

As anchor point \mathbf{p} we can just choose a multiple $c\mathbf{w}$ that is on the plane, i.e., that satisfies $\langle \mathbf{w}, c\mathbf{w} \rangle + b = 0$. This implies that $c = -b/\|\mathbf{w}\|^2$, and consequently $\mathbf{p} = -(b/\|\mathbf{w}\|^2)\mathbf{w}$. The

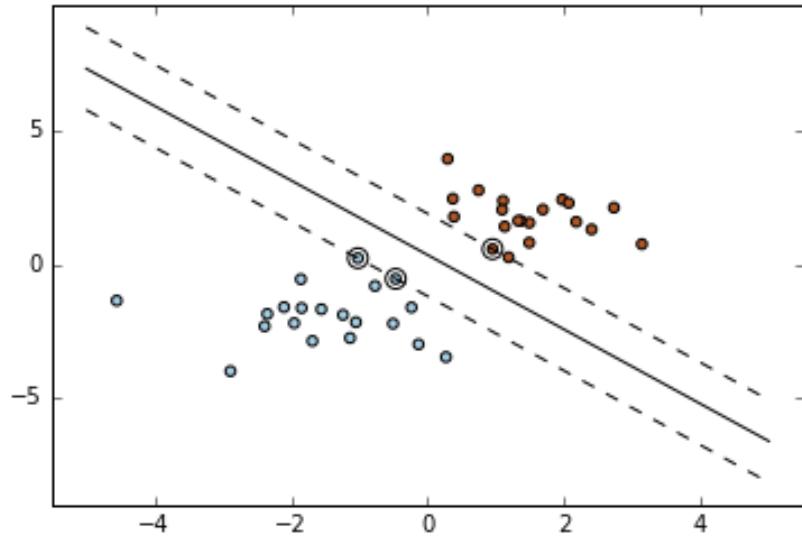


Figure 18.22: A hyperplane separating two sets of points with margin and support vectors.

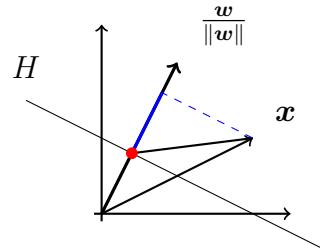


Figure 18.23: Computing the distance to the hyperplane

distance is then

$$\begin{aligned} d_+ &= \left\langle \mathbf{x} + \frac{b}{\|\mathbf{w}\|^2} \mathbf{w}, \frac{\mathbf{w}}{\|\mathbf{w}\|} \right\rangle = \frac{\langle \mathbf{x}, \mathbf{w} \rangle}{\|\mathbf{w}\|} + \frac{b}{\|\mathbf{w}\|^2} \langle \mathbf{w}, \frac{\mathbf{w}}{\|\mathbf{w}\|} \rangle \\ &= \frac{1-b}{\|\mathbf{w}\|} + \frac{b}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}. \end{aligned}$$

Similarly, we get $d_- = 1/\|\mathbf{w}\|$. The margin of this particular separating hyperplane is thus $d = 2/\|\mathbf{w}\|$. If we want to find a hyperplane with *largest* margin, we thus have to solve the quadratic optimization problem

$$\begin{aligned} \text{minimize } & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{subject to } & y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1 \geq 0, \quad 1 \leq i \leq n. \end{aligned}$$

Note that b is also an unknown variable in this problem! The factor $1/2$ in the objective function is just to make the gradient look nicer. The Lagrangian of this problem is

$$\begin{aligned}\mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) &= \frac{1}{2}\|\mathbf{w}\|^2 - \sum_{i=1}^m \lambda_i y_i \mathbf{w}^\top \mathbf{x}_i - \lambda_i y_i b + \lambda_i \\ &= \frac{1}{2}\mathbf{w}^\top \mathbf{w} - \boldsymbol{\lambda}^\top \mathbf{X} \mathbf{w} - b \boldsymbol{\lambda}^\top \mathbf{y} + \sum_{i=1}^m \lambda_i,\end{aligned}$$

where we denote by \mathbf{X} the matrix with the $y_i \mathbf{x}_i^\top$ as rows. We can then write the conditions on the gradient with respect to \mathbf{w} and b of the Lagrangian as

$$\begin{aligned}\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) &= \mathbf{w} - \mathbf{X}^\top \boldsymbol{\lambda} = \mathbf{0} \\ \frac{\partial \mathcal{L}}{\partial b}(\mathbf{w}, b, \boldsymbol{\lambda}) &= \mathbf{y}^\top \boldsymbol{\lambda} = 0.\end{aligned}\tag{18.2}$$

Replacing \mathbf{w} by $\mathbf{X}^\top \boldsymbol{\lambda}$ and $\boldsymbol{\lambda}^\top \mathbf{y}$ by 0 in the Lagrangian function then gives the expression for the Lagrange dual $g(\boldsymbol{\lambda})$,

$$g(\boldsymbol{\lambda}) = -\frac{1}{2} \boldsymbol{\lambda}^\top \mathbf{X} \mathbf{X}^\top \boldsymbol{\lambda} + \sum_{i=1}^m \lambda_i.$$

Finally, changing the sign and the maximum with a minimum, we can formulate the Lagrange dual optimization problem as

$$\text{minimize } \frac{1}{2} \boldsymbol{\lambda}^\top \mathbf{X} \mathbf{X}^\top \boldsymbol{\lambda} - \boldsymbol{\lambda}^\top \mathbf{e} \quad \text{subject to } \boldsymbol{\lambda} \geq \mathbf{0},\tag{18.3}$$

where \mathbf{e} is the vector of all ones.

Note that there is one dual variable λ_i per data point \mathbf{x}_i . We can find the optimal value by solving the dual problem (18.3), but that does not give us automatically the weights \mathbf{w} and the bias b . We can find the weights by $\mathbf{w} = \mathbf{X}^\top \boldsymbol{\lambda}$. As for b , this is best determined from the KKT conditions of the problem. These can be written by combining the constraints of the primal problem with the conditions on the gradient of the Lagrangian (18.2), the condition $\boldsymbol{\lambda} \geq \mathbf{0}$, and complementary slackness as

$$\begin{aligned}\mathbf{X} \mathbf{w} + b \mathbf{y} - \mathbf{e} &\geq \mathbf{0} \\ \boldsymbol{\lambda} &\geq \mathbf{0} \\ \lambda_i(1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) &= 0 \text{ for } 1 \leq i \leq n \\ \mathbf{w} - \mathbf{X}^\top \boldsymbol{\lambda} &= \mathbf{0} \\ \mathbf{y}^\top \boldsymbol{\lambda} &= 0.\end{aligned}$$

To get b , we can choose one of the equations in which $\lambda_i \neq 0$, and then find b by setting $b = y_i(1 - y_i \mathbf{w}^\top \mathbf{x}_i)$. With the KKT conditions written down, we can go about solving the problem of finding a maximum margin linear classifier using methods such as the barrier method.

18.2 Extensions

So far we looked at the particularly simple case where (a) the data falls into two classes, (b) the points can actually be well separated, and (c) they can be separated by an affine hyperplane. In reality, these three assumptions may not hold. We briefly discuss extensions of the basic model to account for the three situations just mentioned.

Non-exact separation

What happens when the data can not be separated by a hyperplane? In this case the constraints can not be satisfied: there is no feasible solution to the problem. We can still modify the problem to allow for *misclassification*: we want to find a hyperplane that separates the two point sets as good as possible, but we allow for some mistakes.

One approach is to add an additional set of n *slack variables* s_1, \dots, s_n , and modify the constraints to

$$\mathbf{w}^\top \mathbf{x}_i + b \geq 1 - s_i, \text{ for } y_i = 1, \quad \mathbf{w}^\top \mathbf{x}_j + b \leq -1 + s_j, \text{ for } y_j = -1, \quad s_i \geq 0.$$

The i -th data point can land on the wrong side of the hyperplane if $s_i > 1$, and consequently the sum $\sum_{i=1}^n s_i$ is an upper bound on the number of errors possible. If we want to minimize the number of misclassified points, we may want to minimize this upper bound, so a sensible choice for objective function would be to add a multiple of this sum. The new problem thus becomes

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + \mu \sum_{j=1}^n s_j \\ \text{subject to} \quad & y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1 + s_i \geq 0, \quad 1 \leq i \leq n \\ & s_i \geq 0, \quad 1 \leq i \leq n, \end{aligned}$$

for some parameter μ . The Lagrangian of this problem and the KKT conditions can be derived in a similar way as in the separable case and are left as an exercise.

Non-linear separation and kernels

The key to extending SVMs from linear to non-linear separation is the observation that the dual form of the optimization problem (18.3) depends only on the dot products $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$ of the data points. In fact, the (i, j) -th entry of the matrix $\mathbf{X} \mathbf{X}^\top$ is precisely $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$!

If we map our data into a higher (possibly infinite) dimensional space \mathcal{H} ,

$$\varphi: \mathbb{R}^p \rightarrow \mathcal{H},$$

and consider the data points $\varphi(\mathbf{x}_i)$, $1 \leq i \leq n$, then applying the support vector machine to these higher dimensional vectors will only depend on the dot products

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \varphi(\mathbf{x}_i), \varphi(\mathbf{x}_j) \rangle.$$

The function K is called a **kernel function**. A typical example, often used in practice, is the Gaussian radial basis function (RBF),

$$K(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|^2/2\sigma^2}.$$

Note that we *don't need to know how the function φ looks like!* In the equation for the hyperplane we simply replace $\mathbf{w}^\top \mathbf{x}$ with $K(\mathbf{w}, \mathbf{x})$. The only difference now is that the function ceases to be linear in \mathbf{x} : we get a non-linear decision boundary.

Multiple classes

One is often interested in classifying data into more than two classes. There are two simple ways in which support vector machines can be extended for such problems: one-vs-one and one-vs-rest. In the one-vs-one case, given k classes, we train one classifier for each pair of classes in the training data, obtaining a total of $k(k - 1)/2$ classifiers. When it comes to prediction, we apply each of the classifiers to our test data and choose the class that was chosen the most among all the classifiers. In the one-vs-rest approach, each train k binary classifiers: in each one, one class corresponds to a chosen class, and the second class corresponds to the rest. By associating confidence scores to each classifier, we choose the one with the highest confidence score.

Example 18.1. An example that uses all three extensions mentioned is handwritten digit recognition. Suppose we have a series of pixels, each representing a number, and associated labels $\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$. We would like to train a support vector machine to recognize new digits. Given the knowledge we have, we can implement this task using standard optimization software such as CVXPY. Luckily, there are packages that have this functionality already implemented, such as the SCIKIT-LEARN package for Python. We illustrate its functioning below. The code also illustrates some standard procedures when tackling a machine learning problem:

- **Separate** the data set randomly into *training data* and *test data*;
- **Create** a support vector classifier with optional parameters;
- **Train** (using `FIT`) the classifier with the training data;
- **Predict** the response using the test data and compare with the true response;
- **Report** the results.

An important aspect to keep in mind is that when testing the performance using the test data, we should compare the classification accuracy to a naive baseline: if, for example, 80% of the test data is classified as +1, then making a prediction of +1 for all the data will give us an accuracy of 80%; in this case, we would want our classifier to perform considerably better than getting the right answer 80% of the time!

```
In [28]: import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
from sklearn import svm, datasets, metrics
from sklearn.model_selection import train_test_split
```

```
In [29]: digits = datasets.load_digits()

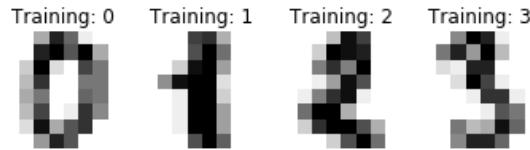
# Display images and labels
images_and_labels = list(zip(digits.images, digits.target))
for index, (image, label) in enumerate(images_and_labels[:4]):
    plt.subplot(2, 4, index + 1)
    plt.axis('off')
    plt.imshow(image, cmap=plt.cm.gray_r, interpolation='nearest')
    plt.title('Training: %i' % label)

# Turn images into 1-D arrays
n_samples = len(digits.images)
data = digits.images.reshape((n_samples, -1))

# Create classifier
svc = svm.SVC(gamma=0.001)

# Randomly split data into train and test set
X_train, X_test, y_train, y_test = train_test_split(data,
                                                    digits.target, test_size = 0.4, random_state=0)
svc.fit(X_train, y_train)
```

```
Out [2]: SVC(C=1.0, cache_size=200, class_weight=None, coef0=0.0,
             decision_function_shape=None, degree=3, gamma=0.001,
             kernel='rbf', max_iter=-1, probability=False, random_state=None,
             shrinking=True, tol=0.001, verbose=False)
```



Now apply prediction to test set and report performance.

```
In [3]: predicted = svc.predict(X_test)
print("Classification report for classifier %s:\n%s\n"
      % (svc, metrics.classification_report(y_test, predicted)))
```

```
Out [3]: Classification report for classifier SVC(C=1.0, cache_size=200,
                                               class_weight=None, coef0=0.0, vdecision_function_shape=None,
                                               degree=3, gamma=0.001, kernel='rbf', max_iter=-1, probability=False,
                                               random_state=None, shrinking=True, tol=0.001, verbose=False):
          precision    recall  f1-score   support

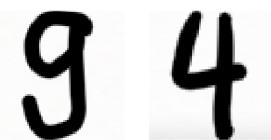
           0       1.00     1.00     1.00      60
           1       0.97     1.00     0.99      73
           2       1.00     0.97     0.99      71
           3       1.00     1.00     1.00      70
           4       1.00     1.00     1.00      63
           5       1.00     0.98     0.99      89
           6       0.99     1.00     0.99      76
           7       0.98     1.00     0.99      65
           8       1.00     0.99     0.99      78
           9       0.99     1.00     0.99      74

      avg / total       0.99     0.99     0.99      719
```

```
In [4]: import skimage
from skimage import data
from skimage.transform import resize
from skimage import io
import os
```

Now try this out on some original data!

```
In [5]: mydigit1 = io.imread('images/digit9.png')
mydigit2 = io.imread('images/digit4.png')
plt.figure(figsize=(8, 4))
plt.subplot(1,2,1)
plt.imshow(mydigit1, cmap=plt.cm.gray_r, interpolation='nearest')
plt.axis('off')
plt.subplot(1,2,2)
plt.imshow(mydigit2, cmap=plt.cm.gray_r, interpolation='nearest')
plt.axis('off')
plt.show()
```



```
In [6]: smalldigit1 = resize(mymain1, (8,8))
smalldigit2 = resize(mymain2, (8,8))
mydigits = np.concatenate((np.round(15*(np.ones((8,8))-
                                         smalldigit1[:, :, 0])).reshape((64,1)).T,
                           np.round(15*(np.ones((8,8))-
                                         smalldigit2[:, :, 0])).reshape((64,1)).T), axis=0)
# After some preprocessing, make prediction
guess = svc.predict(mydigits)
print guess
```

[9 4]

Part V

Semidefinite Programming

Lecture 19

Semidefinite programming (SDP) is a far-reaching generalization of linear programming (LP), in which vectors are replaced by matrices and the non-negativity condition is replaced with the condition that the matrices be positive semidefinite.

19.1 Semidefinite programming

Recall the linear programming formulation

$$\begin{aligned} & \text{minimize} && \langle \mathbf{c}, \mathbf{x} \rangle \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned} \tag{LP}$$

with $\mathbf{c}, \mathbf{x} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{m \times n}$, and $\mathbf{b} \in \mathbb{R}^m$. Semidefinite programming is a far-reaching generalization of linear programming, in which the vector space \mathbb{R}^n is replaced by the vector space of real symmetric matrices,

$$\text{SYM}_n = \{\mathbf{X} \in \mathbb{R}^{n \times n} : x_{ij} = x_{ji}, 1 \leq i < j \leq n\},$$

with the trace inner product,

$$\langle \mathbf{X}, \mathbf{Y} \rangle = \mathbf{X} \bullet \mathbf{Y} = \text{tr}(\mathbf{XY}) = \sum_{i=1}^n \sum_{j=1}^n x_{ij}y_{ij},$$

where $\text{tr}(\mathbf{X}) = \sum_{i=1}^n x_{ii}$ is the *trace* of a matrix \mathbf{X} . The matrix \mathbf{A} in (LP) is replaced by a linear map $\mathcal{A}: \text{SYM}_n \rightarrow \mathbb{R}^m$. So far, this new setting does not fall out of the general framework of linear programming (note that (LP) does not make any reference on the nature of the inner product and the linear map \mathbf{A}). The point of departure of semidefinite programming is that the inequality constraints $\mathbf{x} \geq \mathbf{0}$ are replaced with the constraint $\mathbf{X} \succeq \mathbf{0}$, meaning that \mathbf{X} is positive semidefinite. A *semidefinite programming* (SDP) problem thus has the form

$$\begin{aligned} & \text{minimize} && \mathbf{C} \bullet \mathbf{X} \\ & \text{subject to} && \mathcal{A}(\mathbf{X}) = \mathbf{b} \\ & && \mathbf{X} \succeq \mathbf{0}. \end{aligned} \tag{SDP-P}$$

Remark 19.1. Recall that $\mathbf{X} \succeq \mathbf{0}$ means that for all $\mathbf{v} \in \mathbb{R}^n$, $\mathbf{v}^\top \mathbf{X} \mathbf{v} \geq 0$. Therefore, SDP can be viewed as linear programming with an uncountable number of inequality constraints. Recall also that a symmetric matrix is positive semidefinite if and only if the (necessarily real) eigenvalues are non-negative. While the feasible sets of linear programming are polyhedra, the feasible sets of SDP are called *spectrahedra*, since they are “polyhedra in the spectrum (set of eigenvalues)”.

The vector space SYM_n has dimension $d := n(n + 1)/2$, the number of entries on or above the main diagonal. Therefore, if we identify $\text{SYM}_n \cong \mathbb{R}^d$, the linear map \mathcal{A} can be identified with a “big” $m \times d$ matrix. Alternatively, the condition $\mathcal{A}(\mathbf{X}) = \mathbf{b}$ can be seen as collection of m linear constraints, each of the form $\mathbf{A}_i \bullet \mathbf{X} = b_i$ with a symmetric matrix \mathbf{A}_i , so that we can express (SDP-P) as

$$\begin{aligned} & \text{minimize} && \mathbf{C} \bullet \mathbf{X} \\ & \text{subject to} && \mathbf{A}_i \bullet \mathbf{X} = b_i, \quad 1 \leq i \leq m, \\ & && \mathbf{X} \succeq \mathbf{0}. \end{aligned}$$

Example 19.2. Consider an SDP with $n = 3$ and $m = 2$ (3×3 symmetric matrices and 2 constraints), with the following data

$$\mathbf{A}_1 = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 3 & 7 \\ 1 & 7 & 5 \end{pmatrix}, \quad \mathbf{A}_2 = \begin{pmatrix} 0 & 2 & 8 \\ 2 & 6 & 0 \\ 8 & 0 & 4 \end{pmatrix}, \quad \mathbf{C} = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 9 & 0 \\ 3 & 0 & 7 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 11 \\ 19 \end{pmatrix}.$$

The decision variable is the symmetric matrix

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & x_{13} \\ x_{12} & x_{22} & x_{23} \\ x_{13} & x_{23} & x_{33} \end{pmatrix}.$$

The objective function is

$$\mathbf{C} \bullet \mathbf{X} = x_{11} + 4x_{12} + 6x_{13} + 9x_{22} + 7x_{33},$$

where we used the symmetry $x_{ij} = x_{ji}$ to simplify the expression. Working out the constraints $\mathbf{A}_i \bullet \mathbf{X} = b_i$ in the same way, we arrive at the following form of the SDP:

$$\begin{aligned} & \text{minimize} && x_{11} + 4x_{12} + 6x_{13} + 9x_{22} + 7x_{33} \\ & \text{subject to} && x_{11} + 2x_{13} + 3x_{22} + 14x_{23} + 5x_{33} = 11 \\ & && 4x_{12} + 16x_{13} + 6x_{22} + 4x_{33} = 19 \\ & && \mathbf{X} \succeq \mathbf{0}. \end{aligned}$$

As mentioned in Remark 19.1, were it not for the positive semidefinite constraint, the above would be an old-fashioned LP in the variables $\mathbf{x} = (x_{11}, x_{12}, x_{13}, x_{22}, x_{23}, x_{33})^\top$.

Linear programming is a special case of semidefinite programming. To see this, start with the problem (LP) and create the new matrices

$$\mathbf{A}_i = \text{diag}(a_{i1}, \dots, a_{in}), \quad \mathbf{C} = \text{diag}(c_1, \dots, c_n).$$

Then the problem (LP) is equivalent to the SDP

$$\begin{aligned} & \text{minimize} && \mathbf{C} \bullet \mathbf{X} \\ & \text{subject to} && \mathbf{A}_i \bullet \mathbf{X} = b_i, \quad 1 \leq i \leq m, \\ & && x_{ij} = 0, \quad i < j, \\ & && \mathbf{X} \succeq \mathbf{0}. \end{aligned}$$

In fact, any solution of the above has the form $\mathbf{X} = \text{diag}(x_{11}, \dots, x_{nn})$, which we identify with a vector $\mathbf{x} \in \mathbb{R}^n$, with the semidefinite constraint, applied to the diagonal matrix, translating into $x_i \geq 0$ for $1 \leq i \leq n$.

19.2 Semidefinite programming duality

Just as with linear programming, we can formulate a dual problem

$$\begin{aligned} & \text{maximize} && \langle \mathbf{y}, \mathbf{b} \rangle \\ & \text{subject to} && \sum_{i=1}^m y_i \mathbf{A}_i + \mathbf{S} = \mathbf{C} \\ & && \mathbf{S} \succeq \mathbf{0}. \end{aligned} \tag{SDP-D}$$

Example 19.3. The dual to the problem in Example 19.2 is the problem

$$\begin{aligned} & \text{maximize} && 11y_1 + 19y_2 \\ & \text{subject to} && y_1 \begin{pmatrix} 1 & 0 & 1 \\ 0 & 3 & 7 \\ 1 & 7 & 5 \end{pmatrix} + y_2 \begin{pmatrix} 0 & 2 & 8 \\ 2 & 6 & 0 \\ 8 & 0 & 4 \end{pmatrix} + \mathbf{S} = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 9 & 0 \\ 3 & 0 & 7 \end{pmatrix}, \\ & && \mathbf{S} \succeq \mathbf{0}. \end{aligned}$$

One can also write the constraints more compactly by eliminating \mathbf{S} , in terms of a matrix whose entries are linear functions of y_1 and y_2 and which is required to be positive semidefinite.

We call \mathbf{X} a feasible matrix for (SDP-P) if it satisfies the constraints, and (\mathbf{y}, \mathbf{S}) a feasible pair for (SDP-D) if \mathbf{y} and \mathbf{S} satisfy the constraints. As with the nonlinear optimization setting from Lecture 16, we call the *duality gap* the difference $\mathbf{C} \bullet \mathbf{X} - \langle \mathbf{b}, \mathbf{y} \rangle$ between the primal and dual objective functions at feasible points.

Theorem 19.4. *Let \mathbf{X} be a feasible matrix for (SDP-P) and let (\mathbf{y}, \mathbf{S}) be a feasible pair for (SDP-D). Then the duality gap satisfies*

$$\mathbf{C} \bullet \mathbf{X} - \langle \mathbf{b}, \mathbf{y} \rangle = \mathbf{S} \bullet \mathbf{X} \geq 0.$$

Moreover, if the duality gap is zero, then \mathbf{X} and (\mathbf{y}, \mathbf{S}) are optimal points for (SDP-P) and (SDP-D), respectively.

Recall that $\mathbf{X} \in \text{SYM}_n$ if and only if $\mathbf{X} = \mathbf{Q} \mathbf{D} \mathbf{Q}^\top$, with \mathbf{D} diagonal (containing the eigenvalues) and \mathbf{Q} orthogonal (that is, $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$). If in addition $\mathbf{X} \succeq \mathbf{0}$, then the diagonal entries are nonnegative. We need some more facts about positive semidefinite matrices.

Lemma 19.5. Let $\mathbf{X} \in \mathcal{S}_+^n$. Then:

1. The diagonal entries satisfy $x_{ii} \geq 0$ for $1 \leq i \leq n$;
2. If $x_{ii} = 0$ for some i , then $x_{ij} = x_{ji} = 0$ for all $1 \leq j \leq n$ (if a diagonal is zero, then the whole corresponding row and column is zero).

Proof of Theorem 19.4. (Optional.) Note that

$$\mathbf{C} \bullet \mathbf{X} - \langle \mathbf{b}, \mathbf{y} \rangle = \sum_{i=1}^m y_i \mathbf{A}_i \bullet \mathbf{X} + \mathbf{S} \bullet \mathbf{X} - \langle \mathbf{b}, \mathbf{y} \rangle = \mathbf{S} \bullet \mathbf{X},$$

since $\sum_{i=1}^m y_i \mathbf{A}_i \bullet \mathbf{X} = \sum_{i=1}^m y_i b_i = \langle \mathbf{y}, \mathbf{b} \rangle$. We next show that if $\mathbf{X} \succeq 0$ and $\mathbf{S} \succeq 0$, then $\mathbf{X} \bullet \mathbf{S} \geq 0$. Assume $\mathbf{X} = \mathbf{Q} \mathbf{D} \mathbf{Q}^\top$ and $\mathbf{S} = \mathbf{P} \mathbf{E} \mathbf{P}^\top$, with \mathbf{Q}, \mathbf{P} orthogonal and \mathbf{D}, \mathbf{E} diagonal. Then

$$\mathbf{X} \bullet \mathbf{S} = \text{tr}(\mathbf{S} \mathbf{X}) = \text{tr}(\mathbf{Q} \mathbf{D} \mathbf{Q}^\top \mathbf{P} \mathbf{E} \mathbf{P}^\top) = \text{tr}(\mathbf{D} \mathbf{Q}^\top \mathbf{P} \mathbf{E} \mathbf{P}^\top \mathbf{Q}),$$

where we used that $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$. The last expression equals

$$\sum_{i=1}^n \mathbf{D}_{ii} (\mathbf{Q}^\top \mathbf{P} \mathbf{E} \mathbf{P}^\top \mathbf{Q})_{ii} \geq 0, \quad (19.1)$$

since $\mathbf{D}_{ii} \geq 0$ and $\mathbf{Q}^\top \mathbf{P} \mathbf{E} \mathbf{P}^\top \mathbf{Q} \succeq 0$, which implies that the diagonal entries are also nonnegative by Lemma 19.5.

Now if $\mathbf{X} \bullet \mathbf{S} = 0$, then by (19.1), $\sum_{i=1}^n \mathbf{D}_{ii} (\mathbf{Q}^\top \mathbf{P} \mathbf{E} \mathbf{P}^\top \mathbf{Q})_{ii} = 0$. This means that for each i , either $\mathbf{D}_{ii} = 0$, or $(\mathbf{Q}^\top \mathbf{P} \mathbf{E} \mathbf{P}^\top \mathbf{Q})_{ii} = 0$. In the latter case, by Lemma 19.5, the whole j -th row of this matrix is zero. It follows that

$$\mathbf{D}(\mathbf{Q}^\top \mathbf{P} \mathbf{E} \mathbf{P}^\top \mathbf{Q}) = (\mathbf{D} \mathbf{Q}^\top)(\mathbf{P} \mathbf{E} \mathbf{P}^\top \mathbf{Q}) = \mathbf{0}.$$

Since $\mathbf{AB} = \mathbf{0}$ implies $\mathbf{BA} = \mathbf{0}$, we get that $\mathbf{XS} = (\mathbf{P} \mathbf{E} \mathbf{P}^\top \mathbf{Q})(\mathbf{D} \mathbf{Q}^\top) = \mathbf{0}$. \square

As with nonlinear convex optimization, some mild conditions (also known as Slater's condition) ensure that the primal and dual optimal values coincide.

Theorem 19.6. Let p^* be the optimal value of (SDP-P) and d^* the optimal value of (SDP-D). If there exists a feasible matrix $\mathbf{X} \succ 0$ for (SDP-P), and a feasible pair (\mathbf{y}, \mathbf{S}) for (SDP-D) with $\mathbf{S} \succ 0$, then $p^* = d^*$.

The theory of interior point methods carries over to semidefinite programming, giving us efficient (in theory and practice) algorithms for solving such problem. In the next lecture we will address some concrete problems that can be solved efficiently using semidefinite programming, and for which linear or quadratic programming are not enough.

Lecture 20

The standard form of a semidefinite programming problem is

$$\begin{aligned} & \text{minimize} && \mathbf{C} \bullet \mathbf{X} \\ & \text{subject to} && \mathbf{A}_i \bullet \mathbf{X} = b_i, \quad 1 \leq i \leq m, \\ & && \mathbf{X} \succeq 0, \end{aligned} \tag{SDP-P}$$

where \mathbf{C} and \mathbf{A}_i , $1 \leq i \leq m$, are symmetric $n \times n$ matrices (elements of the vector space SYM_n). At first sight, it is not clear why solving such problem should be of interest. In this lecture we discuss a useful application to *discrete* optimization problems that are usually difficult in practice.

20.1 Semidefinite relaxation

A graph is a pair $G = (V, E)$, where $V = \{1, \dots, n\}$ consists of the *vertices*, and $E \subseteq V \times V$ consists of *edges* connecting some of the vertices. We will consider undirected edges, i.e., (i, j) will be the same edge as (j, i) .

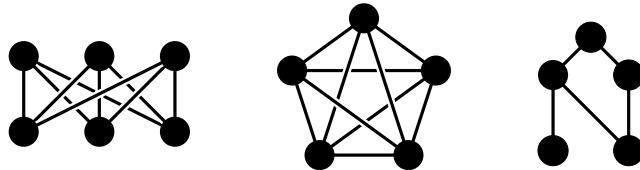


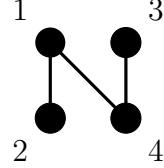
Figure 20.24: Some graphs

Note that for graphs, we only care about the incidences (which node is connected to which other node), the images are just visualizations and the distances and placements of nodes are pure convenience. Graphs are an important tool in areas such as network analysis, and many discrete computational problems can be formulated as graph problems. One such problem is MAXCUT. Given a graph $G = (V, E)$, a *cut* is a pair $(S, V \setminus S)$, where $S \subseteq V$ is a subset of vertices. The *edge set* of the cut is the set

$$E(S, V \setminus S) = \{e \in E : e \text{ connects a vertex in } S \text{ with a vertex in } V \setminus S\}.$$

The size of a cut is the number of edges in it, and the MAXCUT problem is the problem of finding a cut of maximal size.

Example 20.1. Consider the graph with $V = \{1, 2, 3, 4\}$ and $E = \{(1, 2), (3, 4), (1, 4)\}$. There



are 7 nontrivial cuts, namely

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

Note that we don't list all the non-empty subsets of $\{1, 2, 3, 4\}$, as that would amount to counting every cut twice (every set would feature as S and as $V \setminus S$). The sizes of these cuts (number of edges joining S with $V \setminus S$) is

$$2, 1, 1, 2, 1, 3, 2.$$

The maximum cut size is therefore 3 (it can't be any bigger, since we only have three edges), and one such cut can be realized by taking the set $S = \{1, 3\}$ and $V \setminus S = \{2, 4\}$.

The problem MAXCUT (like many other graph problems) has the property that it is NP-hard. This means that the decision version of this problem (given a cut, does it have size at least k ?) is NP-complete, which in turns means that checking the statement is easy (just count the number of edges in the cut), but *finding* a cut of size at least k appears to be difficult: there does not seem to be a way of doing this that is fundamentally faster than to test all the possible cuts. A famous conjecture, $P \neq NP$, implies that there does not exist an efficient (polynomial time) algorithm that would, in general, be able to find a cut of a given size, or even a maximal cut.

Two common approaches to deal with difficult problems are *approximation* and *randomization*. In approximation, one gives up on finding the best solution and concentrates on finding one that is good enough. In randomization, one takes into account that an algorithm may fail with small probability. Both approximation and randomization may lead to efficient algorithms that can solve a problem well enough for “all practical purposes”.

In the case of MAXCUT, an approximation algorithm \mathcal{A} takes as input a graph G and outputs a set $S \subseteq V$, written as $\mathcal{A}(G) = S$. If by $\omega(S)$ we denote the size of the cut induced by S , and by $\text{Opt}(G)$ the largest size of a cut, then the algorithm \mathcal{A} is said to provide an δ -approximation (for some $\delta > 0$), if

$$\omega(\mathcal{A}(G)) \geq \delta \cdot \text{Opt}(G).$$

In words, the algorithms provides a cut that is optimal up to a factor δ . A randomized algorithm is one that is allowed to use internal coin flips. Such an algorithm is a δ -approximation algorithm, if the expected value

$$\mathbb{E}[\omega(\mathcal{A}(G))] \geq \delta \cdot \text{Opt}(G).$$

We are interested in approximation algorithms that run in *polynomial time*. This means that the number of computational steps is a polynomial in the input size. Interior point methods for convex optimization and semidefinite programming are examples of polynomial time algorithms, and it is a remarkable fact that *MaxCut* can be solved approximatively using semidefinite programming.

20.2 The Goemans-Williamson Algorithm

The MAXCUT problem can be formulated as an integer programming problem as follows. Introduce variables x_1, \dots, x_n , where n is the number of vertices. A cut is then an assignment $x_i = 1$ (meaning that $i \in S$) or $x_i = -1$ (meaning that $i \in V \setminus S$). If $e = (i, j)$ is an edge in the cut, then $x_i x_j = -1$, and $x_i x_j = 1$ otherwise. We therefore have

$$\frac{1}{2}(1 - x_i x_j) = \begin{cases} 1 & \text{if } (i, j) \in E(S, V \setminus S) \\ 0 & \text{if } (i, j) \notin E(S, V \setminus S). \end{cases}$$

The size of a cut is then the sum $\sum_{(i,j) \in E} \frac{1}{2}(1 - x_i x_j)$, and the problem of finding the maximum cut can be written as

$$\begin{aligned} & \text{maximize} && \sum_{(i,j) \in E} \frac{1}{2}(1 - x_i x_j) \\ & \text{subject to} && x_i^2 = 1, \quad 1 \leq i \leq n. \end{aligned} \tag{20.1}$$

We now look for a *relaxation* of the problem: a new problem with a bigger constraint set that also contains the solutions of (20.1), but possibly more. As a first step, replace $x_i \in \{-1, 1\}$ with vectors $\mathbf{u}_i \in \{\pm e_1\}$, where $e_1 = (1, 0, \dots, 0)^\top$ is a unit vector, and $x_i x_j$ with $\mathbf{u}_i^\top \mathbf{u}_j$. The optimization is then over sets of vectors $(\mathbf{u}_1, \dots, \mathbf{u}_n)$, and can be written as

$$\begin{aligned} & \text{maximize} && \sum_{(i,j) \in E} \frac{1}{2}(1 - \mathbf{u}_i^\top \mathbf{u}_j) \\ & \text{subject to} && \mathbf{u}_i \in \{\pm e_1\}. \end{aligned}$$

The relaxation now consists in enlarging the constraint set to allow for *any* unit vectors $\mathbf{u}_i \in S^{n-1}$, where $S^{n-1} = \{\mathbf{x} : \|\mathbf{x}\| = 1\}$ is the unit sphere.

$$\begin{aligned} & \text{maximize} && \sum_{(i,j) \in E} \frac{1}{2}(1 - \mathbf{u}_i^\top \mathbf{u}_j) \\ & \text{subject to} && \|\mathbf{u}_i\| = 1. \end{aligned} \tag{20.2}$$

As the constraint set of the relaxation is bigger, the solution will be at least as large as the solution of (20.1) (ideally, the same, but we can't guarantee this).

For any set of unit vectors \mathbf{u}_i , let $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_n)$ be the matrix with unit vectors \mathbf{u}_i as columns, and consider

$$\mathbf{X} = \mathbf{U}^\top \mathbf{U}.$$

Then $x_{ii} = \mathbf{u}_i^\top \mathbf{u}_i = 1$ and the matrix \mathbf{X} is symmetric and positive definite. Conversely, any symmetric positive definite matrix \mathbf{X} can be written as $\mathbf{X} = \mathbf{U}^\top \mathbf{U}$ (Cholesky factorization), and if in addition $x_{ii} = 1$, then the columns of \mathbf{U} necessarily have unit length. We conclude that Problem (20.2) is equivalent to the following SDP:

$$\begin{aligned} & \text{maximize} && \sum_{(i,j) \in E} \frac{1}{2}(1 - x_{ij}) \\ & \text{subject to} && x_{ii} = 1 \\ & && \mathbf{X} \succeq 0. \end{aligned} \tag{20.3}$$

Write $\text{SDP}(G)$ for the optimal value of (20.3). From the above discussion we have

$$\text{SDP}(G) \geq \text{Opt}(G).$$

This is the case, because the solution of (20.1) is contained in the feasible set of (20.3), but the latter has more options.

To recover the cut, we proceed as follows.

1. Solve (20.3) to accuracy ε to obtain a matrix \mathbf{X}^* such that

$$\sum_{(i,j) \in E} \frac{1}{2}(1 - x_{ij}^*) \geq \text{SDP}(G) - \varepsilon.$$

This can be done in polynomial time using, for example, interior point methods.

2. Perform a Cholesky factorization $\mathbf{X}^* = (\mathbf{U}^*)^\top \mathbf{U}^*$ and extract the columns $\mathbf{u}_1^*, \dots, \mathbf{u}_n^*$, which are unit vectors and satisfy

$$\sum_{(i,j) \in E} \frac{1}{2}(1 - (\mathbf{u}_i^*)^\top \mathbf{u}_j^*) \geq \text{Opt}(G) - \varepsilon.$$

3. Finally, we want a way to recover a partition of the graph without loosing too much. We thus need to assign to each \mathbf{u}_i a value $x_i \in \{-1, 1\}$ according so some rule, and declare the result to be our partition. To do so, we *randomly* choose a vector $\mathbf{p} \in S^{n-1}$ and then set

$$x_i = \begin{cases} 1 & \text{if } \mathbf{u}_i^\top \mathbf{p} \geq 0, \\ -1 & \text{else.} \end{cases}$$

There we have our algorithm: it generates a cut. Let's call the above algorithm, which takes a graph G and produces a set S by the above steps, \mathcal{A} .

Theorem 20.2. *The Goemans-Williamson algorithms generates a cut $\mathcal{A}(G) = S$ and satisfies*

$$\mathbb{E}[\omega(\mathcal{A}(G))] \geq 0.8785672 \cdot \text{Opt}(G).$$

That is, the algorithm is expected to generate a cut that is at least as 0.8785672 times as good as the optimal cut (that is, could be about 10 percent smaller but not more). By the law of large numbers, repeating the algorithm will likely give a result close to the expectation.

Appendix

This appendix contain a summary of background material from linear algebra and calculus. Much of the content should be familiar to some degree, and the purpose is to bring it back to attention. Important concepts are **highlighted** in the notes.

1 Asymptotic notation

An **algorithm** is a sequence of instructions carried out by a computer. Important features of an algorithm are computation time (measured as the number of operations or the number of iterations needed to reach a solution) and accuracy. One is mainly interested in the **orders of magnitude** of these quantities, and not so much in their exact values. A convenient notation for this purpose is the asymptotic O-notation.

Let $f, g: \mathbb{R} \rightarrow \mathbb{R}$ be two functions. Then:

- $f(n) \in O(g(n))$ as $n \rightarrow \infty$ if there exists a constant $C > 0$ and an integer n_0 such that $|f(n)| \leq C|g(n)|$ for $n > n_0$;
- $f(x) \in O(g(x))$ as $x \rightarrow 0$ if there exists a constant $C > 0$ and a real number $\varepsilon > 0$ such that $|f(x)| \leq C|g(x)|$ for $|x| < \varepsilon$.

We omit “ $n \rightarrow \infty$ ” or “ $x \rightarrow 0$ ” when it is clear from the context. One often finds statements such as $f(n) = O(g(n))$ or $f(x) = 1 + x + O(x^2)$; the first is equivalent to $f(n) \in O(g(n))$, while the second should be read as $f(x) = 1 + x + g(x)$ for a function $g(x) \in O(x^2)$. The following examples illustrate the O-notation.

- $\sqrt{n} + n^2 \in O(n^2)$ as $n \rightarrow \infty$
- $n^5 \in O(e^n)$ as $n \rightarrow \infty$
- $10^{100} \in O(1)$ as $n \rightarrow \infty$
- $x^3 \in O(x^2)$ as $x \rightarrow 0$
- $\sin(x) \in O(x)$ as $x \rightarrow 0$
- $e^x = 1 + x + O(x^2)$ as $x \rightarrow 0$

The notation $f(n) \in \Omega(g(n))$ as $n \rightarrow \infty$ means that $g(n) \in O(f(n))$ as $n \rightarrow \infty$, and $f(n) \in o(g(n))$ as $n \rightarrow \infty$ means that for all $C > 0$ there exists n_0 such that $|f(n)| < C|g(n)|$ for $n > n_0$. If $g(n) \neq 0$ for sufficiently large n , this is equivalent to $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0$. One defines $\Omega(g(x))$ and $o(g(x))$ as $x \rightarrow 0$ analogously.

2 Linear Algebra

We restrict to linear algebra over the field of real numbers \mathbb{R} , as this is the setting that is of most interest in optimization. A **vector** in \mathbb{R}^n and its **transpose** are written as

$$\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = (x_1, \dots, x_n)^\top, \quad \mathbf{x}^\top = (x_1, \dots, x_n),$$

with **coordinates** $x_i \in \mathbb{R}$ for $1 \leq i \leq n$. The zero vector is denoted by $\mathbf{0}$, while \mathbf{e} is the vector with every coordinate equal to 1. If $\lambda_1, \lambda_2 \in \mathbb{R}$ and $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, then $\lambda_1\mathbf{x} + \lambda_2\mathbf{y}$ is the vector with coordinates $\lambda_1x_i + \lambda_2y_i$ for $1 \leq i \leq n$.

In \mathbb{R}^n we have the Euclidean (or standard) **inner product** (or scalar product)

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^n x_i y_i.$$

The Euclidean inner product is **bilinear**: for $\mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2 \in \mathbb{R}^n$ and $\alpha, \beta \in \mathbb{R}$,

$$\langle \alpha\mathbf{x}_1 + \beta\mathbf{x}_2, \mathbf{y} \rangle = \alpha\langle \mathbf{x}_1, \mathbf{y} \rangle + \beta\langle \mathbf{x}_2, \mathbf{y} \rangle, \quad \langle \mathbf{x}, \alpha\mathbf{y}_1 + \beta\mathbf{y}_2 \rangle = \alpha\langle \mathbf{x}, \mathbf{y}_1 \rangle + \beta\langle \mathbf{x}, \mathbf{y}_2 \rangle,$$

symmetric ($\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$) and satisfies $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$, with equality if and only if $\mathbf{x} = \mathbf{0}$. Two vectors \mathbf{x} and \mathbf{y} are called **orthogonal**, if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$.

Example 1.3. The vectors $(1, 1)^\top$ and $(1, -1)^\top$ are orthogonal in \mathbb{R}^2 , while $(1, 1)^\top$ and $(2, -1)^\top$ are not.

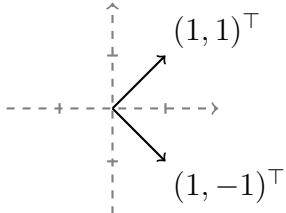


Figure 2.25: Orthogonal vectors

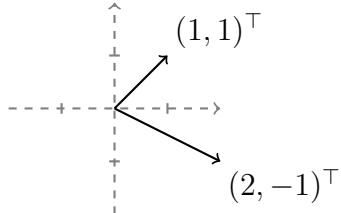


Figure 2.26: Non-orthogonal vectors

Linear subspaces

A **linear subspace** is a subset $V \subseteq \mathbb{R}^n$ such that for any $\mathbf{x}, \mathbf{y} \in V$ and for all $\alpha, \beta \in \mathbb{R}$, $\alpha\mathbf{x} + \beta\mathbf{y} \in V$. In particular, the sets $\{\mathbf{0}\}$ and \mathbb{R}^n are linear subspaces.

Example 1.4. The linear subspaces of \mathbb{R}^2 are $\{\mathbf{0}\}$, lines through the origin, and \mathbb{R}^2 . The linear subspaces of \mathbb{R}^3 are $\{\mathbf{0}\}$, lines and planes through the origin, and \mathbb{R}^3 .

A **linear combination** of vectors $\mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^n$ is an expression of the form $\mathbf{x} = \sum_{i=1}^k \lambda_i \mathbf{x}_i$, where $\lambda_i \in \mathbb{R}$ for $1 \leq i \leq k$. The set of linear combinations

$$V = \text{span} \{ \mathbf{x}_1, \dots, \mathbf{x}_k \} := \left\{ \sum_{i=1}^k \lambda_i \mathbf{x}_i : \lambda_i \in \mathbb{R} \right\}$$

forms a linear subspace of \mathbb{R}^n . It is the intersection of all linear subspaces that contain $\mathbf{x}_1, \dots, \mathbf{x}_k$. The vectors $\mathbf{x}_1, \dots, \mathbf{x}_k$ are **linearly independent** if $\sum_{i=1}^k \lambda_i \mathbf{x}_i = 0$ implies $\lambda_1 = \dots = \lambda_k = 0$. A minimal set of vectors that span a linear subspace V is called a **basis** of this subspace, and the number of elements in a basis is the **dimension** of the linear subspace. The elements of a basis are always linearly independent, and a maximal linearly independent set in a vector subspace V is a basis. If $\{\mathbf{b}_1, \dots, \mathbf{b}_k\}$ is a basis of a subspace V , then every $\mathbf{x} \in V$ has a *unique* representation $\mathbf{x} = \sum_{i=1}^k \lambda_i \mathbf{b}_i$. A basis is **orthogonal** if $\langle \mathbf{b}_i, \mathbf{b}_j \rangle = 0$ for $i \neq j$, and **orthonormal** if in addition $\langle \mathbf{b}_i, \mathbf{b}_i \rangle = 1$ for $1 \leq i \leq k$. The unique expression of $\mathbf{x} \in V$ as linear combination of an orthonormal basis $\{\mathbf{b}_1, \dots, \mathbf{b}_k\}$ of V is given by

$$\mathbf{x} = \sum_{i=1}^k \langle \mathbf{x}, \mathbf{b}_i \rangle \mathbf{b}_i.$$

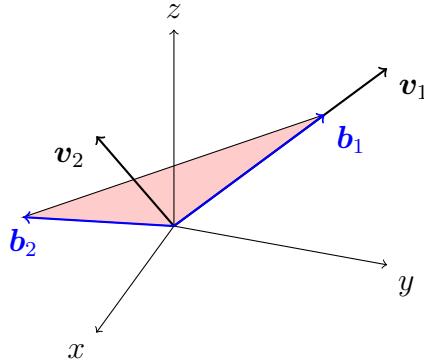
The **standard basis** of \mathbb{R}^n is the orthonormal basis $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$, where \mathbf{e}_i has a 1 in the i -th coordinate and 0 elsewhere.

Example 1.5. The vectors $\mathbf{v}_1 = (0, 1, 1)^\top$ and $\mathbf{v}_2 = (1, 0, 1)^\top$ span a linear subspace of \mathbb{R}^3 of dimension 2, but they are not orthogonal. The vectors

$$\mathbf{b}_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}, \mathbf{b}_2 = \frac{1}{\sqrt{3}} \begin{pmatrix} \sqrt{2} \\ -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix}$$

form an orthonormal basis of V . The vector $\mathbf{x} = (1, 1, 2)^\top$ lives in V , and its representation in terms of $\{\mathbf{b}_1, \mathbf{b}_2\}$ is

$$\mathbf{x} = \frac{3}{\sqrt{2}} \mathbf{b}_1 + \sqrt{\frac{3}{2}} \mathbf{b}_2.$$



The **direct sum** of two vector subspaces $V, W \subset \mathbb{R}^n$ with $V \cap W = \{0\}$ is

$$V \oplus W = \{\mathbf{v} + \mathbf{w} : \mathbf{v} \in V, \mathbf{w} \in W\}.$$

The **orthogonal complement** of a subspace $V \subseteq \mathbb{R}^n$ is the set

$$V^\perp = \{\mathbf{x} \in \mathbb{R}^n : \forall \mathbf{y} \in V : \langle \mathbf{x}, \mathbf{y} \rangle = 0\}.$$

The vector space \mathbb{R}^n is the direct sum of V and its orthogonal complement,

$$\mathbb{R}^n = V \oplus V^\perp. \quad (1)$$

If $V = \text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$, then $V^\perp = \{\mathbf{y} : \langle \mathbf{x}_1, \mathbf{y} \rangle = \dots = \langle \mathbf{x}_k, \mathbf{y} \rangle = 0\}$; to check whether a vector \mathbf{y} is in the orthogonal complement of V we therefore only need to check whether \mathbf{y} is orthogonal to a spanning set (for example, a basis) of V .

Example 1.6. The vector $\mathbf{x} = (-1, -1, 1)^\top$ is orthogonal to the basis $\{\mathbf{b}_1, \mathbf{b}_2\}$ from Example 1.5. It is therefore orthogonal to the whole plane $V = \text{span}\{\mathbf{b}_1, \mathbf{b}_2\}$ spanned by these vectors. The orthogonal complement of V is the line $\{\lambda \mathbf{x} : \lambda \in \mathbb{R}\}$.

The **direct product** of two vector subspaces $V \subseteq \mathbb{R}^n$, $W \subseteq \mathbb{R}^m$ is defined as

$$V \times W = \{(\mathbf{v}, \mathbf{w}) \in \mathbb{R}^{n+m} : \mathbf{v} \in V, \mathbf{w} \in W\}.$$

where (\mathbf{v}, \mathbf{w}) is the vector whose first n coordinates coincide with \mathbf{v} , and the last m coordinates coincide with \mathbf{w} . In particular, $\mathbb{R}^n \times \mathbb{R}^m = \mathbb{R}^{n+m}$.

Linear maps

An $m \times n$ matrix

$$\mathbf{A} = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix},$$

represents a **linear map** from \mathbb{R}^n to \mathbb{R}^m by means of

$$\mathbf{y} = \mathbf{Ax}, \quad y_i = \sum_{j=1}^n a_{ij}x_j.$$

For example,

$$\begin{pmatrix} 2 & 1 & 0 \\ 1 & 0 & 2 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix} = \begin{pmatrix} 3 \\ 5 \end{pmatrix}.$$

The columns of a matrix are vectors, and we sometimes write

$$\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_n)$$

for the matrix whose columns are given by the vectors \mathbf{a}_i . If $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $n = m + k$, then \mathbf{A} can be written as **block matrix**,

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix},$$

with $\mathbf{A}_{11} \in \mathbb{R}^{m \times m}$, $\mathbf{A}_{22} \in \mathbb{R}^{k \times k}$, $\mathbf{A}_{12} \in \mathbb{R}^{m \times k}$ and $\mathbf{A}_{21} \in \mathbb{R}^{k \times m}$. The sum and difference of matrices of the same size are defined component-wise.

The $n \times n$ matrix **1** is the matrix with 1 on the diagonal and 0 elsewhere, while **0** is the matrix consisting of only zeros. A matrix is **diagonal** if all the off-diagonal elements are 0,

lower-triangular if all the elements above the diagonal are 0, and **upper-triangular** if all the elements below the diagonal are 0. A **block-diagonal** matrix is a block matrix, with all blocks outside the main diagonal consisting of zero-matrices $\mathbf{0}$.

The **transpose** \mathbf{A}^\top is the matrix with entries $a'_{ij} := a_{ji}$. It is the matrix \mathbf{A} mirrored on the diagonal from top left to bottom right. A matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is called **symmetric** if $\mathbf{A}^\top = \mathbf{A}$. The set of symmetric matrices in $\mathbb{R}^{n \times n}$ is denoted by \mathcal{S}^n .

The **product** of an $m \times p$ matrix \mathbf{A} with a $p \times n$ matrix \mathbf{B} ,

$$\mathbf{C} = \mathbf{AB},$$

is the $m \times n$ matrix \mathbf{C} whose (i, j) -th entry is given by

$$c_{ij} = \sum_{k=1}^p a_{ik} b_{kj}.$$

It represents a composition of maps $\mathbb{R}^n \rightarrow \mathbb{R}^p \rightarrow \mathbb{R}^m$. The number of columns of \mathbf{A} has to equal the number of rows of \mathbf{B} for this definition to make sense. Products of block matrices or of block matrices with vectors can be carried out block-wise. If, for example, $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)^\top$ with $\mathbf{x}_1 \in \mathbb{R}^{1 \times m}$ and $\mathbf{x}_2 \in \mathbb{R}^{1 \times k}$, then

$$\begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{A}_{11}\mathbf{x}_1 + \mathbf{A}_{12}\mathbf{x}_2 \\ \mathbf{A}_{21}\mathbf{x}_1 + \mathbf{A}_{22}\mathbf{x}_2 \end{pmatrix}.$$

The matrix $\mathbf{1}$ satisfies $\mathbf{1}\mathbf{A} = \mathbf{A}$ and $\mathbf{A}\mathbf{1} = \mathbf{A}$, whenever the dimensions are such that this is defined. In general, even if $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$, $\mathbf{AB} \neq \mathbf{BA}$.

Example 1.7. Let

$$\mathbf{A} = \begin{pmatrix} 1 & 2 \\ 1 & 4 \end{pmatrix}, \quad \mathbf{B} = \begin{pmatrix} 2 & 3 \\ 3 & 2 \end{pmatrix}.$$

Then

$$\mathbf{AB} = \begin{pmatrix} 8 & 7 \\ 14 & 11 \end{pmatrix}, \quad \mathbf{BA} = \begin{pmatrix} 5 & 16 \\ 5 & 14 \end{pmatrix}.$$

If we consider a vector $\mathbf{x} \in \mathbb{R}^n$ as an $n \times 1$ matrix and the transpose as an $1 \times n$ matrix, then for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n \cong \mathbb{R}^{n \times 1}$ we have

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y}.$$

The transpose of a product satisfies $(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$. From this it follows that for any matrix, $\mathbf{A}^\top \mathbf{A}$ is symmetric. For any matrix \mathbf{A} we have

$$\langle \mathbf{x}, \mathbf{Ax} \rangle = \mathbf{x}^\top \mathbf{Ax} = (\mathbf{A}^\top \mathbf{x})^\top \mathbf{x} = \langle \mathbf{A}^\top \mathbf{x}, \mathbf{x} \rangle.$$

It follows from this that if a matrix is symmetric, then it is also self-adjoint, which means that $\langle \mathbf{x}, \mathbf{Ax} \rangle = \langle \mathbf{Ax}, \mathbf{x} \rangle$.

The **rank** of a matrix \mathbf{A} , $\text{rk}(\mathbf{A})$, is the maximum number of linearly independent rows or columns of \mathbf{A} . The **kernel** and **image** of \mathbf{A} are the linear subspaces

$$\ker \mathbf{A} := \{\mathbf{x} \in \mathbb{R}^n : \mathbf{Ax} = 0\}, \quad \text{im } \mathbf{A} = \{\mathbf{Ax} : \mathbf{x} \in \mathbb{R}^n\}.$$

The dimensions are given by $\dim \ker \mathbf{A} = n - \text{rk}(\mathbf{A})$ and $\dim \text{im } \mathbf{A} = \text{rk}(\mathbf{A})$. While $\mathbf{A} \in \mathbb{R}^{m \times n}$ represents a linear map from \mathbb{R}^n to \mathbb{R}^m , the transpose \mathbf{A}^\top represents a map in the other direction, and the image of \mathbf{A}^\top coincides with the orthogonal complement of the kernel of \mathbf{A} , $(\ker \mathbf{A})^\perp = \text{im } \mathbf{A}^\top$. In particular, in view of (1) we have the direct sum decomposition

$$\mathbb{R}^n = \ker \mathbf{A} \oplus \text{im } \mathbf{A}^\top.$$

A system of linear equations

$$\begin{aligned} a_{11}x_1 + \cdots + a_{1n}x_n &= b_1 \\ \vdots &\quad \vdots \quad \vdots \\ a_{m1}x_1 + \cdots + a_{mn}x_n &= b_m \end{aligned}$$

is written as a matrix vector product

$$\mathbf{A}\mathbf{x} = \mathbf{b}, \tag{2}$$

where the $m \times n$ matrix \mathbf{A} is defined as above, and $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^m$. If the columns of \mathbf{A} are linearly independent, then the system of equations can have at most one solution, and otherwise it has infinitely many solutions (this is the case if $n > m$). If $n = m$, then the system (2) has a unique solution if and only if the matrix \mathbf{A} is **invertible** or **non-singular**. This is the case if the rows of \mathbf{A} (or equivalently, the columns of \mathbf{A}) are linearly independent. If \mathbf{A} is not invertible, it is called **singular**.

If \mathbf{A} is invertible, there exists a matrix \mathbf{A}^{-1} (the **inverse**) such that

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}.$$

The solution of (2) is then given by $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$. The following conditions on a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ are equivalent:

1. \mathbf{A} is invertible,
2. $\text{rk}(\mathbf{A}) = n$,
3. $\ker \mathbf{A} = \{\mathbf{0}\}$,
4. $\text{im } \mathbf{A} = \mathbb{R}^n$,
5. the rows of \mathbf{A} are linearly independent,
6. the columns of \mathbf{A} are linearly independent,
7. $\det(\mathbf{A}) \neq 0$,

where the determinant is

$$\det(\mathbf{A}) = \sum_{\sigma \in S_n} \text{sgn}(\sigma) a_{1\sigma(1)} \cdots a_{n\sigma(n)},$$

and S_n is the group of permutations of $[n] = \{1, \dots, n\}$, with $\text{sgn}(\sigma)$ the sign of the permutation (parity of the number of inversions).

Example 1.8. For two- and three-dimensional matrices

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}, \quad \mathbf{B} = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix},$$

the determinants are

$$\det(\mathbf{A}) = a_{11}a_{22} - a_{12}a_{21},$$

$$\det(\mathbf{B}) = a_{11}(a_{22}a_{33} - a_{23}a_{32}) - a_{12}(a_{21}a_{33} - a_{23}a_{31}) + a_{13}(a_{21}a_{32} - a_{22}a_{31}).$$

A matrix \mathbf{Q} is **orthogonal** if $\mathbf{Q} = (\mathbf{q}_1, \dots, \mathbf{q}_n)$, with $\langle \mathbf{q}_i, \mathbf{q}_j \rangle = \delta_{ij}$, and

$$\delta_{ij} = \begin{cases} 0 & \text{if } i \neq j, \\ 1 & \text{if } i = j. \end{cases}$$

As the (i, j) -th entry of $\mathbf{Q}^\top \mathbf{Q}$ are given by $\langle \mathbf{q}_i, \mathbf{q}_j \rangle$, the orthogonality of \mathbf{Q} can succinctly be characterized by the requirement $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$. In particular, $\mathbf{Q}^\top = \mathbf{Q}^{-1}$, and the columns (and rows) of an orthogonal matrix form an orthonormal basis of \mathbb{R}^n . Orthogonal matrices have the property that $\langle \mathbf{Q}\mathbf{x}, \mathbf{Q}\mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle$. From this it follows that orthogonality of vectors is preserved under orthogonal transformations. The determinant of an orthogonal matrix is $\det(\mathbf{Q}) = 1$. As the product of orthogonal matrices is again orthogonal, the set of orthogonal $n \times n$ matrices forms a group, commonly denoted by $O(n)$.

Example 1.9. Consider the three matrices,

$$\mathbf{A} = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}, \quad \mathbf{B} = \begin{pmatrix} 2 & 3 \\ 3 & 2 \end{pmatrix}, \quad \mathbf{C} = \begin{pmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}.$$

The matrices \mathbf{A} and \mathbf{B} are symmetric, while \mathbf{C} is not. The matrices \mathbf{B} and \mathbf{C} are invertible, with inverse

$$\mathbf{B}^{-1} = \begin{pmatrix} -0.4 & 0.6 \\ 0.6 & -0.4 \end{pmatrix}, \quad \mathbf{C}^{-1} = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}.$$

The matrix \mathbf{A} is not invertible, since the second column is a multiple of the first. The kernel of \mathbf{A} is the linear span of $(-2, 1)^\top$. The matrix \mathbf{C} is orthogonal (this can be seen by checking that the columns or rows are orthonormal, or looking at the expression of the inverse above).

Eigenvalues

A vector $\mathbf{u} \neq 0$ is an **eigenvector** of $\mathbf{A} \in \mathbb{R}^{n \times n}$, if there exists a $\lambda \in \mathbb{C}$ such that

$$\mathbf{A}\mathbf{u} = \lambda\mathbf{u}.$$

Such a number λ is called an **eigenvalue** of \mathbf{A} . Note that the eigenvectors are only defined up to scaling: if \mathbf{u} is an eigenvector, then so is $\lambda\mathbf{u}$ for any non-zero $\lambda \in \mathbb{R}$.

From the definition of the determinant, the function $\lambda \mapsto \det(\lambda \mathbf{1} - \mathbf{A})$ is a polynomial of degree at most n , called the **characteristic polynomial** of \mathbf{A} . The eigenvalues are the roots of this polynomial,

$$\det(\lambda \mathbf{1} - \mathbf{A}) = 0.$$

The eigenvalues can be complex numbers, and appear in complex conjugate pairs. If the matrix \mathbf{A} is symmetric, then the eigenvalues are all real numbers. Two important quantities, the **determinant** and the **trace** of a matrix (corresponding, up to sign, to the highest and lowest coefficient of the characteristic polynomial) can be expressed in terms of the eigenvalues:

$$\det(\mathbf{A}) = \lambda_1 \cdots \lambda_n, \quad \text{trace}(\mathbf{A}) := a_{11} + \cdots + a_{nn} = \lambda_1 + \cdots + \lambda_n.$$

A matrix has a zero eigenvalue if and only if it is singular. Eigenvalues may occur with multiplicity.

Norms

A **norm** in \mathbb{R}^n is a function $\|\cdot\|$ that satisfies the following three properties

1. $\|\mathbf{x}\| \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$, and $\mathbf{x} = 0$ if and only if $\mathbf{x} = \mathbf{0}$;
2. $\|\lambda \mathbf{x}\| = |\lambda| \|\mathbf{x}\|$ for $\lambda \in \mathbb{R}$ and $\mathbf{x} \in \mathbb{R}^n$;
3. $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$.

Three important examples of norms are the following:

1. The 1-norm: $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$;
2. The 2-norm: $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$;
3. The ∞ -norm: $\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|$.

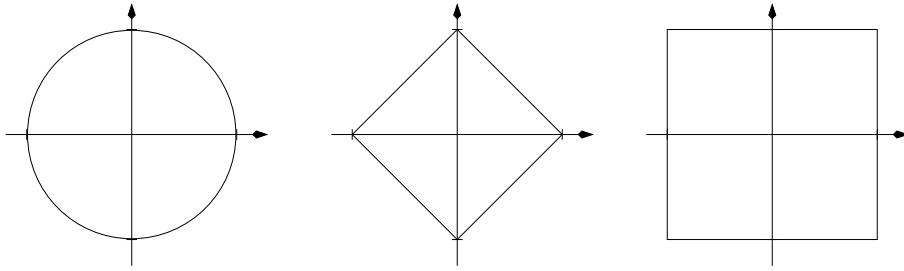
Example 1.10. Let $\mathbf{x} = (2, -3, 4)^\top$. The $\|\mathbf{x}\|_1 = 9$, $\|\mathbf{x}\|_2 = \sqrt{29}$, and $\|\mathbf{x}\|_\infty = 4$.

Note that the 2-norm, also called **Euclidean norm**, can be defined as

$$\|\mathbf{x}\|_2^2 = \mathbf{x}^\top \mathbf{x} = \langle \mathbf{x}, \mathbf{x} \rangle,$$

that is, it is the norm induced by the Euclidean inner product. From this it follows that the 2-norm does not change under orthogonal transformations: if $\mathbf{Q} \in O(n)$ and $\mathbf{x} \in \mathbb{R}^n$, then $\|\mathbf{Q}\mathbf{x}\|_2 = \|\mathbf{x}\|_2$. Orthogonal transformations in \mathbb{R}^2 and \mathbb{R}^3 correspond to rotations and reflections, so it is intuitively clear that these don't change distances.

The **unit sphere** with respect to a norm is the set $\{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| = 1\}$, and the (closed) **unit ball** is the set $\{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| \leq 1\}$. The unit spheres with respect to the 2-norm, the 1-norm and the ∞ -norm in \mathbb{R}^2 are shown in the following diagram.



The unit sphere with respect to the 2-norm in \mathbb{R}^n is usually denoted by S^{n-1} .

The 1-, 2- and ∞ -norms are equivalent, in the sense that they can be bounded in terms of each other. In particular,

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2 \leq \sqrt{n}\|\mathbf{x}\|_\infty, \quad \|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_1 \leq n\|\mathbf{x}\|_\infty. \quad (3)$$

The inner product and the 2-norm are related the **Cauchy-Schwarz** inequality,

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2,$$

with equality if and only if \mathbf{x} and \mathbf{y} are linearly dependent. As a consequence of the Cauchy-Schwarz inequality we get

$$-1 \leq \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2} \leq 1.$$

The **angle** between vectors \mathbf{x} and \mathbf{y} is the number $\theta \in [0, 2\pi]$ such that

$$\cos(\theta) = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2}.$$

If \mathbf{x} and \mathbf{y} are orthogonal, then $\cos(\theta) = 0$ and $\theta \in \{\pi/2, 3\pi/2\}$.

Norms are an important device to measure the size of vectors. In order to measure the amount by which a linear transformation (matrix) distorts vectors, we need the concept of **matrix norms**. A matrix norm is a function on the set of matrices that is a norm when considering a matrix as a vector, and in addition satisfies the condition

$$\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|.$$

The most important examples are given by the **operator norms**. Given a vector norm $\|\mathbf{x}\|$, the associated matrix norm is defined as

$$\|\mathbf{A}\| = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|} = \max_{\mathbf{x}: \|\mathbf{x}\| \leq 1} \|\mathbf{Ax}\|.$$

The matrix norms $\|\mathbf{A}\|_1, \|\mathbf{A}\|_2, \|\mathbf{A}\|_\infty$ are the operator norms that arise when using the 1-, 2- and ∞ -norms. They can conveniently characterized as follows

- $\|\mathbf{A}\|_1 = \max_j \sum_{i=1}^n |a_{ij}|;$
- $\|\mathbf{A}\|_\infty = \max_i \sum_{j=1}^n |a_{ij}|;$

- $\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^\top \mathbf{A})}$.

Here, λ_{\max} denotes the largest eigenvalue of the symmetric matrix $\mathbf{A}^\top \mathbf{A}$. If \mathbf{A} is symmetric, then $\mathbf{A}^\top \mathbf{A} = \mathbf{A}^2$, and the eigenvalues of \mathbf{A}^2 are the squares of the eigenvalues of \mathbf{A} . It follows that for symmetric \mathbf{A} , $\|\mathbf{A}\|_2 = \lambda_{\max}(\mathbf{A})$.

Example 1.11. Let

$$\mathbf{A} = \begin{pmatrix} 1 & -2 \\ -2 & 1 \end{pmatrix}.$$

Then $\|\mathbf{A}\|_1 = \|\mathbf{A}\|_\infty = 3$ and $\|\mathbf{A}\|_2 = 3$.

A special case is the **dual norm** of a vector norm: to a given norm is defined as

$$\|\mathbf{x}\|^* = \max_{\mathbf{y}: \|\mathbf{y}\| \leq 1} \langle \mathbf{x}, \mathbf{y} \rangle.$$

This is the operator norm of \mathbf{x}^\top , considered as a $1 \times d$ matrix. The dual of the dual norm is the norm itself. The dual norm of the 2-norm is again the 2-norm, while the 1-norm and the ∞ -norm are dual to each other.

In addition to the operator norms, an important matrix norm is the **Frobenius norm** of a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$,

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i,j=1}^n a_{ij}^2}.$$

This is just the 2-norm of \mathbf{A} interpreted as a vector in \mathbb{R}^{n^2} . The 2-norm and the Frobenius norm have the important property of being **orthogonal invariant**, which means that for any $\mathbf{Q} \in O(n)$,

$$\|\mathbf{Q}\mathbf{A}\|_2 = \|\mathbf{A}\mathbf{Q}\|_2 = \|\mathbf{A}\|_2, \quad \|\mathbf{Q}\mathbf{A}\|_F = \|\mathbf{A}\mathbf{Q}\|_F = \|\mathbf{A}\|_F.$$

Orthogonal invariance allows to simplify a matrix without changing the norm.

Positive semidefinite matrices

If $\mathbf{A} \in \mathcal{S}^n$ is a symmetric matrix and $\mathbf{u} \in \mathbb{R}^n$ an eigenvector with $\|\mathbf{u}\|_2 = 1$ and corresponding eigenvalue λ , then $\mathbf{u}^\top \mathbf{A} \mathbf{u} = \lambda \mathbf{u}^\top \mathbf{u} = \lambda$. In particular, the largest and smallest values of an eigenvalue are given by

$$\lambda_1 = \max_{\mathbf{u}: \|\mathbf{u}\|_2=1} \mathbf{u}^\top \mathbf{A} \mathbf{u}, \quad \lambda_n = \min_{\mathbf{u}: \|\mathbf{u}\|_2=1} \mathbf{u}^\top \mathbf{A} \mathbf{u}.$$

A symmetric matrix \mathbf{A} is called **positive semidefinite**, written $\mathbf{A} \succeq 0$, if for all non-zero $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$, and **positive definite**, written $\mathbf{A} > 0$, if $\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0$ for all $\mathbf{x} \neq 0$. Equivalently, a symmetric matrix is positive semidefinite if all its eigenvalues are non-negative, and positive definite if they are all positive. The set of positive semidefinite symmetric matrices in $\mathbb{R}^{n \times n}$ is denoted by \mathcal{S}_+^n , while the set of positive definite matrices is \mathcal{S}_{++}^n .

An **inner product** (or scalar product) on $\mathbb{R}^{n \times n}$ is a function

$$\langle \cdot, \cdot \rangle: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}, \quad (\mathbf{x}, \mathbf{y}) \mapsto \langle \mathbf{x}, \mathbf{y} \rangle$$

that is bilinear (linear in each of the two arguments), symmetric ($\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$), and satisfies $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$, with $\langle \mathbf{x}, \mathbf{x} \rangle = 0$ if and only if $\mathbf{x} = 0$. The standard inner product $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y}$ is an example, and the notation $\langle \cdot, \cdot \rangle$ usually refers to this product. More generally, every matrix $\mathbf{A} \in \mathcal{S}_{++}^n$ defines an inner product by

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{A}} := \langle \mathbf{x}, \mathbf{A}\mathbf{y} \rangle = \mathbf{x}^\top \mathbf{A}\mathbf{y}.$$

The associated norm is $\|\mathbf{x}\|_{\mathbf{A}} = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle_{\mathbf{A}}}$. The unit sphere with respect to this norm,

$$E = \{ \mathbf{x} \in \mathbb{R}^n : \mathbf{x}^\top \mathbf{A}\mathbf{x} = 1 \},$$

is an **ellipsoid**, where the 2-norms of the largest and smallest axes are the largest and smallest eigenvalues of \mathbf{A}^{-1} .

Matrix decompositions

Matrices can be represented as products of simpler matrices. Important examples are:

1. **QR decomposition.** A matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ can be written as

$$\mathbf{A} = \mathbf{Q}\mathbf{R},$$

where $\mathbf{Q} \in \mathbb{R}^{m \times m}$ is an orthogonal, and $\mathbf{R} \in \mathbb{R}^{m \times n}$ an upper triangular matrix. Gram-Schmidt orthogonalisation produces such a decomposition.

2. **LU decomposition.** A square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ can be written as

$$\mathbf{A} = \mathbf{L}\mathbf{U}$$

where $\mathbf{L} \in \mathbb{R}^{n \times n}$ is a lower triangular, and $\mathbf{U} \in \mathbb{R}^{n \times n}$ an upper triangular matrix. Gaussian eliminations produces such a decomposition.

3. **Symmetric eigenvalue decomposition.** A symmetric matrix $\mathbf{A} \in \mathcal{S}^n$ can be written as

$$\mathbf{A} = \mathbf{Q}\Lambda\mathbf{Q}^\top,$$

where \mathbf{Q} is an orthogonal matrix, with the eigenvectors as rows, and Λ a *diagonal* matrix with the eigenvalues on the diagonal.

4. **Cholesky decomposition.** A positive definite symmetric matrix $\mathbf{A} \in \mathcal{S}_{++}^n$ can be factored as

$$\mathbf{A} = \mathbf{L}\mathbf{L}^\top,$$

with $\mathbf{L} \in \mathbb{R}^{n \times n}$ lower-triangular and with strictly positive diagonal entries.

One of the most powerful matrix decompositions is the **singular value decomposition** (SVD). It states that any $m \times n$ matrix \mathbf{A} can be written as

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top,$$

where \mathbf{U} is a $m \times m$ orthogonal matrix, \mathbf{V} an $n \times n$ orthogonal matrix, and Σ is a diagonal matrix with the **singular values** $\sigma_1 \geq \dots \geq \sigma_{\min\{m,n\}}$ on the diagonal. The singular values are the square roots of the eigenvalues of $\mathbf{A}^\top \mathbf{A}$. The singular values are related to the matrix 2-norm and Frobenius norm of $\mathbf{A} \in \mathbb{R}^{n \times n}$ as follows:

$$\|\mathbf{A}\|_2 = \sigma_1(\mathbf{A}), \quad \|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^n \sigma_i^2(\mathbf{A})}.$$

If \mathbf{A} is symmetric, the singular values are the absolute values of the eigenvalues of \mathbf{A} .

Matrix decompositions can help reduce a problem into one involving simpler (orthogonal, triangular) matrices. For example, to solve $\mathbf{Ax} = \mathbf{b}$, one can first compute $\mathbf{A} = \mathbf{QR}$, solve the simpler system of equations $\mathbf{Qy} = \mathbf{b}$ by computing $\mathbf{y} = \mathbf{Q}^\top \mathbf{b}$, and then solve the triangular system $\mathbf{Rx} = \mathbf{y}$ by back-substitution.

Example 1.12. Consider the matrix $\begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix}$. The QR decomposition is

$$\mathbf{Q} = \frac{1}{\sqrt{5}} \begin{pmatrix} -2 & 1 \\ 1 & 2 \end{pmatrix}, \quad \mathbf{R} = \frac{1}{\sqrt{5}} \begin{pmatrix} -5 & 4 \\ 0 & 3 \end{pmatrix},$$

the LU decomposition is

$$\mathbf{L} = \begin{pmatrix} 1 & 0 \\ -1/2 & 1 \end{pmatrix}, \quad \mathbf{U} = \begin{pmatrix} 2 & -1 \\ 0 & 3/2 \end{pmatrix},$$

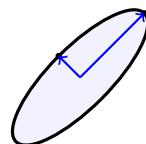
the symmetric eigenvalue decomposition and the SVD are given by

$$\Lambda = \Sigma = \begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{Q} = \mathbf{U} = \mathbf{V} = \begin{pmatrix} \cos(\pi/4) & \sin(\pi/4) \\ -\sin(\pi/4) & \cos(\pi/4) \end{pmatrix},$$

and the Cholesky decomposition is given by

$$\mathbf{L} = \frac{1}{\sqrt{2}} \begin{pmatrix} 2 & 0 \\ -1 & \sqrt{3} \end{pmatrix}.$$

The eigenvalue decomposition shows how to visualize the ellipse $\{\mathbf{Ax} : \|\mathbf{x}\|_2 = 1\}$:



Applying the transformation $\mathbf{Ax} = \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{x}$ corresponds to rotating the vector \mathbf{x} clockwise by an angle of $\pi/4$, then stretching by a factor of 3 in the x -direction, and then rotating back.

3 Calculus

We write $C([a, b]) = C^0([a, b])$ for the set of continuous functions on an interval $[a, b]$, and for $k \geq 1$ we write $C^k([a, b])$ for the set of functions continuous on $[a, b]$, and whose first k derivatives $f', \dots, f^{(k)}$ exist and are continuous on (a, b) . In the above definition we allow $a = -\infty$ or $b = \infty$. If $a, b \in \mathbb{R}$, then any function $f \in C([a, b])$ is bounded. The **infimum** (largest lower bound) and **supremum** (smallest upper bound) of a function f on an interval $[a, b]$ are defined as

$$\inf_{x \in [a, b]} f(x) = \max\{y \in \mathbb{R} : \forall x \in [a, b], f(x) \geq y\},$$

$$\sup_{x \in [a, b]} f(x) = \min\{y \in \mathbb{R} : \forall x \in [a, b], f(x) \leq y\}.$$

Again, we allow the “values” $-\infty$ and ∞ . If the infimum is attained (i.e., there exists x^* such that $f(x^*) = \inf_{x \in [a, b]} f(x)$), then we write $\min_{x \in [a, b]} f(x)$, and similarly max if the supremum is attained. Any $f \in C([a, b])$ for $a, b \in \mathbb{R}$ attains its infimum and supremum on $[a, b]$.

Three important concepts are the **Intermediate Value Theorem**, the **Mean Value Theorem**, and the **Taylor expansion**.

Theorem (Intermediate Value Theorem). If $f \in C([a, b])$ and if y satisfies

$$\inf_{x \in [a, b]} f(x) \leq y \leq \sup_{x \in [a, b]} f(x),$$

then there exists $\xi \in [a, b]$ such that $f(\xi) = y$. In particular, the infimum and supremum are attained.

Theorem (Mean Value Theorem). Let $f \in C^1([a, b])$ and let $x, x_0 \in (a, b)$ with $x \neq x_0$. Then there exists a number $\xi \in (x_0, x)$ (or (x, x_0) if $x < x_0$) such that

$$f(x) = f(x_0) + f'(\xi)(x - x_0).$$

This can also be written as

$$f'(\xi) = \frac{f(x) - f(x_0)}{x - x_0}.$$

The Mean Value Theorem is a special case of the Taylor expansion.

Theorem (Taylor expansion). Let $f \in C^{(n)}([a, b])$ and let $x, x_0 \in (a, b)$ with $x \neq x_0$. Then there exists $\xi \in (x, x_0)$ (or (x, x_0) if $x < x_0$) such that

$$\begin{aligned} f(x) &= f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 + \dots \\ &\quad + \frac{f^{(n)}(x_0)}{n!}(x - x_0)^n + \frac{f^{(n+1)}(\xi)}{(n+1)!}(x - x_0)^{n+1} \end{aligned}$$

The first $(n+1)$ terms of the above sum can be seen as an approximation to the function f that becomes more accurate as n increases. The last term is known as the *truncation error* in numerical approximation.

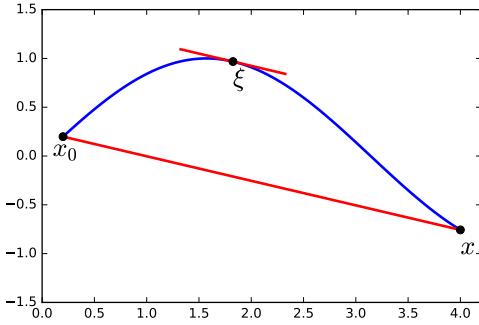


Figure 3.1: MVT: there exists a point at which the slope (derivative) is the same as that of the secant connecting $(x_0, f(x_0))$ and $(x, f(x))$.

As an example, consider the Taylor expansion of the sine function at $x_0 = 0$,

$$\sin(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots$$

The Taylor approximation to different orders is illustrated in the following figure.

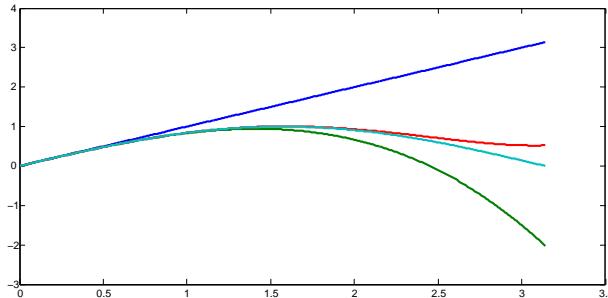


Figure 3.2: Taylor expansion of $\sin(x)$.

A special case of the Mean Value Theorem is **Rolle's Theorem**.

Theorem (Rolle's Theorem). Let $f \in C^1([a, b])$ with $f(a) = f(b)$. Then there exists a number $\xi \in (a, b)$ such that $f'(\xi) = 0$.

The intuition is that if you walk across mountains and arrive at a point with the same elevation as where you started, then there must be places on the way where the slope is 0, that is, a local maximum or minimum. Of importance is also the following variant of the the Mean Value Theorem.

Theorem (Integral Mean Value Theorem). Let $f, g \in C([a, b])$ and assume that $f(x)$ does not change sign on $[a, b]$. Then there exists a $\xi \in (a, b)$ such that

$$\int_a^b f(x)g(x) dx = g(\xi) \int_a^b f(x) dx.$$

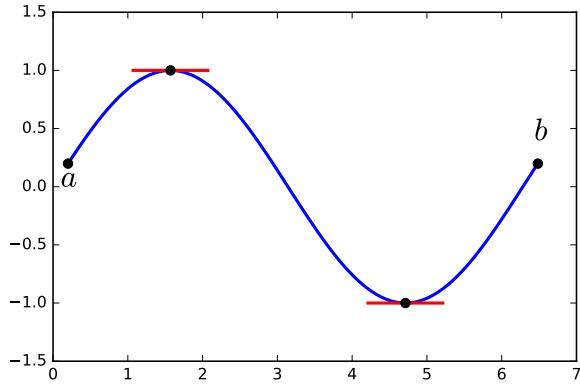


Figure 3.3: Rolle's Theorem: there exist points with “flat” slope (derivative zero).

Topology

The **open ball** of radius ε around $\mathbf{p} \in \mathbb{R}^n$ is defined as

$$B^n(\mathbf{p}, \varepsilon) = \{\mathbf{x} : \|\mathbf{x} - \mathbf{p}\|_2 < \varepsilon\}.$$

We write $B^n := B^n(\mathbf{0}, 1)$ for the (open) **unit ball**. A subset $U \subseteq \mathbb{R}^n$ is called **open** if for every $p \in U$ there exists an $\varepsilon > 0$ such that $B(\mathbf{p}, \varepsilon) \subset U$. A set C is **closed** if $\overline{C} = \mathbb{R}^n \setminus C$ is open. The **closure** $\text{cl } S$ of a set $S \subseteq \mathbb{R}^n$ is the intersection of all closed sets containing S , while the **interior** $\text{int } S$ is the union of all open sets contained in S . The **boundary** of S is defined as $\text{bd } S = \text{cl } S \setminus \text{int } S$. For example, the boundary of the open unit ball is the **unit sphere**

$$\text{bd } B^n = S^{n-1} = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_2 = 1\}.$$

The superscript $n - 1$ refers to the fact that this set is a manifold of dimension $n - 1$. A **neighbourhood** N of a point $\mathbf{x} \in \mathbb{R}^n$ is a set such that there exists an open set U with $\mathbf{x} \in U \subseteq N$. An **open neighbourhood** is a neighbourhood that is open. Note that if $U_1 \subseteq \mathbb{R}^n$ and $U_2 \subseteq \mathbb{R}^m$ are open sets, then the product $U_1 \times U_2 \subseteq \mathbb{R}^{n+m}$ is also open.

Any subset $S \subseteq \mathbb{R}^n$ inherits a topological structure from \mathbb{R}^n , where the open sets in S are the sets of the form $U \cap S$, with $U \subseteq \mathbb{R}^n$ open. If a set S is contained in a lower-dimensional linear subspace $V \subset \mathbb{R}^n$, say, with $\dim V = k < n$, then S is always closed. However, it can be open *relative to its linear span*,

$$\text{span}(S) = \left\{ \sum_{i=1}^k \lambda_i \mathbf{x}_i : \lambda_1, \dots, \lambda_k \in \mathbb{R}, \mathbf{x}_1, \dots, \mathbf{x}_k \in S \right\}.$$

We call a set S **relatively open** or **relatively closed** if it is open or closed in the induced topology on $\text{span}(S)$. Another way of defining this notion goes as follows. Let $\{\mathbf{b}_1, \dots, \mathbf{b}_k\}$ be an orthonormal basis of $\text{span}(S)$, with $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_k)$, and consider the map

$$\varphi_{\mathbf{B}} : \mathbb{R}^k \rightarrow \text{span}(S), \quad \varphi_{\mathbf{B}}(\mathbf{x}) = \sum_{i=1}^k x_i \mathbf{b}_i.$$

Then a set S is relatively open or relatively closed if the preimage

$$\varphi_B^{-1}(S) = \{\mathbf{x} : \varphi_B(\mathbf{x}) \in S\}$$

is open or closed in \mathbb{R}^k . Based on these notions, one defines the **relative closure** $\text{relcl } S$ and **relative interior** $\text{relint } S$ just as before.

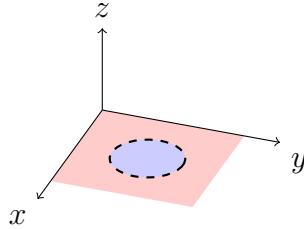


Figure 3.4: The disk without boundary on the xy -plane is relatively open, and is the relative interior of the disk with boundary.

A subset $S \subseteq \mathbb{R}^n$ is **bounded** if there exists number $M > 0$ such that $\|\mathbf{x}\|_2 < M$ for all $\mathbf{x} \in S$. Invoking the equivalence of norms (3) one sees that this definition does not depend on the norm chosen. A set $K \subseteq \mathbb{R}^n$ is called **compact** if it is closed and bounded. Equivalently, every cover of K with open sets contains a finite subcover. A function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is **continuous** if for every open set $U \subset \mathbb{R}^m$,

$$f^{-1}(U) := \{\mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) \in U\}$$

is an open subset of \mathbb{R}^n . A function defined on a subset $S \subseteq \mathbb{R}^n$ is said to be continuous if it is continuous on the induced topology. The set of continuous functions $f: S \rightarrow \mathbb{R}^m$ is denoted by $C(S, \mathbb{R}^m) = C^0(S, \mathbb{R}^m)$. If $f \in C(K, \mathbb{R})$, where K is compact, then f is bounded, and attains its infimum and supremum there: there exist $\mathbf{x}_*, \mathbf{x}^* \in K$ such that

$$\inf_{\mathbf{x} \in K} f(\mathbf{x}) = f(\mathbf{x}_*), \quad \sup_{\mathbf{x} \in K} f(\mathbf{x}) = f(\mathbf{x}^*).$$

A weaker notion is that of a **Lipschitz continuous** function. A function $f: S \rightarrow \mathbb{R}^m$ is called Lipschitz continuous with Lipschitz constant $L > 0$, if for all $\mathbf{x}, \mathbf{y} \in S$,

$$\|f(\mathbf{x}) - f(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2.$$

Notions about continuity of functions can be conveniently stated in terms of sequences. A **sequence** of points $\{\mathbf{x}_k\}_{k \in \mathbb{N}} \subset \mathbb{R}^n$ (for short, $\{\mathbf{x}_k\}$) **converges** to $\mathbf{x} \in \mathbb{R}^n$ as $k \rightarrow \infty$ with respect to a norm $\|\cdot\|$, written $\mathbf{x}_k \rightarrow \mathbf{x}$, if the sequence of numbers $\|\mathbf{x}_k - \mathbf{x}\|$ converges to 0,

$$\lim_{k \rightarrow \infty} \|\mathbf{x}_k - \mathbf{x}\| = 0.$$

Formally, this means that for every $\varepsilon > 0$ there exists an index k_0 , such that for all $k > k_0$, $\|\mathbf{x}_k - \mathbf{x}\| < \varepsilon$. From the equivalence of norms (3) it follows that if a sequence converges with respect to one norm, it also converges with respect to the other norms.

A **subsequence** of a sequence $\{\mathbf{x}_k\}$ is an infinite subset of S . A **limit point** for a sequence $S = \{\mathbf{x}_k\}$ is a point \mathbf{x} that is the limit of an subsequence of S . Formally, for every $\varepsilon > 0$ there exists a k_0 such that $\|\mathbf{x}_k - \mathbf{x}\| < \varepsilon$ for some $k > k_0$. A sequence $\{\mathbf{x}_k\}$ is called a **Cauchy sequence** if for every $\varepsilon > 0$ there exists an index $k_0 > 0$, such that for all $k, \ell > k_0$, $\|\mathbf{x}_k - \mathbf{x}_\ell\|_2 < \varepsilon$. The vector space \mathbb{R}^n with the 2-norm (or any other norm) is a **Banach space**, which means that every Cauchy sequence contains a convergent subsequence.

All the topological notions discussed earlier have an interpretation in terms of sequences and limits:

1. A set C is closed if and only if for every sequence $\{\mathbf{x}_k\} \subset C$, all limit points are in C ;
2. The closure of a set S is the set of all limit points of sequences in S ;
3. A set K is compact, if and only if every sequence of points $\{\mathbf{x}_k\}$ in K has a limit point in K .

Given a function $f: \Omega \rightarrow \mathbb{R}^m$, where $\Omega \subseteq \mathbb{R}^n$, and $\mathbf{x} \in \Omega$, then f is **continuous at \mathbf{x}^*** if

$$\lim_{\mathbf{x} \rightarrow \mathbf{x}^*} f(\mathbf{x}) = f(\mathbf{x}^*).$$

Formally, for every $\varepsilon > 0$ there exists a $\delta > 0$ such that whenever $\|\mathbf{x} - \mathbf{x}^*\| < \delta$, $\|f(\mathbf{x}) - f(\mathbf{x}^*)\| < \varepsilon$. This means that for every sequence of points $\{\mathbf{x}_k\}$ with $\lim_{k \rightarrow \infty} \mathbf{x}_k = \mathbf{x}^*$, the sequence $f(\mathbf{x}_k) \rightarrow f(\mathbf{x}^*)$ as $k \rightarrow \infty$ with respect to some norm on \mathbb{R}^m .

Differentiable functions

A function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is called (Fréchet) **differentiable** at $\mathbf{x}_0 \in \mathbb{R}^n$ if there exists a linear map $\mathbf{J}f(\mathbf{x}_0): \mathbb{R}^n \rightarrow \mathbb{R}^m$, such that

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{\|f(\mathbf{x}_0 + \mathbf{h}) - f(\mathbf{x}_0) - \mathbf{J}f(\mathbf{x}_0)\mathbf{h}\|_2}{\|\mathbf{h}\|_2} = 0.$$

A function is differentiable on an open subset $U \subseteq \mathbb{R}^n$ if it is differentiable at every $\mathbf{x} \in U$. If $f(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_m(\mathbf{x}))^\top$ is differentiable, then all the partial derivatives exist, and $\mathbf{J}f(\mathbf{x}_0)$ is represented by the **Jacobian matrix**

$$\mathbf{J}f(\mathbf{x}_0) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix},$$

where the partial derivatives are evaluated at \mathbf{x}_0 . If all the partial derivatives exist and are continuous in a neighbourhood of \mathbf{x}_0 (called **continuously differentiable**), then f is differentiable at \mathbf{x}_0 .

If $m = 1$, then $J(\mathbf{x}_0)$ is the transpose of the **gradient** $\nabla f(\mathbf{x}_0)$ of f at \mathbf{x}_0 ,

$$\nabla f(\mathbf{x}_0) = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)^\top.$$

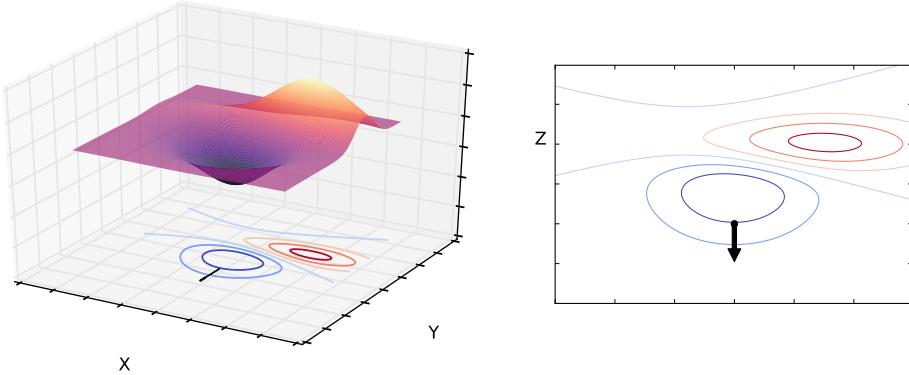


Figure 3.5: A surface, level sets, and the gradient

The gradient points in the direction in which f increases the most.

A convenient way to visualise a function $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ is through **level sets** $\{x \in \mathbb{R}^2 : f(x) = c\}$. For each $c \in \mathbb{R}$, such a level set defines a curve in \mathbb{R}^2 , the curve on which the function value does not change. The gradient is always orthogonal to the level set, pointing in the direction in which f increases the most (see Figure 3.5).

If the gradient, considered as a map $\mathbb{R}^n \rightarrow \mathbb{R}^n$, is itself differentiable at \mathbf{x}_0 , then the Jacobian matrix of the gradient is called the **Hessian matrix**,

$$\nabla^2 f(\mathbf{x}_0) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{pmatrix}$$

Since

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

the Hessian is a symmetric matrix.

The **directional derivative** $D_{\mathbf{x}_0} f(\mathbf{x}_0)$ of a function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ in direction $\mathbf{v} \in \mathbb{R}^n$ at \mathbf{x}_0 is defined as

$$D_{\mathbf{v}} f(\mathbf{x}_0) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x}_0 + h\mathbf{v}) - f(\mathbf{x}_0)}{h}.$$

In the special case where $\mathbf{v} = \mathbf{e}_i$, we obtain the partial derivative with respect to x_i ,

$$\frac{\partial f}{\partial x_i}(\mathbf{x}_0) = D_{\mathbf{e}_i} f(\mathbf{x}_0).$$

If $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable with continuous derivative near \mathbf{x}_0 , then

$$D_{\mathbf{v}} f(\mathbf{x}_0) = \nabla f(\mathbf{x}_0)^{\top} \mathbf{v} = \langle \nabla f(\mathbf{x}_0), \mathbf{v} \rangle.$$

If $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $g: \mathbb{R}^m \rightarrow \mathbb{R}^p$ are differentiable in a neighbourhood of $\mathbf{x}_0 \in \mathbb{R}^n$ and $f(\mathbf{x}_0) \in \mathbb{R}^m$, respectively, then the composition $h = g \circ f: \mathbb{R}^n \rightarrow \mathbb{R}^p$ is continuously differentiable in a neighbourhood of \mathbf{x}_0 , and the Jacobian matrix is defined by the **chain rule**:

$$Jh(\mathbf{x}_0) = Jg(f(\mathbf{x}_0)) Jf(\mathbf{x}_0).$$

If $n = 1$, then $f: \mathbb{R} \rightarrow \mathbb{R}^n$ is called a **curve**, and we write

$$Jf = \frac{df}{dt} = \dot{f} = (\dot{f}_1, \dots, \dot{f}_n)^\top \in \mathbb{R}^n$$

for the derivative of the curve. If $\dot{f}(t_0) = \mathbf{v} \in \mathbb{R}^n$ and $g: \mathbb{R}^n \rightarrow \mathbb{R}$, then by the chain rule, the derivative of $g \circ f: \mathbb{R} \rightarrow \mathbb{R}$ is the directional derivative of g in the direction \mathbf{v} ,

$$\frac{dg \circ f(t_0)}{dt} = \langle \nabla g(f(t_0)), \mathbf{v} \rangle.$$

Before going on to deal with higher derivative, we state the generalisation of the Mean Value Theorem to higher dimensions.

Theorem (Multivariate Mean Value Theorem). Let $f \in C^1(U)$ for an open set U with $\mathbf{x}_0, \mathbf{x} \in U$, $\mathbf{x} \neq \mathbf{x}_0$. Then there exists $t \in (0, 1)$ such that

$$f(\mathbf{x}) - f(\mathbf{x}_0) = \langle \nabla f(t\mathbf{x} + (1-t)\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle.$$

Note that $t\mathbf{x} + (1-t)\mathbf{x}_0$ parametrises the line segment connecting \mathbf{x} and \mathbf{x}_0 .

For a tuple of natural number $\alpha = (\alpha_1, \dots, \alpha_n)$, set $|\alpha| = \sum_{i=1}^n \alpha_i$, and define the higher order partial derivative

$$D^\alpha f(\mathbf{x}) = \frac{\partial^{|\alpha|} f(\mathbf{x})}{\partial^{\alpha_1} x_1 \cdots \partial^{\alpha_n} x_n}.$$

For a set $S \subseteq \mathbb{R}^n$, denote by $C^k(S, \mathbb{R}^m)$ the set of functions f such that all partial derivatives $D^\alpha f$ with $|\alpha| \leq k$ exists and are continuous on $\text{int } S$. If $m = 1$, we write $C^k(S) := C^k(S, \mathbb{R})$.

Define, for a vector \mathbf{x} and multi-index α ,

$$\mathbf{x}^\alpha := x_1^{\alpha_1} \cdots x_n^{\alpha_n}.$$

We then have the **Taylor expansion** around a point \mathbf{x}_0 ,

$$f(\mathbf{x}) = \sum_{|\alpha| \leq k} D^\alpha f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0)^\alpha + \sum_{|\alpha|=k} r_\alpha(\mathbf{x})(\mathbf{x} - \mathbf{x}_0)^\alpha,$$

with $r_\alpha(\mathbf{x}) \rightarrow 0$ as $\mathbf{x} \rightarrow \mathbf{x}_0$.

If a differentiable function $f(\mathbf{x})$ has a local minimum or maximum at a point \mathbf{x} , then this point satisfies $\nabla f(\mathbf{x}) = 0$, that is, it is a critical point. The **Lagrange multiplier theorem** says something about local extrema under certain constraints.

Theorem (Lagrange multipliers). Let \mathbf{x}^* be maximum of $f(\mathbf{x})$ under the constraint $g(\mathbf{x}) = c$ (that is, a maximum among all points \mathbf{x} such that $g(\mathbf{x}) = c$). Then there exist a $\lambda \in \mathbb{R}$ such that

$$\nabla f(\mathbf{x}^*) = \lambda \nabla g(\mathbf{x}^*).$$

The **Lagrangian** of a function $f(\mathbf{x})$ with constraint $g(\mathbf{x}) = c$ is the function $\Lambda: \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}$ defined by

$$\Lambda(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda(g(\mathbf{x}) - c).$$

The Lagrange multiplier theorem then says that if \mathbf{x}^* is a maximum point of $f(\mathbf{x})$ under the constraint $g(\mathbf{x}) = c$, then there exists $\lambda \in \mathbb{R}$ such that the pair (\mathbf{x}^*, λ) is a critical point of the Lagrangian $\Lambda(\mathbf{x}, \lambda)$.

The **Implicit Function Theorem** is one of the most important results in analysis, and underlies much of differential geometry and physics. Let $F: \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ by a function that is continuously differentiable in a neighbourhood of a point $(\mathbf{x}_0, \mathbf{y}_0)$, with $\mathbf{x}_0 \in \mathbb{R}^n$ and $\mathbf{y}_0 \in \mathbb{R}^m$. The Jacobian $J_{\mathbf{x}}(\mathbf{x}_0, \mathbf{y}_0)$ with respect to the first set of n coordinates consists of the first n columns of the Jacobian matrix,

$$J_{\mathbf{x}}(\mathbf{x}_0, \mathbf{y}_0) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{x}_0, \mathbf{y}_0) & \cdots & \frac{\partial f_1}{\partial x_n}(\mathbf{x}_0, \mathbf{y}_0) \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{x}_0, \mathbf{y}_0) & \cdots & \frac{\partial f_m}{\partial x_n}(\mathbf{x}_0, \mathbf{y}_0) \end{pmatrix}$$

The interpretation is that we consider f as a function in only the first set of coordinates, with the remaining ones (denoted by \mathbf{y}) considered as parameters.

Theorem (Implicit Function Theorem). Let $f: \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ be k times continuously differentiable in an open neighbourhood of $(\mathbf{x}_0, \mathbf{y}_0) \in \mathbb{R}^n \times \mathbb{R}^m$, and assume that $f(\mathbf{x}_0, \mathbf{y}_0) = \mathbf{0}$. Assume further that the Jacobian $J_{\mathbf{x}}f(\mathbf{x}_0, \mathbf{y}_0) \in \mathbb{R}^{n \times n}$ in the first n coordinates is *non-singular* at $(\mathbf{x}_0, \mathbf{y}_0)$. Then there exists an open neighbourhood $\mathbf{y}_0 \in U_{\mathbf{y}} \subseteq \mathbb{R}^m$, and a function $h \in C^k(U_{\mathbf{y}}, \mathbb{R}^n)$ such that

- $h(\mathbf{y}_0) = \mathbf{x}_0$,
- $f(h(\mathbf{y}), \mathbf{y}) = \mathbf{0}$ for $\mathbf{y} \in U_{\mathbf{y}}$.

Moreover, the Jacobian of h is given by

$$Jh(\mathbf{y}) = -J_{\mathbf{y}}f(h(\mathbf{y}), \mathbf{y})(J_{\mathbf{x}}f(h(\mathbf{y}), \mathbf{y}))^{-1}$$

for all $\mathbf{y} \in U_{\mathbf{y}}$.

Example 1.13. Let $f(x, y) = x^2 + y^2 - 1$ and $(x_0, y_0) = (1, 0)$. The Jacobian of f in the first coordinate is

$$\frac{\partial f}{\partial x}(1, 0) = 2 \neq 0,$$

which is non-singular. Choosing the neighbourhood $U_{\mathbf{y}} = (-1, 1)$, the open interval between -1 and 1 , we get the function $h: U_{\mathbf{y}} \rightarrow \mathbb{R}$ as

$$h(y) = \sqrt{1 - y^2}.$$

This function is defined and differentiable on $(-1, 1)$, and satisfies

$$f(h(y), y) = h(y)^2 + y^2 - 1 = (1 - y^2) + y^2 - 1 = 0, \quad y \in U_{\mathbf{y}}.$$

The derivative of h can be computed using the chain rule:

$$\frac{df(h(y), y)}{dy} = \frac{\partial f}{\partial x}(h(y), y) \frac{dh}{dy}(y) + \frac{\partial f}{\partial y}(h(y), y) \frac{dy}{dy}$$

from which we get

$$\frac{dh}{dy}(y) = -\frac{\partial f}{\partial y}(h(y), y) \left(\frac{\partial f}{\partial x}(h(y), y) \right)^{-1} = -y(1-y^2)^{-3/2}.$$

In this example, the implicit function theorem just gives the usual way of parametrising part of the circle.

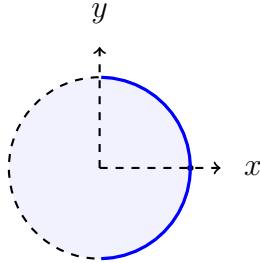


Figure 3.6: The blue arc is parametrised by $h(y) = \sqrt{1-y^2}$ for $y \in (-1, 1)$.

4 Finite precision arithmetic

In practical applications, one often cannot simply plug numbers into formulae and get all the exact results. Most numerical data also requires an infinite amount of storage (just try to store π on a computer!), but a piece of paper or a computer only has limited space. These are some of the reasons that lead us to work with **approximations**.

Measuring errors

To measure the quality of approximations, we use the concept of **relative error**. Given a quantity x and a computed approximation \hat{x} , the **absolute error** is given by

$$E_{\text{abs}}(\hat{x}) = |x - \hat{x}|,$$

while the *relative error* is given as

$$E_{\text{rel}}(\hat{x}) = \frac{|x - \hat{x}|}{|x|}.$$

The benefit of working with relative errors is that they are scale invariant. Absolute errors can be meaningless at times: for example, an error of one hour is irrelevant when estimating the age of Stan the Tyrannosaurus Rex at Manchester Museum, but it is crucial when determining the time of a lecture. That is because in the former one hour corresponds to a relative error is of the order 10^{-11} , while in the latter it is of the order 1.

Floating point and significant figures

The established way of representing real numbers on computers is using **floating-point arithmetic**. In the double precision version of the IEEE standard for floating-point arithmetic, a number is represented using 64 bits, where a bit is either 1 or 0. A number is written

$$x = \pm f \times 2^e,$$

where f is a fraction in $[0, 1]$, represented using 52 bits, and e is the exponent, using 11 bits, and one bit is for the sign. There are largest possible numbers, and there are gaps between representable numbers. The largest and smallest numbers representable in this form are of the order of $\pm 10^{308}$, enough for most practical purposes. A bigger concern are the gaps, which means that the results of many computations almost always have to be rounded to the closest floating-point number.

When going through calculations without using a computer, we usually use the terminology of **significant figures** (s.f.) and work with 4 significant figures in base 10. For example, in base 10, $\sqrt{3}$ equals 1.732 to 4 significant figures. To count the number of significant figures in a given number, start with the first non-zero digit from the left and, moving to the right, count all the digits thereafter, counting final zeros if they are to the right of the decimal point. For example, 1.2048, 12.040, 0.012048, 0.0012040 and 1204.0 all have 5 significant figures (s.f.). In rounding or truncation of a number to n s.f., the original is replaced by the closest number with n s.f. An approximation \hat{x} of a number x is said to be **correct to n significant figures** if both \hat{x} and x round to the same n s.f. number.

Remark 1.14. Note that final zeros to left of the decimal point may or may not be significant: the number 1204000 has at least 4 significant figures, but without any more information there is no way of knowing whether or not any more figures are significant. When 1203970 is rounded to 5 significant figures to give 1204000, an explanation that this has 5 significant figures is required. This could be made clear by writing it in scientific notation: 1.2040×10^6 . In some cases we also have to agree whether to round up or round down: for example, 1.25 could equal 1.2 or 1.3 to two significant figures. If we agree on rounding up, then to say that $a = 1.2048$ to 5 s.f. means that the exact value of a satisfies $1.20475 \leq a < 1.40485$.

Example 1.15. Suppose we want to find the solution to the quadratic equation

$$ax^2 + bx + c = 0.$$

The two solutions to this problem are given by

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a}, \quad x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}. \quad (1.1)$$

In principle, to find x_1 and x_2 one only needs to evaluate the expressions for given a, b, c . Assume, however, that we are only allowed to compute to four significant figures, and consider the particular equation

$$x^2 + 39.7x + 0.13 = 0.$$

Using the formula 1.1, we have, always rounding to four significant figures,

$$a = 1, b = 39.7, c = 0.13,$$

$$b^2 = 1576.09 = 1576 \text{ (to 4 s.f.)}, 4ac = 0.52 \text{ (to 4 s.f.)},$$

$$b^2 - 4ac = 1575.48 = 1575 \text{ (to 4 s.f.)}, \sqrt{b^2 - 4ac} = 39.69.$$

Hence, the computed solutions (to 4 significant figures) are given by

$$\bar{x}_1 = -0.005, \bar{x}_2 = -39.69$$

The exact solutions, however, are

$$x_1 = -0.0032748..., x_2 = -39.6907...$$

The solution x_1 is completely wrong, at least if we look at the relative error:

$$\frac{|\bar{x}_1 - x_1|}{|x_1|} = 0.5268.$$

While the accuracy can be increased by increasing the number of significant figures during the calculation, such effects happen all the time in scientific computing and the possibility of such effects has to be taken into account when designing numerical algorithms.

By analysing what causes the error it is sometimes possible to modify the method of calculation in order to improve the result. In the present example, the problems are being caused by the fact that $b \approx \sqrt{b^2 - 4ac}$, and therefore

$$\frac{-b + \sqrt{b^2 - 4ac}}{2a} = \frac{-39.7 + 39.69}{2}$$

causes what is called “catastrophic cancellation”. A way out is provided by the observation that the two solutions are related by

$$x_1 \cdot x_2 = \frac{c}{a}. \quad (1.2)$$

When $b > 0$, the calculation of x_2 according to (1.1) shouldn't cause any problems, in our case we get -39.69 to four significant figures. We can then use (1.2) to derive $\bar{x}_1 = c/(a\bar{x}_2) = -0.00327$.

There are other potential sources of error besides those introduced by rounding operations.

1. Overflow
2. Errors in the model
3. Human or measurements errors
4. Truncation or approximation errors

The first is rarely an issue, as we can represent numbers of order 10^{308} on a computer. The second two are important factors that need to be addressed when working on real-world problems. The fourth has to do with the fact that many computations are done approximately rather than exactly. For computing the exponential, for example, we might use a method that gives the approximation

$$e^x \approx 1 + x + \frac{x^2}{2}.$$

As it turns out, many optimization problems work with approximations of the functions of interest, and the solution found is only an approximation to the “true” solution of the problem. Quantifying the quality of such an approximation is an important aspect in the design and analysis of optimization algorithms.