

Mathematical Methods in Data Science

Paul Breiding, Alessio D'Alì and Samantha Fairchild

Author's addresses:

Paul Breiding, University of Osnabrück + MPI MiS Leipzig, pbreiding@uni-osnabrueck.de.

Alessio D'Alì, University of Osnabrück, alessio.dali@uni-osnabrueck.de.

Samantha Fairchild, University of Osnabrück + MPI MiS Leipzig, samantha.fairchild@mis.mpg.de

These lecture notes were written in the Summer Semester 2022, when the authors gave the class "Mathematische Grundlagen der Datenanalyse" at the University of Osnabrück. The notes are intended for a course on an advanced Bachelor level. A main theme of these notes is studying the **geometry of data**.

The first and third authors have been funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – Projektnummer 445466444.

Contents

1	The Basics	1
1.1	Linear Algebra	1
1.2	Probability Theory	9
2	Network Analysis	18
2.1	Graphs and the Laplace Matrix	18
2.2	The Spectrum of a Graph	25
2.3	Markov Processes in Networks	37
2.4	Centrality Measures	46
3	Machine Learning	52
3.1	Data, Models, and Learning	52
3.2	Nonlinear Regression and Neural Networks	56
3.3	Support Vector Machines	66
3.4	Principal Component Analysis	74

1 The Basics

Many mathematical methods in data analysis rely on linear algebra and probability. In the first two lectures we will recall basic concepts from these fields.

1.1 Linear Algebra

This lecture is based on the article *The Fundamental Theorem of Linear Algebra* by Gilbert Strang [Str93]. We will use the following notation:

$$A = (a_{ij}) \in \mathbb{R}^{m \times n} \text{ (resp. } \mathbb{C}^{m \times n} \text{)}$$

is an $m \times n$ **matrix** with real (resp. complex) entries a_{ij} for $1 \leq i \leq m$, $1 \leq j \leq n$. The column vectors are

$$a_j := (a_{ij})_{i=1}^m.$$

A matrix $A \in \mathbb{R}^{m \times n}$ can be viewed as a **list** of vectors in \mathbb{R}^m which we denote by

$$A = [a_1, \dots, a_n].$$

For $x = (x_1, \dots, x_n)^T \in \mathbb{R}^n$,

$$Ax = x_1 a_1 + \dots + x_n a_n$$

is a **linear combination** of the columns of A . Other interpretations of A are

1. a list of n vectors in \mathbb{R}^m
2. a list of m vectors in \mathbb{R}^n
3. a linear map $\mathbb{R}^n \rightarrow \mathbb{R}^m$ given by $x \mapsto Ax$
4. a linear map $\mathbb{R}^m \rightarrow \mathbb{R}^n$ given by $y \mapsto A^T y$
5. a bilinear map $\mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ given by $(x, y) \mapsto y^T Ax$.

1 The Basics



Figure 1.1: The meaning of the inner product between u and v is illustrated in this picture: let $t \in \mathbb{R}$ such that $v = tu + u'$, where u' is orthogonal to u . Then, $\langle u, v \rangle = u^T(tu + u') = tu^T u = t\langle u, u \rangle$. In particular, if $\langle u, u \rangle = \langle v, v \rangle = 1$, then $t = \langle u, v \rangle$ is the arccosine of the angle between u and v .

All of these viewpoints are best understood by considering four subspaces (two subspaces of \mathbb{R}^n and two of \mathbb{R}^m).

Definition 1.1 (Four Subspaces). Let $A \in \mathbb{R}^{m \times n}$. The **image** and **kernel** of A and A^T are

1. $\text{Im}(A) := \{Ax \mid x \in \mathbb{R}^n\} \subseteq \mathbb{R}^m$,
2. $\text{Im}(A^T) := \{A^T y \mid y \in \mathbb{R}^m\} \subseteq \mathbb{R}^n$,
3. $\ker(A) := \{x \in \mathbb{R}^n \mid Ax = 0\} \subseteq \mathbb{R}^n$,
4. $\ker(A^T) := \{y \in \mathbb{R}^m \mid A^T y = 0\} \subseteq \mathbb{R}^m$

We give the \mathbb{R} -vector spaces \mathbb{R}^n and \mathbb{R}^m the structure of a Euclidean space by defining the positive definite form $\langle a, b \rangle := a^T b$. For $A = [a_1, \dots, a_n]$, we have $A^T y = [\langle a_i, y \rangle]_{i=1}^n$.

Definition 1.2. Let $U, V \subseteq \mathbb{R}^n$ be subspaces. Then U is **perpendicular to** V (denoted $U \perp V$) when for all $u \in U$ and $v \in V$, $\langle u, v \rangle = 0$.

In the following we will denote $r = r(A)$ to be the **rank** of A .

Theorem 1.3. Let $A \in \mathbb{R}^{m \times n}$. Then

1. $\text{Im}(A) \oplus \ker(A^T) = \mathbb{R}^m$,
2. $\text{Im}(A) \perp \ker(A^T)$,
3. $\text{Im}(A^T) \oplus \ker(A) = \mathbb{R}^n$,
4. $\text{Im}(A^T) \perp \ker(A)$.

Proof of (1) and (2). From linear algebra, we know that

$$r(A) = \dim(\text{Im}(A)) = \dim(\text{Im}(A^T))$$

1 The Basics



Figure 1.2: The situation when $b \in \text{Im}(A)$: in this case, $Ax = b$ has a unique solution $x \in \text{Im}(A^T)$ and the solution space for $Ax = b$ is $x + \ker(A)$.

and by the Rank-Nullity theorem

$$\dim(\text{Im}(A^T)) + \dim(\ker(A^T)) = m.$$

Therefore

$$\dim(\text{Im}(A)) + \dim(\ker(A^T)) = m.$$

Moreover for $y \in \ker(A^T)$ and $Ax \in \text{Im}(A)$,

$$\langle y, Ax \rangle = y^T Ax = (A^T y)^T x = 0.$$

Thus $\text{Im}(A) \perp \ker(A^T)$ and in particular $\text{Im}(A) \cap \ker(A^T) = \{0\}$ and

$$\dim(\text{Im}(A) + \ker(A^T)) = \dim(\text{Im}(A)) + \dim(\ker(A^T)) = m.$$

Thus $\text{Im}(A) \oplus \ker(A^T) = \mathbb{R}^m$. The proof of (3) and (4) follows similarly. \square

We now want to understand the solution of the system of linear equations $Ax = b$ in the context of Theorem 1.3. Namely, let $b \in \text{Im}(A)$ and let $r = \dim(\text{Im}(A)) = \dim(\text{Im}(A^T))$. First, we observe that $Ax = b$ has a solution $x \in \mathbb{R}^n$, if and only if $b \in \text{Im}(A)$. Suppose that x is such a solution. This situation is depicted in Fig. 1.2. From Theorem 1.3 we know that $\text{Im}(A^T) \oplus \ker(A) = \mathbb{R}^n$. So, there exist uniquely determined $x_0 \in \text{Im}(A^T)$ and $x_1 \in \ker(A)$ with $x = x_0 + x_1$ and we have

$$b = Ax = A(x_0 + x_1) = Ax_0 + Ax_1 = Ax_0.$$

1 The Basics



Figure 1.3: Visualization of the proof of Lemma 1.4: b_0 minimizes the distance from b to $\text{Im}(A)$.

Therefore, $Ax = b$ has a **unique** solution in $\text{Im}(A^T)$. Consequently, A restricted to $\text{Im}(A^T)$ is a **linear isomorphism**.

When $b \notin \text{Im}(A)$ there is no solution to $Ax = b$. We can however find the point $b_0 \in \text{Im}(A)$ which minimizes the Euclidean distance $\|b - b_0\| = \sqrt{\langle b - b_0, b - b_0 \rangle}$. We use the notation

$$b_0 = \operatorname{argmin}_{y \in \text{Im}(A)} \|b - y\|$$

to denote the argument (i.e. the value $y = b_0$) which minimizes the function $\|b - y\|$. The solution to this minimization problem and the fact that b_0 is uniquely determined is given by the next lemma.

Lemma 1.4. *Let $b \in \mathbb{R}^m$ and $A \in \mathbb{R}^{m \times n}$. The point $b_0 = \operatorname{argmin}_{y \in \text{Im}(A)} \|b - y\|$ is determined by*

1. *the decomposition from Theorem 1.3, which gives $b = b_0 + c$ for $c \in \ker(A^T)$; or*
2. $A^T b = A^T b_0$.

Proof. (2. \Rightarrow 1.) This direction follows because $e \in \ker(A^T)$. This also shows that b_0 is uniquely determined.

So now it suffices to prove (2.) Let $A = [a_1, \dots, a_n]$. Since $b_0 \in \text{Im}(A)$, set $b_0 = Ax_0$ for some x_0 . Define the map $\phi(x) = Ax - b$. Suppose that we write the output vector in \mathbb{R}^m by $\phi = [\phi_1, \dots, \phi_m]^T$. Then, we minimize the scalar function $\|\phi(x)\|$ by taking the

1 The Basics



Figure 1.4: The pseudoinverse $A^\dagger \in \mathbb{R}^{m \times n}$ of $A \in \mathbb{R}^{m \times n}$ first orthogonally projects $b \in \mathbb{R}^n$ to $b_0 \in \text{Im}(A)$ and then maps b_0 to the unique point $x \in \text{Im}(A^T)$ with $Ax = b_0$.

derivative and setting it equal to zero. Namely we want to compute when the gradient $\frac{d}{dx} \|\phi(x)\| = \left[\frac{\partial}{\partial x_1} \|\phi(x)\|, \dots, \frac{\partial}{\partial x_n} \|\phi(x)\| \right] \in \mathbb{R}^n$ is equal to zero. We compute

$$\frac{\partial}{\partial x_i} \|\phi(x)\| = \frac{1}{2\|\phi(x)\|} \sum_{j=1}^m 2\phi_j(x) \frac{\partial \phi_j}{\partial x_i}(x) = \frac{1}{\|\phi(x)\|} a_i^T (Ax - b).$$

If $\|\phi(x_0)\| = 0$, then $b_0 = b$ and we are done. Otherwise we must have that x_0 satisfies $A^T(Ax - b) = 0$. This implies $A^T Ax_0 = A^T b$, and so $A^T b_0 = A^T b$. \square

Lemma 1.4 implies that the map which projects b to the point on $\text{Im}(A)$ minimizing the distance to b in the Euclidean norm is linear: call it $\Pi_{\text{Im}(A)}$. Furthermore, recall that A restricted to $\text{Im}(A^T)$ is a linear isomorphism, hence invertible. Consequently, we have a well-defined linear map $(A|_{\text{Im}(A^T)})^{-1} \circ \Pi_A : \mathbb{R}^m \rightarrow \mathbb{R}^n$, shown in Fig. 1.4. The matrix representation of this linear map is called the **pseudoinverse** of A .

Definition 1.5. Let $A \in \mathbb{R}^{m \times n}$. The **pseudoinverse** $A^\dagger \in \mathbb{R}^{n \times m}$ is the matrix such that

$$A^\dagger b = x$$

for $x \in \text{Im}(A^T)$, $Ax = b_0$ and $b_0 = \arg\min_{y \in \text{Im}(A)} \|b - y\|$.

Note when consulting other texts sometimes A^+ is used instead A^\dagger .

Let us first notice two properties of the pseudoinverse, which follow immediately from the definition.

1 The Basics

Corollary 1.6. Let $A \in \mathbb{R}^{m \times n}$ and $A^\dagger \in \mathbb{R}^{n \times m}$ be its pseudoinverse.

1. If A is invertible, then $A^\dagger = A^{-1}$.
2. AA^\dagger is the orthogonal projection onto $\text{Im}(A)$; i.e., $AA^\dagger = \Pi_{\text{Im}(A)}$.

In the case when $A \in \mathbb{R}^{m \times n}$ has full rank, which means that $r(A) = \min\{m, n\}$, the pseudoinverse has the following properties.

Proposition 1.7. Let $A \in \mathbb{R}^{m \times n}$ have full rank.

1. If $r(A) = n$,

$$A^\dagger = (A^T A)^{-1} A^T$$

and $A^\dagger A = \mathbf{1}_n$. So A is left-invertible.

2. If $r(A) = m$,

$$A^\dagger = A^T (AA^T)^{-1}$$

and $AA^\dagger = \mathbf{1}_m$. So A is right-invertible.

Proof. Let $b \in \mathbb{R}^m$ and $A^\dagger b = x$. By Lemma 1.4 we have $A^T A x = A^T b$, which implies

$$A^T A A^\dagger b = A^T b.$$

Since $r(A) = n$, the matrix $A^T A \in \mathbb{R}^{n \times n}$ is invertible, so that

$$A^\dagger b = (A^T A)^{-1} A^T b.$$

This shows $A^\dagger = (A^T A)^{-1} A^T$ and it also shows $A^\dagger A = (A^T A)^{-1} A^T A = \mathbf{1}_n$. For the second part, see Exercise 1.2. \square

In closing of this lecture we want to discuss an important choice of bases for $\text{Im}(A)$ and $\text{Im}(A^T)$. Let $r := r(A)$. The matrix $A^T A \in \mathbb{R}^{n \times n}$ is symmetric. By the spectral theorem, there is a basis $\{v_1, \dots, v_n\}$ of eigenvectors of $A^T A$ such that $\langle v_i, v_j \rangle = \delta_{i,j}$ (called an **orthonormal basis**). Let λ_i be the eigenvalue corresponding to v_i ; i.e., $A^T A v_i = \lambda_i v_i$. Since $A^T A$ is positive semidefinite, it has only real nonnegative eigenvalues. We can assume that $\lambda_1 \geq \dots \geq \lambda_r > 0$ and $\lambda_{r+1} = \dots = \lambda_n = 0$. We have

$$\text{span}\{v_1, \dots, v_r\} = \text{span}\{v_{r+1}, \dots, v_n\}^\perp = \text{Im}(A^T) = \text{ker}(A)^\perp$$

1 The Basics

(the last equality because of Theorem 1.3), so $\{v_1, \dots, v_r\}$ is an orthonormal basis for $\text{Im}(A^T)$ and $\{v_{r+1}, \dots, v_n\}$ is an orthonormal basis for $\ker(A)$. Let

$$u_i := \frac{1}{\sqrt{\lambda_i}} A v_i, \quad i = 1, \dots, r.$$

Then, by construction we have

$$\langle u_i, u_j \rangle = \frac{1}{\sqrt{\lambda_i \lambda_j}} v_i^T A^T A v_j = \delta_{i,j},$$

which shows that $\{u_1, \dots, u_r\}$ is an orthonormal basis for $\text{Im}(A)$. We have

$$A v_i = \sigma_i u_i, \quad \sigma_i = \sqrt{\lambda_i}.$$

For $U = [u_1, \dots, u_r] \in \mathbb{R}^{m \times r}$, $V = [v_1, \dots, v_r] \in \mathbb{R}^{n \times r}$ and $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r)$ we have

$$A = U \Sigma V^T.$$

This decomposition is called the **singular value decomposition** (SVD) of A and $\sigma_1, \dots, \sigma_r$ are called the **singular values** of A . The next theorem shows that the SVD is essentially unique.

Theorem 1.8. *Let $A \in \mathbb{R}^{m \times n}$ and $r = r(A)$. Then, there exist matrices $U \in \mathbb{R}^{m \times r}$ and $V \in \mathbb{R}^{n \times r}$ with $U^T U = V^T V = \mathbf{1}_r$ and uniquely determined numbers $\sigma_1, \dots, \sigma_r > 0$ such that*

$$A = U \Sigma V^T, \quad \Sigma = \text{diag}(\sigma_1, \dots, \sigma_r).$$

We have $\text{Im}(A) = \text{Im}(U)$ and $\text{Im}(A^T) = \text{Im}(V)$. If the σ_i are pairwise distinct and ordered $\sigma_1 > \dots > \sigma_r > 0$, the matrices U and V are uniquely determined up to the signs of their columns.

Proof. Existence of the SVD and $\text{Im}(A) = \text{Im}(U)$ and $\text{Im}(A^T) = \text{Im}(V)$ follow from the discussion above. We have to show uniqueness of singular values, and in the case when the singular values are pairwise distinct uniqueness of U and V (up to sign). Suppose

$$A = U \Sigma V^T = \tilde{U} \tilde{\Sigma} \tilde{V}^T$$

are two SVDs of A with $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r)$ and $\tilde{\Sigma} = \text{diag}(\tilde{\sigma}_1, \dots, \tilde{\sigma}_r)$. Then, we have

$$A A^T = U \Sigma V^T V \Sigma U^T = U \Sigma^2 U^T \quad \text{and} \quad A A^T = \tilde{U} \tilde{\Sigma} \tilde{V}^T \tilde{V} \tilde{\Sigma} \tilde{U}^T = \tilde{U} \tilde{\Sigma}^2 \tilde{U}^T,$$

1 The Basics

because $V^T V = \tilde{V}^T \tilde{V} = \mathbf{1}_r$. Let us write $U = [u_1, \dots, u_r]$ and $\tilde{U} = [\tilde{u}_1, \dots, \tilde{u}_r]$. The above equations imply that $AA^T u_i = \sigma_i^2 u_i$ and $AA^T \tilde{u}_i = \tilde{\sigma}_i^2 \tilde{u}_i$ for $i = 1, \dots, r$. Using that $r = r(A) = r(AA^T)$ we conclude that both $\sigma_1, \dots, \sigma_r$ and $\tilde{\sigma}_1, \dots, \tilde{\sigma}_r$ are the nonzero eigenvalues of AA^T . Since eigenvalues are unique, this implies $\sigma_i = \tilde{\sigma}_i$, $i = 1, \dots, r$. Therefore, the singular values are uniquely determined.

Let us now assume that the σ_i are pairwise distinct. Then, since every σ_i is positive, also the σ_i^2 are pairwise distinct for $i = 1, \dots, r$. This means that the nonzero eigenvalues of AA^T are all simple, which implies that the eigenvector of σ_i is unique up to sign, hence $u_i = \pm \tilde{u}_i$. Repeating the same argument for $A^T A$ shows that the columns of V and \tilde{V} also coincide up to sign. \square

An alternative definition of the SVD is $A = U S V^T$ for $U \in \mathbb{R}^{m \times k}$ and $V \in \mathbb{R}^{n \times k}$ with $k = \min\{m, n\}$ and $U^T U = V^T V = \mathbf{1}_k$, and $S = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0)$. This is sometimes called the **non-compact SVD**, while the decomposition in Theorem 1.8 is called **compact SVD**. The difference between the two is that the compact SVD involves orthonormal bases of $\text{Im}(A)$ and $\text{Im}(A^T)$, while the non-compact SVD appends orthonormal vectors from $\ker(A^T)$ and $\ker(A)$. The way one should think about the SVD (compact or non-compact) is that it provides particular orthonormal bases that reveal essential information about the matrix A .

The final result of this lecture is the connection between SVD and pseudoinverse.

Lemma 1.9. *Let $A \in \mathbb{R}^{m \times n}$ and $A = U \Sigma V^T$ be the SVD of A as in Theorem 1.8. Then,*

$$A^\dagger = V \Sigma^{-1} U^T.$$

Exercise 1.1. Show that $\sum_{j=1}^m \phi_j(x) \frac{\partial \phi_j}{\partial x_i}(x) = a_i^T (Ax - b)$ as in the proof of Lemma 1.4.

Exercise 1.2. Prove part 2 of Proposition 1.7.

Exercise 1.3. Prove Lemma 1.9.

Exercise 1.4. Consider $A = \begin{bmatrix} 1 & 0 & -2 \\ -\frac{1}{2} & 0 & 1 \end{bmatrix}$.

- (a) Compute by hand a singular value decomposition $U \Sigma V^T$ and the pseudoinverse A^\dagger of A .
- (b) Now try to do the same using the LinearAlgebra library in Julia [BEKS17] (or any other numerical linear algebra implementation). Do you get what you expected? What happens if you compare the pseudoinverse obtained via the command `pinv` to the one obtained by taking $V \Sigma^{-1} U^T$?

1.2 Probability Theory

Using probability theory we can model uncertainty and randomness in data. The basic idea is to assign to an event A a probability $P(A) \in [0, 1]$. It measures how likely it is that A happens.

There are two main interpretations of $P(A)$.

1. The first interpretation is that $P(A)$ should be approximately equal to the relative frequency of the event A happening in n experiments. That is, $P(A) \approx \frac{k}{n}$, where k is the number of times A happened in n experiments. Furthermore, as $n \rightarrow \infty$ the \approx should become an equality. This point of view is called **frequentist probability**.
2. The second interpretation is that $P(A)$ is a value based on experience or knowledge inferred from data. In particular, this means that unlike in the frequentist's view $P(A)$ is not independent of the observed data and can be updated when new data is available. Furthermore, we can model incomplete information about deterministic processes. This point of view is called **Bayesian probability**.

For data analysis Bayesian probability is more relevant. However, both points are only interpretations of the abstract mathematical definitions in probability! We discuss the theory in this lecture. For more details see, for instance, the (freely available) textbook [Ash70].

Definition 1.10. Let Ω be a nonempty set and $\mathcal{A} \subset 2^\Omega$ be a subset of the power set of Ω . We call \mathcal{A} a **σ -algebra**, if it satisfies the following properties

1. $\Omega \in \mathcal{A}$;
2. if $A \in \mathcal{A}$, then $\Omega \setminus A \in \mathcal{A}$;
3. if $A_n \in \mathcal{A}, n \in \mathbb{N}$, then $\bigcup_{n \in \mathbb{N}} A_n \in \mathcal{A}$.

Definition 1.11. A **probability space** is a triple (Ω, \mathcal{A}, P) , where

1. Ω is a nonempty set,
2. $\mathcal{A} \subset 2^\Omega$ is a σ -algebra, and
3. $P : \mathcal{A} \rightarrow [0, 1]$ is a probability measure. This means that

$$P(\Omega) = 1 \quad \text{and} \quad P\left(\bigcup_{n \in \mathbb{N}} A_n\right) = \sum_{n \in \mathbb{N}} P(A_n), \text{ if } A_i \cap A_j = \emptyset \text{ for } i \neq j.$$

1 The Basics

Every set $A \in \mathcal{A}$ is called an **event**, Ω is called the **space of events**, and $P(A)$ is the **probability** of A . The map P is called a **(probability) distribution**.

The restriction that \mathcal{A} is a σ -algebra is crucial: without this assumption a probability might not even exist. However, if Ω is discrete or even finite we can always take $\mathcal{A} = 2^\Omega$ as σ -algebra. In the case $\Omega = \mathbb{R}$ we have the **Borel σ -algebra**. This is the smallest σ -algebra (by inclusion) that contains every interval in \mathbb{R} .

Definition 1.12. Let \mathcal{A} be the Borel σ -algebra in \mathbb{R} . We call a function $g : \mathbb{R} \rightarrow \mathbb{R}$ **measurable**, if for all $A \in \mathcal{A}$ we have $g^{-1}(A) \in \mathcal{A}$.

Example 1.13. Let $\Omega = \{0, 1\}$ and $\mathcal{A} = \{\emptyset, \{0\}, \{1\}, \Omega\} = 2^\Omega$. Suppose $P(\{1\}) = p$. Then, we have

$$P(\{0\}) = P(\Omega) - P(\{1\}) = 1 - p.$$

This probability distribution is called **Bernoulli distribution** with parameter p . It models the probability of an experiment with two outcomes.

Often Ω is complicated, but at the same time we don't want to know every information about events in Ω , just some particular pieces of information. This motivates the definition of random variables.

Definition 1.14. A **random variable** X is a map $X : (\Omega', \mathcal{A}', P') \rightarrow (\Omega, \mathcal{A}, P)$ between probability spaces, such that for all events $A \in \mathcal{A}$ it holds that

$$X^{-1}(A) \in \mathcal{A}' \quad \text{and} \quad P(A) = P'(X^{-1}(A)).$$

We also write $P(X \in A) := P'(X^{-1}(A))$ and call it the **probability distribution** of X .

If $\Omega = \mathbb{R}$ and \mathcal{A} is the Borel σ -algebra, we call X a **continuous real random variable**. If $\Omega \subset \mathbb{R}$ is discrete and $\mathcal{A} = 2^\Omega$, we call X a **discrete real random variable**.

The definition of a random variable X is rather technical. What is it good for? The definition of a probability space in Definition 1.11 introduces the probability measure of *sets* in Ω . By contrast, one should think of a random variables as *random elements* in Ω . Often Ω is \mathbb{R} or \mathbb{R}^n so that a random variable $X \in \mathbb{R}$ represents a random real number and $X \in \mathbb{R}^n$ is a random real vector.

Example 1.15. Suppose that Ω is the set of all coin tosses. Let $X : \Omega \rightarrow \{0, 1\}$ be a random variable with $P(X = 0) = P(X = 1) = \frac{1}{2}$. Then, $P(X = 0)$ can be interpreted as the probability that the coin lands on heads, and $P(X = 1)$ as the probability that the coin lands on tails.

1 The Basics

Given a continuous real random variable $X \in \mathbb{R}$ every measurable function $g : \mathbb{R} \rightarrow \mathbb{R}$ induces another random variable $Y := g(X)$ with $P(Y \in A) := P(X \in g^{-1}(A))$.

In the following, we fix a probability space (Ω, \mathcal{A}, P) . Let $A, B \in \mathcal{A}$. We want to describe the probability of A in the situation when we already know that B has happened. This probability is denoted $P(A | B)$. It is reasonable to require $P(A | B)$ to be proportional to $P(A \cap B)$ and $P(B | B) = 1$. This motivates the following definition.

Definition 1.16. Assume $P(B) > 0$. The **conditional probability** of A given B is

$$P(A | B) := \frac{P(A \cap B)}{P(B)}.$$

Example 1.17. Let $\Omega = \{1, \dots, 6\}$, $A = \{2\}$ and $B = \{2, 4, 6\}$. Suppose that $P(\{k\}) = \frac{1}{6}$ for $k = 1, \dots, 6$. Then:

$$P(A) = \frac{1}{6} \quad \text{and} \quad P(A | B) = \frac{1}{3}.$$

In other words, when all the 6 numbers are equally likely, it is more likely to draw number 2, if we know that only even numbers will be drawn.

Theorem 1.18 (Bayes' theorem). Let $A, B \in \mathcal{A}$ with $P(A), P(B) > 0$. Then,

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

Proof. By Definition 1.16 we have $P(A | B) = \frac{P(A \cap B)}{P(B)}$ and $P(B | A) = \frac{P(A \cap B)}{P(A)}$. This implies $P(A | B)P(B) = P(B | A)P(A)$, from which the statement follows. \square

There is an interesting consequence of Bayes' theorem. Namely, $P(A | B) > P(A)$, if and only if $P(B | A) > P(B)$. In other words, B makes A more likely, if and only if A makes B more likely.

Definition 1.19. Let $A, B \in \mathcal{A}$ with $P(B) \neq 0$. We call A and B **independent**, if

$$P(A | B) = P(A).$$

We call two continuous (resp. discrete) real random variables X and Y **independent**, if

$$P(X \in A \text{ and } Y \in A) = P(X \in A)P(Y \in A)$$

for all events $A \in \mathcal{A}$. We say that a sequence of continuous (resp. discrete) real random variables $(X_n)_{n \in \mathbb{N}}$ are **independent and identically distributed** (abbreviated as **i.i.d.**), if the X_n are pairwise independent and all have the same probability distribution.

1 The Basics

Let now $X \in \mathbb{R}$ be a real random variable. If X is discrete and $X(\Omega) = \{x_1, x_2, \dots\}$ is the range of discrete values that X can admit, its probability distribution P is completely determined by the values $P(X = x_i)$ for $i = 1, \dots, n$. If X is continuous, the probability distribution of X is not so easy to describe. In many cases, however, the probability distribution can be given by a so-called probability density.

Definition 1.20. Let $X \in \mathbb{R}$ be a continuous real random variable. An integrable function $f : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$ is called a **probability density** of X , if for all events A we have

$$P(X \in A) = \int_A f(x) dx.$$

In particular, $\int_{\mathbb{R}} f(x) dx = 1$. If $X = (X_1, \dots, X_n) \in \mathbb{R}^n$ is a continuous real random vector, we call a function $f : \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$ with

$$P(X \in A) = \int_A f(x) dx$$

the **joint density** of the X_i , or simply the probability density of X .

Again recall all random variables have a probability distribution, but not all random variables have densities.

The interpretation of a probability density is that $f(x)$ measures the “infinitesimal probability” of $x \in \mathbb{R}$. We will often denote the probability density by $P_X(x) := f(x)$ or $P(x) := f(x)$. The only time this becomes confusing is when we have the singleton $A = \{x_0\}$, in which case the probability of any single event occurring for a continuous random variable is always zero:

$$P(\{x_0\}) = \int_{\{x_0\}} P(x) dx = \int_{\{x_0\}} f(x) dx = 0,$$

while $P(x)$ does not need to be zero.

Suppose now that $X \in \mathbb{R}^n$ is a continuous random variable with a density. While the probability that $X = x$ for a single point $x_0 \in \mathbb{R}^n$ is zero, we can still express the conditional probability distribution given $X = x$.

Definition 1.21. Let $(X, Y) \in \mathbb{R}^n \times \mathbb{R}^m$ be a random variable with a probability density $P_{(X,Y)}$ and $x \in \mathbb{R}^n$. The **conditional density** of Y given $X = x$ is

$$P_{Y|X=x}(y) = \frac{P_{(X,Y)}(x, y)}{P_X(x)}.$$

We write $Y | X = x$ for the random variable with this density.

1 The Basics

To see that the right-hand side in Definition 1.21 is indeed a density observe that

$$P_X(x) = \int_{\mathbb{R}^m} P_{(X,Y)}(x,y) dy, \quad (1.2.1)$$

since $P(X \in A) = P((X,Y) \in A \times \mathbb{R})$. Here, P_X is called the **marginal density**.

Theorem 1.22 (Bayes' theorem for densities). *Let $(X,Y) \in \mathbb{R}^n \times \mathbb{R}^m$ be a random variable with a probability density $P_{(X,Y)}$ and $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^m$. Then:*

$$P_{Y|X=x}(y) = P_{X|Y=y}(x) \cdot \frac{P_Y(y)}{P_X(x)}.$$

Proof. This follows immediately from Definition 1.21. □

Next, we introduce several important properties of real random variables.

Definition 1.23. Let $X \in \{x_1, x_2, \dots\}$ be a discrete real random variable. The **expected value** of X is

$$\mathbb{E}X := \sum_{i=1}^{\infty} x_i \cdot P(X = x_i).$$

If $X \in \mathbb{R}$ is a continuous real random variable with a density P its **expected value** is

$$\mathbb{E}X := \int_{\mathbb{R}} x \cdot P(x) dx.$$

In both cases, the **variance** is defined as

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

The **standard deviation** is $s(X) := \sqrt{\text{Var}(X)}$. Let X and Y be two continuous (resp. discrete) real random variables. The **covariance** of X and Y is

$$\text{Cov}(X,Y) := \mathbb{E}(X - \mathbb{E}X)(Y - \mathbb{E}Y).$$

In particular, $\text{Var}(X) = \text{Cov}(X,X)$.

Lemma 1.24 (Linearity of the expected value). *Let X and Y be two real random variables with finite expected values: $\mathbb{E}X, \mathbb{E}Y < \infty$. Then, for all $a, b \in \mathbb{R}$ we have*

$$\mathbb{E}(aX + bY) = a\mathbb{E}X + b\mathbb{E}Y.$$

Proof. See, e.g., [Ash70, Section 3.3]. See also Exercise 1.5. □

1 The Basics

Linearity of the expected value implies

$$\text{Var } X = \mathbb{E}X^2 - (\mathbb{E}X)^2 \quad \text{and} \quad \text{Cov}(X, Y) = \mathbb{E}XY - \mathbb{E}X\mathbb{E}Y.$$

Lemma 1.25. *Let $X \in \mathbb{R}^n$ be a random variable and $g : \mathbb{R}^n \rightarrow \mathbb{R}$.*

1. *If $X \in \{x_1, x_2, \dots\}$ is discrete, then*

$$\mathbb{E}g(X) = \sum_{i=1}^{\infty} g(x_i) \cdot P(X = x_i).$$

2. *If X is continuous with density P and g is measurable, then*

$$\mathbb{E}g(X) = \int_{\mathbb{R}^n} g(x) \cdot P(x) \, dx,$$

provided $\int_{\mathbb{R}^n} |g(x)| \cdot P(x) \, dx < \infty$.

Proof. We denote the random variable $Z := g(X)$. In the discrete case we set $z_i = g(x_i)$.

Then, $P(Z = z_i) = \sum_{k: g(x_k)=z_i} P(X = x_k)$ and therefore

$$\mathbb{E}Z = \sum_{i=1}^{\infty} z_i \cdot P(Z = z_i) = \sum_{i=1}^{\infty} \sum_{k: g(x_k)=z_i} g(x_k) P(X = x_k) = \sum_{k=1}^{\infty} g(x_k) \cdot P(X = x_k).$$

The continuous case requires some ideas from measure theory, which we skip here. We refer to [Ash70, Section 3, Theorem 2] for a proof. \square

Lemma 1.25 implies the following expressions for covariance of random variables $(X, Y) \in \mathbb{R}^2$ with joint density P :

$$\text{Cov}(X, Y) = \int_{\mathbb{R}^2} xy P(x, y) \, d(x, y) - \mathbb{E}X\mathbb{E}Y.$$

Lemma 1.26. *Let $X \in \mathbb{R}^n$ and $Y \in \mathbb{R}^m$ be random variables, and suppose that Y has a density $P_Y(y)$ and that $(X | Y)$ has a density $P_{X|Y=y}(x)$. Then,*

$$\mathbb{E}_X X = \mathbb{E}_Y \mathbb{E}_{X|Y=y} X.$$

Proof. By Eq. (1.2.1), the density of X is given by $P_X(x) = \int_{\mathbb{R}^m} P_{X|Y=y}(x) P_Y(y) \, dy$. This implies

$$\begin{aligned} \mathbb{E}_X X &= \int_{\mathbb{R}^n} x \cdot P_X(x) \, dx \\ &= \int_{\mathbb{R}^n} \int_{\mathbb{R}^m} x \cdot P_{X|Y=y}(x) \cdot P_Y(y) \, dy \, dx \\ &= \int_{\mathbb{R}^m} \left(\int_{\mathbb{R}^n} x \cdot P_{X|Y=y}(x) \, dx \right) \cdot P_Y(y) \, dy \\ &= \mathbb{E}_Y \mathbb{E}_{X|Y=y} X. \end{aligned}$$

□

Example 1.27. The following list of random variables describes important distributions.

1. **Bernoulli distribution:** $X \in \{0, 1\}$ and $P(X = 0) = p$.

We write $X \sim \text{Ber}(p)$.

2. **Binomial distribution:** $X \in \{0, \dots, n\}$ and $P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$.

We write $X \sim \text{Bin}(n, p)$.

Notice that $P(X = k) = P(\#\{i \mid Z_i = 0, 1 \leq i \leq n\} = k)$ for $Z_1, \dots, Z_n \stackrel{\text{i.i.d.}}{\sim} \text{Ber}(p)$.

3. **Discrete uniform distribution:** $X \in \{a_1, \dots, a_n\}$ and $P(X = k) = \frac{1}{n}$.

We write $X \sim \text{Unif}(\{a_1, \dots, a_n\})$.

4. **Continuous uniform distribution:** $X \in [a, b]$ and $P(A) = \int_A \frac{1}{b-a} dx$ for $A \subseteq [a, b]$.

We write $X \sim \text{Unif}([a, b])$.

5. **Normal distribution:** $X \in \mathbb{R}$ and

$$P(A) = \frac{1}{\sqrt{2\pi\sigma^2}} \int_A \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right) dx,$$

where $\sigma^2 > 0$ and $\mu \in \mathbb{R}$.

We write $X \sim N(\mu, \sigma^2)$.

6. **Multivariate normal distribution:** $X \in \mathbb{R}^n$ and

$$P(A) = \frac{1}{\sqrt{(2\pi)^n \det(\Sigma)}} \int_A \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) dx$$

for $\Sigma \in \mathbb{R}^{n \times n}$ symmetric positive definite and $\mu \in \mathbb{R}^n$.

We write $X \sim N(\mu, \Sigma)$.

We further write

$$\Phi(x \mid \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^n \det(\Sigma)}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) \quad (1.2.2)$$

for the density of X .

The next lemma explains why Σ is called **covariance matrix**.

1 The Basics

Lemma 1.28. Let $X \sim N(\mu, \sigma^2)$ for $\mu \in \mathbb{R}$ and $\sigma^2 > 0$. Then,

$$\mathbb{E}X = \mu \quad \text{and} \quad \text{Var}(X) = \sigma^2.$$

Furthermore, if $Y \sim N(\nu, \Sigma)$ with $\nu \in \mathbb{R}^n$ and Σ symmetric positive definite, we have

$$\text{Cov}(Y_i, Y_j) = \Sigma_{i,j}$$

for all $1 \leq i, j \leq n$.

Lemma 1.29. Let $X \sim N(\mu, \Sigma)$ be a Gaussian random variable in \mathbb{R}^n , and let $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$. Then,

$$AX + b \sim N(A\mu + b, A\Sigma A^T).$$

Exercise 1.5. Prove Lemma 1.24 (for the continuous case you can assume that (X, Y) has a joint density). **Hint:** Use Lemma 1.25 for $g(X, Y) = X + Y$ and that for every random variable $\mathbb{E}|X| < \infty$ if and only if $\mathbb{E}X < \infty$ (see [Ash70, Eq. (3.1.7)]).

Exercise 1.6. The element Caesium-137 has a half-life of about 30.17 years. In other words, a single atom of Caesium-137 has a 50 percent chance of surviving after 3.17 years, a 25 percent chance of surviving after 60.34 years, and so on.

- (a) Determine the probability that a single atom of Caesium-137 decays (i.e., does not survive) after a single day. How would you model the random variable X that takes the value 1 when the atom decays and 0 otherwise?
- (b) Using Julia, simulate 1000 times the behaviour of a collection C of 10^6 Caesium-137 atoms in a single day. How would you model the following random variable?

$$Y = \# \text{ atoms in } C \text{ decaying after a single day}$$

- (c) The *Poisson distribution* with parameter λ is a discrete probability distribution that is used to “model rare events”. When $Z \sim \text{Pois}(\lambda)$, one has that

$$P\{Z = k\} = \frac{\lambda^k}{k!} e^{-\lambda}.$$

Plot the Poisson distribution with $\lambda = 10^6 \cdot p$, where p is the probability computed in part (a).

1 The Basics

- (d) Compare the empirical distribution in part (b) to the theoretical distribution in (c).

Some Julia packages that might be useful: `Distributions`, `StatsPlots`.

Exercise 1.7. Let $\Omega := \{x_1, \dots, x_n\}$ and $p_1, \dots, p_n \geq 0$ with $p_1 + \dots + p_n = 1$. Prove that the following algorithm generates a random variable $X \in \Omega$ with $P(X = x_i) = p_i$:

1. define the numbers $w_k := \sum_{i=1}^k p_i$, $1 \leq k \leq n$, and $w_0 := 0$ and $w_{n+1} := 1$;
2. draw $Y \sim \text{Unif}([0, 1])$
(for instance, in Julia one can draw Y using the command `rand()`);
3. let k such that $w_{k-1} \leq Y < w_k$;
4. return x_k .

Exercise 1.8. Prove Lemma 1.28.

Exercise 1.9. Prove Lemma 1.29. **Hint:** Compute the density of $AX + b$.

2 Network Analysis

After the preliminaries we will now start the first chapter on mathematical methods in data science. Our first goal is to analyze structures of networks using spectral methods. We will mostly follow the book by Chung [Chu97], and the lecture notes by Guruswami and Kannan [GK12], and by Sauerwald and Sun [SS11]. For more context we also recommend [Chu10].

2.1 Graphs and the Laplace Matrix

In this section we follow the first chapter in [Chu97].

A network consist of a number of entities that are in relation to each other. Think of users in a social network that are connected, or airports for which there is a direct flight from one to another. The mathematical model for networks is a *graph*.

Definition 2.1. A graph $G = (V, E)$ is a pair consisting of a finite number of **vertices** given by

$$V = \{1, \dots, n\}$$

and a finite number of **edges** between pairs of vertices

$$E \subseteq \{\{i, j\} : i, j \in V, i \neq j\}.$$

When $v \in V$ and $e = \{u, v\} \in E$ for some $u \in V$ we say that u is **adjacent** to v . The **adjacency matrix** of G is

$$A(G) = (a_{ij}) \in \mathbb{R}^{n \times n}, \quad \text{where } a_{ij} = \begin{cases} 1 & \{i, j\} \in E \\ 0 & \text{otherwise.} \end{cases}$$

The adjacency matrix $A(G)$ can be understood as a **data structure** for a graph.

Given a vertex $v \in V$, the **degree of** v is the number of vertices adjacent to v denoted

$$\deg(v) := \#\{u \in V \mid \{u, v\} \in E\}.$$

2 Network Analysis

In the following, we will only consider graphs $G = (V, E)$ that have no isolated vertices; i.e., we assume

$$\deg(u) > 0, \quad \text{for all } v \in V.$$

Isolated vertices do not contribute to the network structure we want to analyze, which is why we want to ignore them. Detecting isolated graphs from the adjacency matrix $A(G)$ is straightforward, so that we can remove columns and rows corresponding to isolated vertices from $A(G)$ immediately.

Remark 2.2. The notation $\{i, j\}$ is used to denote an unordered set, so in particular $\{i, j\} = \{j, i\}$ which means we will be working with **simple** and **undirected** graphs.

Example 2.3. Consider $G = (V, E)$ for $V = \{1, 2, 3\}$ and $E = \{\{1, 2\}, \{1, 3\}\}$



The adjacency matrix of this graph is

$$A(G) = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix}.$$

The degrees are $\deg(1) = 2$ and $\deg(2) = \deg(3) = 1$.

Definition 2.4. We say $G = (V, E)$ is **complete** when $E = \{\{i, j\} : i, j \in V, i \neq j\}$.

Example 2.5. The following is a complete graph on 6 nodes.



2 Network Analysis

Definition 2.6. A graph $G = (V, E)$ is said to be **bipartite** if the vertex set V can be subdivided into two disjoint subsets V_1 and V_2 so that every edge in E has an endpoint in V_1 and the other in V_2 . If moreover every possible edge between V_1 and V_2 is present, G is a **complete bipartite** graph (note that there are several possible complete bipartite graphs on the same vertex set V).

Definition 2.7. A **walk** in G is a sequence of vertices

$$P = (v_0, v_1, \dots, v_D),$$

such that $\{v_{i-1}, v_i\} \in E$ for all $1 \leq i \leq D$. In this case, we say that P is a walk from v_0 to v_D . The edges of P are

$$E(P) = \{\{v_{i-1}, v_i\} \mid 1 \leq i \leq D\}.$$

The **length** of P is D . If $v_i \neq v_j$ for $i \neq j$ (no repeated vertices), we call P a **path**. We say that G is **connected**, if for every $v, w \in V$ there is a walk from v to w in G . A **connected component** of G is a maximal connected subgraph of G .

Lemma 2.8. Let A be the adjacency matrix of a graph $G = (V, E)$, and let $v, w \in V$. Then the number of walks from v to w of length k is given by $(A^k)_{v,w}$.

Proof. See Exercise 2.2. □

Next, we introduce the **Laplace matrix** or **Laplacian** of a graph G . We will see that its eigenvalues provide essential information about the network structure of G .

Definition 2.9. Let $G = (V, E)$ be a graph. The **Laplace Matrix** of G is

$$L(G) = (\ell_{ij}) \in \mathbb{R}^{|V| \times |V|},$$

where

$$\ell_{ij} = \begin{cases} 1 & i = j \\ \frac{-1}{\sqrt{\deg(i)\deg(j)}} & i \neq j \text{ and } \{i, j\} \in E \\ 0 & \text{otherwise} \end{cases}$$

Example 2.10. Consider the graph from Example 2.3 with $G = (V, E)$ and $V = \{1, 2, 3\}$ and $E = \{\{1, 2\}, \{1, 3\}\}$:

2 Network Analysis



Then

$$L(G) = \begin{pmatrix} 1 & -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & 1 & 0 \\ -\frac{1}{\sqrt{2}} & 0 & 1 \end{pmatrix}$$

In the following, we fix a graph $G = (V, E)$ and denote $A := A(G)$ and $L := L(G)$.

Definition 2.11. We define the following diagonal matrix

$$T = (t_{uv}) \in \mathbb{R}^{|V| \times |V|}, \quad t_{uv} = \begin{cases} \deg(u) & u = v \\ 0 & \text{otherwise.} \end{cases}$$

Remark 2.12. Another common definition of the Laplacian of a graph is $\mathcal{L} := T - A$, where T is as in Definition 2.11 and A is the adjacency matrix of G . In fact, we have $L = T^{-1/2} \mathcal{L} T^{-1/2}$ (as shown in the next lemma). Compared to \mathcal{L} our Laplacian is also called the **normalized Laplacian**. In our lecture we follow the definition in [Chu97] using L . In [Chu97, Section 1.2] Chung discusses that preferring L over \mathcal{L} can be helpful in the context of stochastic processes - a topic that we will cover later in our lectures.

Lemma 2.13. *The following holds*

$$L = \mathbf{1}_{|V|} - T^{-1/2} A T^{-1/2}.$$

Proof. For $u \in V = \{1, \dots, n\}$ let $e_u = (0, \dots, 0, 1, 0, \dots, 0)^T$ the u -th standard basis vector. We compute for $u, v \in V$, and using the fact that T is symmetric so $T = T^T$

$$\begin{aligned} (T^{-1/2} A T^{-1/2})_{uv} &= e_u^T T^{-1/2} A T^{-1/2} e_v \\ &= (T^{-1/2} e_u)^T A (T^{-1/2} e_v) \\ &= \frac{1}{\sqrt{\deg(u) \deg(v)}} e_u^T A e_v \\ &= \begin{cases} \frac{1}{\sqrt{\deg(u) \deg(v)}} & \{u, v\} \in E \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

Since $\{u, u\} \notin E$, we conclude $L = \mathbf{1}_{|V|} - T^{-1/2} A T^{-1/2}$. □

2 Network Analysis

The vector space $\mathbb{R}^{|V|}$ can be interpreted as the space of functions

$$\mathcal{F}(V) := \{f : V \rightarrow \mathbb{R}\}$$

with the correspondence given by

$$x = (x_1, \dots, x_{|V|}) \leftrightarrow f : V \rightarrow \mathbb{R}, f(i) = x_i. \quad (2.1.1)$$

Then $L = L(G) \in \mathbb{R}^{|V| \times |V|}$ induces a linear mapping $\mathcal{F}(V) \rightarrow \mathcal{F}(V), f \mapsto Lf$. In this way, L is the linear map given by

$$Lf(v) = \sum_{j=1}^n \ell_{vj} f(j)$$

for $v \in V$.

Lemma 2.14. *The map L induced by the Laplacian of a graph $G = (V, E)$ is given by*

$$Lf(u) = \frac{1}{\sqrt{\deg(u)}} \sum_{v \in V: \{u,v\} \in E} \frac{f(u)}{\sqrt{\deg(u)}} - \frac{f(v)}{\sqrt{\deg(v)}}.$$

Proof. Let us write

$$g := T^{-1/2} f.$$

By Lemma 2.13 we have for $u \in V$:

$$Lf(u) = f(u) - (T^{-1/2} A g)(u) = f(u) - \frac{1}{\sqrt{\deg(u)}} \sum_{v \in V} A_{uv} g(v).$$

Moreover,

$$\sum_{v \in V} A_{uv} g(v) = \sum_{v \in V: \{u,v\} \in E} g(v) = \frac{1}{\sqrt{\deg(v)}} \sum_{v \in V: \{u,v\} \in E} f(v).$$

This shows,

$$Lf(u) = f(u) - \sum_{v \in V: \{u,v\} \in E} \frac{f(v)}{\sqrt{\deg(u) \deg(v)}}. \quad (2.1.2)$$

We can write $\deg(v) = \sum_{u \in V: \{u,v\} \in E} 1$. Thus multiplying and dividing by $\deg(v)$ we can write

$$f(u) = \deg(u) \frac{f(u)}{\deg(u)} = \sum_{v \in V: \{u,v\} \in E} \frac{f(u)}{\deg(u)}. \quad (2.1.3)$$

2 Network Analysis

Combining with Eq. (2.1.2) we then have

$$\begin{aligned}
 Lf(u) &= f(u) - \sum_{v \in V: \{u,v\} \in E} \frac{f(v)}{\sqrt{\deg(u) \deg(v)}} && \text{(by Eq. (2.1.2))} \\
 &= \sum_{v \in V: \{u,v\} \in E} \frac{f(u)}{\deg(u)} - \frac{f(v)}{\sqrt{\deg(u) \deg(v)}} && \text{(by Eq. (2.1.3))} \\
 &= \frac{1}{\sqrt{\deg(u)}} \sum_{v \in V: \{u,v\} \in E} \frac{f(u)}{\sqrt{\deg(u)}} - \frac{f(v)}{\sqrt{\deg(v)}}
 \end{aligned}$$

□

The Laplace Matrix is real and symmetric, $L = L^T$. Thus all eigenvalues of L are real.

Definition 2.15. The eigenvalues of L ,

$$\lambda_0 \leq \dots \leq \lambda_{|V|-1},$$

are called the **spectrum** of G . We define

$$\lambda_G = \lambda_1.$$

Example 2.16. The Laplace matrix from Example 2.10 has spectrum 0, 1, 2.

The spectrum of a graph G encodes information about the structure of G as we will see in the following.

Definition 2.17. We define the following inner product on $\mathcal{F}(V)$:

$$\langle f, g \rangle := \sum_{u \in V} f(u)g(u).$$

We first investigate how L behaves relative to this inner product.

Theorem 2.18. The *Rayleigh quotient* of L for $f \in \mathcal{F}(V)$ is

$$\frac{\langle f, Lf \rangle}{\langle f, f \rangle} = \frac{1}{\sum_{u \in V} f(u)^2} \sum_{\{u,v\} \in E} \left(\frac{f(u)}{\sqrt{\deg(u)}} - \frac{f(v)}{\sqrt{\deg(v)}} \right)^2.$$

Proof. By Lemma 2.14, we have

$$\langle f, Lf \rangle = \sum_{u \in V} f(u) \cdot Lf(u) = \sum_{u \in V} \frac{f(u)}{\sqrt{\deg(u)}} \sum_{v \in V: \{u,v\} \in E} \frac{f(u)}{\sqrt{\deg(u)}} - \frac{f(v)}{\sqrt{\deg(v)}}.$$

2 Network Analysis

As above we set

$$g := T^{-1/2}f,$$

so that

$$\langle f, Lf \rangle = \sum_{u \in V} g(u) \sum_{v \in V: \{u,v\} \in E} g(u) - g(v).$$

We order the sum on the right as follows:

$$\begin{aligned} \langle f, Lf \rangle &= \frac{1}{2} \left(\sum_{u \in V} g(u) \sum_{v \in V: \{u,v\} \in E} g(u) - g(v) \right) \\ &\quad - \frac{1}{2} \left(\sum_{v \in V} g(v) \sum_{u \in V: \{u,v\} \in E} g(u) - g(v) \right) \\ &= \sum_{\{u,v\} \in E} (g(u) - g(v))^2. \end{aligned} \tag{2.1.4}$$

Passing back to f coordinates, where we have $\sqrt{\deg(u)}g(u) = f(u)$, finally yields

$$\frac{\langle f, Lf \rangle}{\langle f, f \rangle} = \frac{1}{\sum_{u \in V} f(u)^2} \sum_{\{u,v\} \in E} \left(\frac{f(u)}{\sqrt{\deg(u)}} - \frac{f(v)}{\sqrt{\deg(v)}} \right)^2$$

as claimed. □

Theorem 2.18 shows that L defines a bilinear form $(f, g) \mapsto \langle f, Lg \rangle$ that is positive semi-definite. Consequently, the spectrum of G is always nonnegative. We give a formal proof for this observation.

Corollary 2.19. *We have $\lambda_i \geq 0$ for $i = 0, \dots, |V| - 1$, and $\lambda_0 = 0$.*

Proof. Let λ be an eigenvalue of L with eigenvector $f \neq 0$. Then, by Theorem 2.18,

$$\lambda = \frac{\langle f, Lf \rangle}{\langle f, f \rangle} \geq 0.$$

Furthermore, let us consider the vector with $f(u) = \sqrt{\deg(u)}$. Then, again by Theorem 2.18, we have $\langle f, Lf \rangle = 0$, which shows that at least one eigenvalue is zero, so that $\lambda_0 = 0$. □

2 Network Analysis

The proof shows that we always have

$$f = T^{1/2}e \in \ker L, \quad (2.1.5)$$

where $e \in \mathcal{F}(V)$ is the constant one function (in the identification from Eq. (2.1.1) this is $e = (1, \dots, 1)$).

Next, we give the spectra of some example graphs.

Proposition 2.20. *Let G be a graph with $n = |V|$ vertices.*

1. *If G is the complete graph, then $\lambda_k = \frac{n}{n-1}$ for $k \geq 1$.*
2. *If G is a complete bipartite graph, then $\lambda_k = 1$ for $1 \leq k \leq |V| - 2$ and $\lambda_{|V|-1} = 2$.*
3. *If G is a path, then $\lambda_k = 1 - \cos \frac{\pi k}{n-1}$.*
4. *If G is a cycle, then $\lambda_k = 1 - \cos \frac{2\pi k}{n}$.*

Exercise 2.1. Consider the complete graph on 6 vertices from Example 2.5. Construct the adjacency matrix and the Laplace matrix for this graph. What are the adjacency matrix and the Laplace matrix for a complete graph on n vertices?

Exercise 2.2. Prove Lemma 2.8.

Exercise 2.3. Let

$$A = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix}.$$

Compute the $(1, 1)$ -entry of A^k for any $k \geq 1$ without computing the matrix power A^k explicitly.

Exercise 2.4. For the graph in Example 2.5, compute the number of paths of length 3 from vertex 1 to vertex 2.

Exercise 2.5. Prove Proposition 2.20.

2.2 The Spectrum of a Graph

In this lecture, G is a fixed graph with $n = |V|$ vertices and $L = L(G)$ is its Laplacian. Recall from the previous lecture that the spectrum of a graph $G = (V, E)$ is given by the

2 Network Analysis

eigenvalues of its Laplacian $L(G)$. We proved in Corollary 2.19 that these eigenvalues are nonnegative.

The main goal of this lecture is to prove the following theorem.

Theorem 2.21. *Let $G = (V, E)$ be a graph with $n = |V| \geq 2$, and let*

$$0 = \lambda_0 \leq \lambda_1 \leq \dots \leq \lambda_{n-1}$$

be the spectrum of G . We also denote $\lambda_G := \lambda_1$. The following holds.

1. $\lambda_0 + \lambda_1 + \dots + \lambda_{n-1} = n$.
2. $\lambda_G \leq \frac{n}{n-1} \leq \lambda_{n-1}$.
3. *If G is not complete, $\lambda_G \leq 1$. Otherwise, $\lambda_G = \frac{n}{n-1}$.*
4. $\lambda_i = 0$ and $\lambda_{i+1} > 0$, if and only if G has exactly $i + 1$ connected components.
5. We have $\lambda_{n-1} \leq 2$. Furthermore, $\lambda_{n-1} = 2$, if and only if there is a connected component in G that is bipartite.
6. *The spectrum of G is the union of the spectra of its connected components.*

Example 2.22. Before we prove this theorem, let us recall the graph from Example 2.3:



From Example 2.16 we know that the spectrum of this graph is $\lambda_0 = 0, \lambda_1 = 1, \lambda_2 = 2$.

First, G is not complete, which can also be seen from $\lambda_1 = 1$ (see Theorem 2.21 1.). We have one connected component corresponding to $\lambda_0 = 0 < \lambda_1$ (see Theorem 2.21 4.). Finally, $\lambda_2 = 2$ as G is bipartite (see Theorem 2.21 5.).

Let us now prove Theorem 2.21.

Proof of Theorem 2.21. Let $L := L(G)$ be the Laplacian of G .

The first item follows because the diagonal entries of L are all equal to 1 (see Definition 2.9), so that $\lambda_0 + \lambda_1 + \dots + \lambda_{n-1} = \text{Trace}(L) = n$. Using that $\lambda_0 = 0$ this implies

$$n = \lambda_1 + \dots + \lambda_{n-1} \geq (n-1)\lambda_G,$$

2 Network Analysis

so that $\lambda_G \leq \frac{n}{n-1}$. In the same spirit,

$$n = \lambda_1 + \cdots + \lambda_{n-1} \leq (n-1)\lambda_{n-1},$$

so that $\lambda_{n-1} \geq \frac{n}{n-1}$. This proves the second item.

For the third item we recall from Proposition 2.20 that, if G is complete, $\lambda_G = \frac{n}{n-1}$. We show that otherwise $\lambda_G \leq 1$.

Recall from Eq. (2.1.5) that $T^{1/2}e \in \ker L$. As a consequence, using the Rayleigh quotient, λ_G can be written in the following way:

$$\lambda_G = \min_{g \in \mathcal{F}(V) \setminus \{0\}: \langle g, T^{1/2}e \rangle = 0} \frac{\langle g, Lg \rangle}{\langle g, g \rangle}. \quad (2.2.1)$$

If G is not complete, there exist $u, v \in V$ with $\{u, v\} \notin E$. Let us define $f \in \mathcal{F}(V)$ with

$$f(i) = \begin{cases} \sqrt{\deg(v)}, & \text{if } i = u \\ -\sqrt{\deg(u)}, & \text{if } i = v \\ 0, & \text{else.} \end{cases}$$

The function f satisfies

$$\langle f, T^{1/2}e \rangle = \sqrt{\deg(u)\deg(v)} - \sqrt{\deg(u)\deg(v)} = 0.$$

By Theorem 2.18, we have

$$\begin{aligned} \frac{\langle f, Lf \rangle}{\langle f, f \rangle} &= \frac{1}{\sum_{a \in V} f(a)^2} \sum_{\{a, b\} \in E} \left(\frac{f(a)}{\sqrt{\deg(a)}} - \frac{f(b)}{\sqrt{\deg(b)}} \right)^2 \\ &= \frac{1}{\deg(u) + \deg(v)} \left(\sum_{i \in V: \{u, i\} \in E} \frac{\deg(v)}{\deg(u)} + \sum_{i \in V: \{v, i\} \in E} \frac{\deg(u)}{\deg(v)} \right) \\ &= \frac{1}{\deg(u) + \deg(v)} \left(\deg(u) \frac{\deg(v)}{\deg(u)} + \deg(v) \frac{\deg(u)}{\deg(v)} \right) = 1. \end{aligned}$$

This shows $\lambda_G \leq 1$.

For item 4. we first make the following observation: let $f \in \ker L$. Then $\langle f, Lf \rangle = 0$. Writing $g := T^{-1/2}f$ and recalling Eq. (2.1.4), we infer that $g(u) = g(v)$ for all edges $\{u, v\} \in E$. Let now i, j be two vertices in G and P be a path in G from i to j . Since for all edges $\{u, v\}$ in G we have that $g(u) = g(v)$, it follows that $g(i) = g(j)$ and g is constant on the given path.

2 Network Analysis

Assume now that G is connected: then, for every $i, j \in G$ we can find a path from i to j . It follows that g is a multiple of the constant one function e , and f is a multiple of $T^{1/2}e$. Consequently, 0 is a simple eigenvalue of L and $\lambda_1 > 0$. Conversely, if $\lambda_1 = 0$, then there exists a nonzero function f in the kernel of L which is not a multiple of $T^{1/2}e$. But then there must exist vertices $i, j \in G$ such that G contains no path from i to j , and hence G is not connected.

The statement for multiple connected components follows from this and item 6.

To prove item 5. note first that, for every $a, b \in \mathbb{R}$, one has that $(a - b)^2 \leq 2(a^2 + b^2)$, and equality holds if and only if $b = -a$.

Now, setting $g = T^{-1/2}f$ and using again the expression of the Rayleigh quotient in Theorem 2.18, we get that

$$\lambda_{n-1} = \max_{f \in \mathcal{F}(V) \setminus \{0\}} \frac{\langle f, Lf \rangle}{\langle f, f \rangle} = \max_{f \in \mathcal{F}(V) \setminus \{0\}} \frac{1}{\sum_{u \in V} \deg(u) g(u)^2} \sum_{\{u, v\} \in E} (g(u) - g(v))^2.$$

Combining this with the above inequality and $\deg(u) = \sum_{v \in V: \{u, v\} \in E} 1$ yields

$$\lambda_{n-1} \leq \frac{2}{\sum_{u \in V} \deg(u) g(u)^2} \sum_{\{u, v\} \in E} (g(u)^2 + g(v)^2) = 2.$$

The only inequality used in the argument was $(g(u) - g(v))^2 \leq 2(g(u)^2 + g(v)^2)$. Therefore, $\lambda_{n-1} = 2$ if and only if there is a function $g \in \mathcal{F}(V) \setminus \{0\}$ with $g(u) = -g(v)$ for all $\{u, v\} \in E$. If G has a bipartite component $H = (V', E')$, we can write $V' = V'_1 \sqcup V'_2$ and choose

$$g(u) = \begin{cases} 1, & \text{if } u \in V'_1 \\ -1, & \text{if } u \in V'_2 \\ 0, & \text{if } u \in V \setminus V' \end{cases}$$

to see that $\lambda_{n-1} = 2$. Conversely, if there exists a nonzero function $g \in \mathcal{F}(V)$ with $g(u) = -g(v)$ for all $\{u, v\} \in E$, let $H = (V', E')$ be a connected component of G on which g does not vanish. We define the subsets of vertices

$$W_1 := \{w \in W \mid g(w) > 0\} \quad \text{and} \quad W_2 := \{w \in W \mid g(w) < 0\}.$$

Then $V' = W_1 \cup W_2$ and, for every edge $\{u, v\} \in E'$, the endpoints u and v must lie in different W_i 's. Therefore, H is bipartite.

Finally, for the last item we denote the connected components of G by G_1, \dots, G_k . Let us write $G_i = (V_i, E_i)$, so that $V = \bigcup_{i=1}^k V_i$. We can reenumerate the vertices to

2 Network Analysis

have $V_i = \{n_{i-1} + 1, \dots, n_i\}$ with $0 = n_0 < n_1 < \dots < n_k = n$. Let also L_i be the Laplacian of G_i . Then, the Laplace matrix of G is a block diagonal matrix:

$$L(G) = \begin{bmatrix} L_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & L_k \end{bmatrix}.$$

This shows that the eigenvalues of $L(G)$ are given by the eigenvalues of the L_i . \square



Figure 2.1: The network from Example 2.25. One can see 7 connected components. The 5 vertices with highest degree are labelled. The labels are the hashtags the vertices correspond to.

Theorem 2.21 shows how we can obtain information about the structure of a graph by computing its spectrum. However, often networks are almost disconnected or almost bipartite rather than having exactly this property. Such a scenario is also reflected in the spectrum. We need another definition for formulating results in this direction.

2 Network Analysis

Definition 2.23. Let $G = (V, E)$ be a graph. The **volume** of G is

$$\text{vol}(G) := \sum_{v \in V} \deg(v) = 2|E|.$$

Proposition 2.24. Let $G = (V, E)$ be a graph, $n = |V|$, with λ_G, λ_{n-1} as in Theorem 2.21. Let $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$ be two subgraphs with $V = V_1 \cup V_2$, $V_1 \cap V_2 = \emptyset$ and $E_j = \{\{u, v\} \in E \mid u, v \in V_j\}, j = 1, 2$. Denote

$$\varepsilon := \frac{|\{\{u, v\} \in E \mid u \in V_1, v \in V_2\}|}{|E|}.$$

the fraction of edges between V_1 and V_2 . Then,

$$\lambda_G \leq \varepsilon \frac{2}{(\text{vol}(G_1)/|E| + \varepsilon) \cdot (\text{vol}(G_2)/|E| + \varepsilon)} \leq \lambda_{n-1}.$$

The meaning of the proposition is that, if λ_G is large, G can't be almost disconnected, and if λ_{n-1} is small, G can't be almost bipartite.



Figure 2.2: This graph shows the big component in Fig. 2.1, called G .

Let us also see how Proposition 2.24 is related to Theorem 2.18. The graph G is bipartite with components G_1 and G_2 , if and only if $\text{vol}(G_1) = \text{vol}(G_2) = 0$ and $\varepsilon = 1$. In this

2 Network Analysis

case, the bound in Proposition 2.24 becomes $\lambda_G \leq 2 \leq \lambda_{n-1}$, similar to Theorem 2.21 5. Furthermore, the components G_1 and G_2 are disconnected, if and only if $\varepsilon = 0$, in which case we have $\lambda_G = 0$.

Example 2.25. We illustrate Proposition 2.24 in an example. We generate a graph with the following data. Using the Julia package `Twitter.jl` we download the 500 most recent tweets featuring the hashtag `#DataScience`. The vertices in this graph are all hashtags used in these tweets. We add an edge between two vertices if the two corresponding hashtags appear together in at least one tweet. This gives the graph on $n = 171$ vertices that can be seen in Fig. 2.1. We have labelled the 5 vertices with highest degrees in this graph with their corresponding hashtag.

The graph in Fig. 2.1 has 7 connected components. We consider the big component and call the underlying graph G . Fig. 2.2 shows G and Fig. 2.3 shows the spectrum of G .



Figure 2.3: Spectrum of the graph from Fig. 2.2. The x-axis represents the index of the λ_i and the y-axis their numerical value. Instead of plotting discrete points we have plotted a piecewise linear curve connecting the discrete values $\lambda_0, \dots, \lambda_{n-1}$.

Fig. 2.3 shows that λ_G is small. Following Proposition 2.24 we therefore anticipate¹ that $G = (V, E)$ has two components $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$ such that there are few edges between them. The proof of Proposition 2.24 motivates the following clustering method: let f be the eigenfunction of λ_G . Then, we take

$$V_1 := \{v \in V \mid f(v) > 0\} \quad \text{and} \quad V_2 := \{v \in V \mid f(v) < 0\}.$$

¹Strictly speaking, in Proposition 2.24 we only showed an upper bound for λ_G in terms of the number of connected edges ε , so that λ_G could be small while ε is big.

2 Network Analysis

These two clusters are shown in Fig. 2.4. The vertices in V_1 are shown in red and the vertices in V_2 are blue. The two clusters indeed seem to give two relatively separated components of G .

We can use the same approach for analyzing the largest eigenvalue λ_{n-1} . Here, however, two clusters are not enough, because the eigenfunction g of λ_{n-1} defines *three* clusters: one cluster, where g is negative, one, where g is positive, and a third cluster, where g is zero. We show these three components in Fig. 2.5, where we have plotted the components corresponding to where g is zero in yellow. We see that the yellow vertices form the major part of G and that the blue and red vertices give a bipartite subgraph. The analysis of three components is the topic of Exercise 2.7.



Figure 2.4: We have partitioned the vertices of the graph in Fig. 2.2 into two classes of vertices corresponding to whether the eigenfunction of λ_G is positive or negative.

Proof of Proposition 2.24. Let us denote $m_i := \text{vol}(G_i) + \varepsilon$. Observe that for $i \neq j$:

$$\sum_{u \in V_i} \deg(u) = \sum_{u \in V_i} \left(\sum_{v \in V_i: \{u,v\} \in E} 1 + \sum_{v \in V_j: \{u,v\} \in E} 1 \right) = \text{vol}(G_i) + \varepsilon|E| = m_i.$$

This also shows $m_1 + m_2 = \text{vol}(G)$.

2 Network Analysis



Figure 2.5: We have partitioned the vertices of the graph in Fig. 2.2 into three classes of vertices corresponding to whether the eigenfunction of λ_{n-1} is positive, negative or zero. The class corresponding to zero is yellow.

Let us define the function

$$f(u) = \begin{cases} m_2 \sqrt{\deg(u)}, & \text{if } u \in V_1 \\ -m_1 \sqrt{\deg(u)}, & \text{if } u \in V_2 \end{cases}.$$

Then,

$$\langle f, T^{1/2}e \rangle = m_2 \sum_{u \in V_1} \deg(u) - m_1 \sum_{u \in V_2} \deg(u) = m_2 m_1 - m_1 m_2 = 0,$$

and, by Theorem 2.18,

$$\begin{aligned} \frac{\langle f, Lf \rangle}{\langle f, f \rangle} &= \frac{1}{\sum_{u \in V} f(u)^2} \sum_{\{u,v\} \in E} \left(\frac{f(u)}{\sqrt{\deg(u)}} - \frac{f(v)}{\sqrt{\deg(v)}} \right)^2 \\ &= \varepsilon |E| \frac{(m_1 + m_2)^2}{m_1 m_2^2 + m_1^2 m_2}. \end{aligned}$$

2 Network Analysis

On the one hand, we have $\lambda_{n-1} = \max_{g \in \mathcal{F}(V) \setminus \{0\}} \frac{\langle g, Lg \rangle}{\langle g, g \rangle}$, and on the other hand, as in Eq. (2.2.1) we have that $\lambda_G = \min_{g \in \mathcal{F}(V) \setminus \{0\}: \langle g, T^{1/2}e \rangle = 0} \frac{\langle g, Lg \rangle}{\langle g, g \rangle}$. This shows that

$$\lambda_G \leq \frac{\langle f, Lf \rangle}{\langle f, f \rangle} = \varepsilon |E| \frac{\text{vol}(G)}{m_1 m_2} \leq \lambda_{n-1}.$$

Finally, $|E| \frac{\text{vol}(G)}{m_1 m_2} = 2(\text{vol}(G_1)/|E| + \varepsilon)^{-1}(\text{vol}(G_2)/|E| + \varepsilon)^{-1}$. \square

The next two results underline that λ_G should be understood as a measure of connectivity for G .

Proposition 2.26. *Let $G = (V, E)$ be a connected graph and let $\text{diam}(G)$ be the **diameter** of G ; i.e., $\text{diam}(G)$ is the maximal length of all shortest paths between vertices $u, v \in V$. Then,*

$$\lambda_G \geq \frac{1}{\text{diam}(G) \cdot \text{vol}(G)}.$$

Proof. Let $f \in \mathcal{F}(V)$ be an eigenvector of λ_G with $\langle f, T^{1/2}e \rangle = 0$. Such an eigenvector exists by Eq. (2.1.5) (and since $L(G)$ is symmetric). Write $g := T^{-1/2}f$. Then,

$$0 = \langle f, T^{1/2}e \rangle = \langle T^{1/2}g, T^{1/2}e \rangle = \sum_{u \in V} \deg(u)g(u). \quad (2.2.2)$$

Let $v_0 \in V$ with $|g(v_0)| = \max_{v \in V} |g(v)|$. Eq. (2.2.2) implies that there exists $u_0 \in V$ with $g(u_0)g(v_0) < 0$ (i.e., they have opposite sign). If P is a shortest path from u_0 to v_0 of length $D > 0$, then

$$\frac{1}{D \cdot \text{vol}(G)} \geq \frac{1}{\text{diam}(G) \cdot \text{vol}(G)}.$$

We show that $\lambda_G \geq (D \cdot \text{vol}(G))^{-1}$. Using Theorem 2.18 we have

$$\begin{aligned} \lambda_G &= \frac{\langle f, Lf \rangle}{\langle f, f \rangle} = \frac{1}{\sum_{u \in V} \deg(u)g(u)^2} \sum_{\{u, v\} \in E} (g(u) - g(v))^2 \\ &\geq \frac{1}{\text{vol}(G)g(v_0)^2} \sum_{\{u, v\} \in E(P)} (g(u) - g(v))^2. \end{aligned}$$

Let us now denote the edges in P by $\{v_i, v_{i+1}\}$ for $i = 0, \dots, D-1$, where $v_D = u_0$. We define the following vectors

$$\begin{aligned} a &= (1, \dots, 1)^T \in \mathbb{R}^D \quad \text{and} \\ b &= (g(v_1) - g(v_0), g(v_2) - g(v_1), \dots, g(v_D) - g(v_{D-1}))^T \in \mathbb{R}^D. \end{aligned}$$

2 Network Analysis

Then, we can use the Cauchy-Schwartz inequality to deduce that

$$D \cdot \sum_{\{u,v\} \in E(P)} (g(u) - g(v))^2 = \|a\|^2 \cdot \|b\|^2 \geq (a^T b)^2 = (g(v_0) - g(u_0))^2.$$

It follows that

$$\lambda_G \geq \frac{1}{D \cdot \text{vol}(G)} \frac{(g(v_0) - g(u_0))^2}{g(v_0)^2},$$

and we have $(g(v_0) - g(u_0))^2 \geq g(v_0)^2$, because $g(u_0)g(v_0) < 0$. \square

The final result of this lecture is a classical result in combinatorics known as Kirchhoff's theorem. From Theorem 2.21 we know that, if G is connected, $L(G)$ has rank $n - 1$. Therefore, all $(n - 1) \times (n - 1)$ submatrices of $L(G)$ are invertible, and so have nonzero determinant. The next theorem shows that this determinant counts spanning trees in G . Recall that a tree is a graph that has no circles, and that a spanning tree τ is a tree-subgraph of $G = (V, E)$, such that the vertices of τ are V .

Theorem 2.27. *Let $G = (V, E)$ be connected, $L = L(G)$ be the Laplacian of G and $u \in V$. Let L_u be the submatrix of L that is obtained by removing the u -th row and u -th columns from L . Then,*

$$\det(L_u) = \frac{\# \text{spanning trees in } G}{\prod_{v \in V: v \neq u} \deg(v)}.$$

Remark 2.28. Let $\mathcal{L} = T - A$ as in Remark 2.12. Theorem 2.27 implies that $\det(\mathcal{L}_u)$ is the number of spanning trees in G for every $u \in V$.

Proof. We follow the proof in [AZ18].

Let us define the matrix $S = (s_{eu}) \in \mathbb{R}^{|V| \times |E|}$ indexed by vertices times edges with

$$s_{u\{i,j\}} = \begin{cases} 0, & \text{if } u \neq i, u \neq j \\ \frac{1}{\sqrt{\deg u}}, & \text{if } u = i < j \\ \frac{-1}{\sqrt{\deg u}}, & \text{if } u = i > j. \end{cases}$$

It follows from Theorem 2.18 that $\langle S^T f, S^T f \rangle = \langle f, Lf \rangle$, which shows that $L = SS^T$. Therefore, if S_u denotes the matrix that is obtained from $S \in \mathbb{R}^{|E| \times (|V|-1)}$ by removing the u -th row, we have $L_u = S_u S_u^T$. As before, we let $n = |V|$. The Cauchy-Binet formula [HJ92, Sec. 0.8.7] implies

$$\det(L_u) = \sum_{\substack{B \in \mathbb{R}^{(n-1) \times (n-1)}; \\ B \text{ is a submatrix of } S_u}} \det(B)^2. \quad (2.2.3)$$

2 Network Analysis

Let us consider a fixed $(n-1) \times (n-1)$ submatrix of S_u and call it B . The columns of B are labelled by $n-1$ edges in G . Let these edges be E_1, \dots, E_{n-1} , and let τ be the subgraph of G spanned by the E_i .

If τ is not a spanning tree, it must contain a circle; i.e., a walk (v_0, v_1, \dots) with $v_D = v_0$, $D \leq n-1$. After relabeling we can assume that $E_i = \{v_{i-1}, v_i\}$ for $1 \leq i \leq D$. But then we can find the following linear combination of the first D columns of B :

$$\begin{bmatrix} \frac{1}{\sqrt{\deg(v_0)}} \\ \frac{-1}{\sqrt{\deg(v_1)}} \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ \frac{-1}{\sqrt{\deg(v_1)}} \\ \frac{1}{\sqrt{\deg(v_2)}} \\ \vdots \\ 0 \\ 0 \end{bmatrix} + \dots + \begin{bmatrix} \frac{-1}{\sqrt{\deg(v_0)}} \\ 0 \\ 0 \\ \vdots \\ 0 \\ \frac{-1}{\sqrt{\deg(v_D)}} \end{bmatrix} = 0.$$

Therefore, the columns of B are linearly dependent and we have $\det(B) = 0$.

If, on the other hand, τ is a spanning tree, there must exist a vertex $v_1 \in V \setminus \{u\}$ such that v_1 has degree 1 in τ . After relabeling we can assume that $v_1 \in E_1$ and $v_1 \notin E_2$. After deleting v_1 from τ , we get a tree with $n-2$ edges, for which we can repeat this argument. Therefore, after a permutation of the rows and columns, we can bring B into an upper triangular form with diagonal given by $(\deg(v)^{\frac{1}{2}})_{v \in V \setminus \{u\}}$. This shows

$$\det(B)^2 = \frac{1}{\prod_{v \in V: v \neq u} \deg(v)}.$$

Since all spanning trees must be obtained as some choice of $n-1$ edges, this shows that the right hand side in Eq. (2.2.3) is equal to the number of spanning trees in G divided by $(\prod_{v \in V: v \neq u} \deg(v))^{-1}$. \square

Exercise 2.6. It follows from Theorem 2.21 that G has exactly k connected components, if and only if $\dim \ker L(G) = k$. Can you determine the k components from $\ker L$?

Exercise 2.7. State and prove a version of Proposition 2.24 involving three pairwise distinct components of a graph. Prove bounds for λ_G and λ_{n-1} for a tripartite graph.

Hint: Consider a function that is positive on the first components, negative on the second, and zero on the third component.

Exercise 2.8. Use `Twitter.jl` to generate a network, and analyze it using the spectral methods from this lecture.

2 Network Analysis

Exercise 2.9. Compute the number of spanning trees in a complete graph.

Exercise 2.10. Let G_1, G_2, G_3 be the following graphs on vertex set $\{1, 2, 3, 4\}$:



- (a) List all spanning trees for G_1, G_2 and G_3 . (You can either draw them or list their edges.) How many do you get?
- (b) Compute the number of spanning trees for G_1, G_2 and G_3 using Theorem 2.27.

2.3 Markov Processes in Networks

In this section we partially follow the lecture notes by Guruswami and Kannan [GK12], and by Sauerwald and Sun [SS11].

In some situations the complexity of a network makes it computationally infeasible to treat all vertices and edges of the corresponding graph at once. As an alternative one can explore the network vertex by vertex using **random walks**. In this lecture we analyze networks using a special type of random walks, namely **Markov processes**. Again, eigenvalues will play a central role.

As before, $G = (V, E)$ denotes a graph with $V = \{1, \dots, n\}$. We adopt the point of view from Eq. (2.1.1) that identifies vectors in \mathbb{R}^n with functions $V \rightarrow \mathbb{R}$.

Definition 2.29. A **Markov process** X on G is a sequence of random variables

$$X_0, X_1, \dots \in V,$$

called the **steps** of the walk, such that for all $i \geq 1$ we have:

1. $P(X_i = u \mid X_{i-1} = v, X_{i-2} = v_{i-2}, \dots, X_0 = v_0) = P(X_i = u \mid X_{i-1} = v)$;
2. $P(X_i = u \mid X_{i-1} = v) > 0$ only if $\{u, v\} \in E$ or $u = v$ (i.e., remaining at the current vertex is allowed);
3. $P(X_i = u \mid X_{i-1} = v)$ does not depend on i .

2 Network Analysis

The first item means that the probability law of the i -th step only depends on the position of the $(i - 1)$ -th step, but is independent of what happened before. The third item means that the probability law of a step does not depend on the number of steps that have passed. The second item means that we can only progress along edges in the graph or stand still. In the following, we will denote

$$P(u \mid v) := P(X_i = u \mid X_{i-1} = v) \text{ for } i \geq 1.$$

Definition 2.30. Let X be a Markov process in G . The **transition matrix** of X is

$$P = (p_{uv}) \in \mathbb{R}^{n \times n}$$

with $p_{uv} = P(u \mid v)$. The p_{uv} are called **transition probabilities**.

Example 2.31. Consider $G = (V, E)$ for $V = \{1, 2, 3\}$ and $E = \{\{1, 2\}, \{1, 3\}\}$:



Suppose the we start at the vertex 3. We have $P(2 \mid 3) = 0$, because there is no edge between 2 and 3. We can either move to 1 or stand still. This means that, if $p := P(1 \mid 3)$, then $P(3 \mid 3) = 1 - p$. Similarly, starting at 2 we can't move to 3. The transition matrix of any Markov process on G therefore has the form:

$$P = \begin{bmatrix} r & q & p \\ s & 1 - q & 0 \\ 1 - r - s & 0 & 1 - p \end{bmatrix}$$

with $0 \leq p, q, r, s$ and $p, q, r + s \leq 1$.

Let us make a simple but important observation, implicitly used in Example 2.31.

Lemma 2.32. Let $e = (1, \dots, 1) \in \mathbb{R}^n$ and let $P = (p_{uv}) \in \mathbb{R}^{n \times n}$ be the transition matrix of a Markov process on G . Then, e is an eigenvector of P^T with eigenvalue 1:

$$P^T e = e.$$

Proof. We have $(P^T e)(v) = \sum_{u \in V} p_{uv} = \sum_{u \in V} P(u \mid v) = 1$. □

2 Network Analysis

Definition 2.33. We call X a **uniform** Markov process, if its transition probabilities are

$$p_{uv} = P(u | v) = \begin{cases} \frac{1}{\deg(v)}, & \text{if } u \neq v \text{ and } \{u, v\} \in E \\ 0, & \text{else.} \end{cases}$$

That is, $P = AT^{-1}$, where A is the adjacency matrix of G and T is as in Definition 2.11.

Let us denote

$$\mathcal{F}_+(V) := \{f \in \mathcal{F}(V) \mid f(u) \geq 0 \text{ for all } u \in V\}. \quad (2.3.1)$$

A probability distribution f in V is given by a function $f \in \mathcal{F}_+(V)$ such that $\langle f, e \rangle = 1$, where as before e is the constant one function on V .

Lemma 2.34. Let X be a Markov process on G with transition matrix P . Let also $i \geq 0$ and $f_i := (P(X_i = 1), \dots, P(X_i = n))^T \in \mathcal{F}_+(V)$ be the probability distribution of the i -th step of X . Then,

$$f_{i+k} = P^k f_i$$

for all $k \geq 0$.

Proof. Let $u \in V$. Then, we have

$$\begin{aligned} (Pf_i)(u) &= \sum_{v \in V} p_{uv} \cdot f_i(v) = \sum_{v \in V} P(X_{i+1} = u \mid X_i = v) \cdot P(X_i = v) \\ &= P(X_{i+1} = u) = f_{i+1}(u). \end{aligned}$$

This shows $Pf_i = f_{i+1}$. Consequently,

$$P^k f_i = P^{k-1}(Pf_i) = P^{k-1} f_{i+1} = \dots = Pf_{i+k-1} = f_{i+k}.$$

□

Lemma 2.34 should be interpreted as follows: $(P^k)_{uv}$ is the probability that, if the Markov process starts at v , it reaches u after k steps.

Example 2.35. In Example 2.31 we take the transition matrix of the uniform process:

$$P = \begin{bmatrix} 0 & 1 & 1 \\ \frac{1}{2} & 0 & 0 \\ \frac{1}{2} & 0 & 0 \end{bmatrix}.$$

2 Network Analysis

This means that, if we start at 3, the next vertex is almost surely 1. Then, from 1 with probability $\frac{1}{2}$ we either go back to 3 or go to 2. We also have

$$P^2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix}.$$

The last column of P^2 shows that, starting from 3, after 2 steps we either are at 2 or 3, both with probability $\frac{1}{2}$.

Definition 2.36. Let X be a Markov process on G with transition matrix $P = (p_{uv})$.

1. We call a probability distribution $\pi : V \rightarrow \mathbb{R}$ **stationary** (with respect to X), if

$$P\pi = \pi.$$

2. A probability distribution $f : V \rightarrow \mathbb{R}$ is called **reversible** (with respect to X), if

$$f(v)p_{uv} = f(u)p_{vu}$$

for all $u, v \in V$.

We will spend a great part of the remainder of this lecture to give conditions under which stationary distributions exist and are unique. In the literature, a stochastic process that has a stationary distribution is also called an **ergodic** process. First, we give a sufficient condition for π being a stationary distribution of a Markov process X .

Proposition 2.37. *Let X be a Markov process on G and $\pi : V \rightarrow \mathbb{R}$ be a probability distribution reversible with respect to X . Then, π is a stationary distribution of X .*

Proof. Let P denote the transition matrix of X . For all $u \in V$ we have

$$(P\pi)(u) = \sum_{v \in V} p_{uv}\pi(v) = \sum_{v \in V} p_{vu}\pi(u) = \pi(u)(P^T e)(u) = \pi(u),$$

because $P^T e = e$ by Lemma 2.32. □

Example 2.38. The probability distribution with

$$\pi(u) = \frac{\deg(u)}{\text{vol}(G)}$$

is reversible with respect to the transition matrix of the uniform process from Definition 2.33. For instance, if we consider the transition matrix P from Example 2.35, $P\pi = \pi$ for $\pi = (\frac{1}{2}, \frac{1}{4}, \frac{1}{4})^T$.

2 Network Analysis

Definition 2.39. Let X be a Markov process on G with transition matrix P .

1. X is called **aperiodic**, if

$$\gcd(\{k \in \mathbb{N} \mid (P^k)_{uu} > 0\}) = 1$$

for all $u \in V$.

2. X is called **irreducible**, if for all $u, v \in V$ there exists $k \in \mathbb{N}$ with

$$(P^k)_{uv} > 0.$$

Recall from Lemma 2.34 that $(P^k)_{uv}$ is the probability of being at u after k steps given that we have started from v . Therefore, if on G there exists an irreducible Markov process, G must be connected. On the other hand, if G is connected, this does not imply that any Markov process on G is irreducible: on any graph we can always define the process which does not move with probability one.

We now come to the main theorem of this lecture.

Theorem 2.40. *Let X be an aperiodic and irreducible Markov process on G . Let P be the transition matrix of X . Then:*

1. X has a unique stationary distribution π .
2. $\lim_{k \rightarrow \infty} P^k = \pi e^T$, where $e = (1, \dots, 1)^T \in \mathbb{R}^n$.
3. For all probability distributions $f : V \rightarrow \mathbb{R}$ we have $\lim_{k \rightarrow \infty} P^k f = \pi$.

We need two auxiliary results for the proof of Theorem 2.40. We prove them first, and then we prove Theorem 2.40 towards the end of this section.

Proposition 2.41 (The Perron-Frobenius Theorem). *Let $A = (a_{ij}) \in \mathbb{R}^{n \times n}$ be a matrix with $a_{ij} > 0$ for all i, j and $A^T e = e$. Then, 1 is a simple eigenvalue of A and all other eigenvalues λ satisfy $|\lambda| < 1$.*

Proof. Since $A^T e = e$, the transpose of A has eigenvalue 1, and so A has eigenvalue 1. Fix now $k \geq 1$ and let $M := A^k$. Let us write the entries of M as m_{ij} . Since $M^T e = e$ and since the m_{ij} are all positive, we must have

$$\max_{1 \leq i, j \leq n} \|m_{ij}\| < 1. \tag{2.3.2}$$

2 Network Analysis

If in the Jordan normal form of A we have a Jordan block of the form

$$B = \begin{bmatrix} \lambda & 1 & & \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ & & & \lambda \end{bmatrix},$$

then B^k has $k\lambda^{k-1}$ on the off-diagonal. Since by Eq. (2.3.2) the entries of M must be bounded, there can be such a Jordan block only if $|\lambda| < 1$. This implies that 1 is a simple eigenvalue.

Let now $\lambda \in \mathbb{C}, \lambda \neq 1$, be another eigenvalue of A^T . Let $x = (x_1, \dots, x_n)^T \in \mathbb{C}^n$ be a corresponding eigenvector. Then, we have $\lambda x_i = a_{1i}x_1 + \dots + a_{ni}x_n$ for all $1 \leq i \leq n$. We have $1 = a_{1i} + \dots + a_{ni}$, because $A^T e = e$. Since the a_{ji} are positive, this implies that λx_i is a convex combination of the entries of x . Let us denote by $C \subset \mathbb{C}$ the convex hull of the x_i . Since x must be linearly independent of e , there must be indices i and j with $x_i \neq x_j$. Therefore, C is not a single point. Moreover, since the a_{1i} are all strictly positive, λx_i lies in the relative interior of C . From this we get

$$|\lambda x_i| < \max_{1 \leq j \leq n} |x_j|,$$

which implies $|\lambda| < 1$. □

The second result we need is a lemma from commutative algebra. For this, we recall that a **semigroup** in \mathbb{N} is a subset $S \subset \mathbb{N}$, such that for all $s, r \in S$ also $r + s \in S$.

Lemma 2.42. *Let $S \subset \mathbb{N}$ be semigroup and suppose that $\gcd(S) = 1$ (i.e., the greatest common divisor of all elements in S is 1). Then, $\mathbb{N} \setminus S$ is finite.*

Proof. We refer to [RGS09]. □

We can now prove Theorem 2.40.

Proof of Theorem 2.40. The proof is based on the fact that under the hypothesis of the theorem there exists $M \in \mathbb{N}$ with $(P^m)_{uv} > 0$ for all $m \geq M$ and $u, v \in V$. We show this at the end of this proof. Let us first see how the statement of the theorem follows.

We know that $P^T e = e$ from Lemma 2.32, and so $(P^m)^T e = e$ for all $m \geq 1$. Therefore, if $(P^m)_{uv} > 0$ for all $u, v \in V$, Proposition 2.41 implies that P^m has 1 as simple eigenvalue and that all other eigenvalues λ of P^m satisfy $|\lambda| < 1$. Since μ is an eigenvalue of P , if

2 Network Analysis

and only if $\lambda = \mu^m$ is an eigenvalue of P^m , we see that P has 1 as simple eigenvalue and that all other eigenvalues of P have absolute value strictly less than 1. It follows that P has a unique right-eigenvector $\pi \in \mathcal{F}(V)$ with eigenvalue 1. This is not yet enough to prove that π is a stationary distribution. For this, we have to show $\pi \in \mathcal{F}_+(V)$ (i.e., π is nonnegative). We do this next.

Since $\lim_{k \rightarrow \infty} \mu^k = 0$ for all eigenvalues $\mu \neq 1$ of P , we see that P^k converges to a matrix P^* of rank one. Let us write $P^* = xy^T$. Since $P^k \pi = \pi$ and $(P^k)^T e = e$ for every k , we also have $P^* \pi = \pi$ and $(P^*)^T e = e$. This shows

$$\langle y, \pi \rangle x = \pi \quad \text{and} \quad \langle x, e \rangle y = e.$$

First, this shows that $\langle y, \pi \rangle \neq 0$ and $\langle x, e \rangle \neq 0$. After rescaling y , we can assume that $\langle x, e \rangle = 1$, so that $y = e$. Then, $\pi = \langle y, \pi \rangle x = \langle e, \pi \rangle x$. Furthermore, after rescaling π , we can assume that

$$\langle e, \pi \rangle = 1. \tag{2.3.3}$$

We get

$$\lim_{k \rightarrow \infty} P^k = \pi e^T. \tag{2.3.4}$$

On the one hand, this proves the second item of Theorem 2.40. On the other hand, since π is the limit of the columns of P^k , we also have $\pi \in \mathcal{F}_+(V)$. Together with Eq. (2.3.3) this shows that π indeed gives a stationary distribution. Moreover, Eq. (2.3.4) shows that for all probability distributions $f : V \rightarrow \mathbb{R}$ we have

$$\lim_{k \rightarrow \infty} P^k f = P^* f = \pi e^T f = \pi,$$

because $e^T f = \langle e, f \rangle = 1$.

It remains to show the claim above. For this we pick $u \in V$ and define

$$S_u := \{t \in \mathbb{N} \mid (P^t)_{uu} > 0\}.$$

Recall from Lemma 2.34 that $(P^t)_{uu}$ is the probability of having a walk with t steps that starts and ends at u . If $\alpha, \beta \in S_u$, then we must have $\alpha + \beta \in S_u$, since we can simply join two walks starting and ending at u . This shows that S_u is a semigroup. Since X is aperiodic, $\mathbb{N} \setminus S_u$ is finite by Lemma 2.42. So, there exists $M_u \in \mathbb{N}$ with $(P^m)_{uu} > 0$ for all $m \geq M_u$. We set

$$M' := \max_{u \in V} M_u.$$

2 Network Analysis

Let now $u, v \in V$, $u \neq v$, and define for $t \geq 1$:

$$\ell_{uv}^t := P(X_t = u \text{ and } X_i \neq u \text{ for } i < t \mid X_0 = v). \quad (2.3.5)$$

This is the probability of arriving at u for the first time after t steps, when we have started at v . For every $k \geq 1$ we have

$$(P^k)_{uv} = \sum_{t=1}^k \ell_{uv}^t \cdot (P^{k-t})_{uu}.$$

Since X is irreducible, there exists $r \in \mathbb{N}$ such that $(P^r)_{uv} > 0$. Therefore, $\ell_{uv}^1, \dots, \ell_{uv}^r$ can't all be equal to zero, so that $\sum_{t=1}^r \ell_{uv}^t > 0$. We set

$$M := M' + r$$

and get for $m \geq M$:

$$(P^m)_{uv} = \sum_{t=1}^m \ell_{uv}^t \cdot (P^{m-t})_{uu} \geq \sum_{t=1}^r \ell_{uv}^t \cdot (P^{m-t})_{uu} > \delta \sum_{t=1}^r \ell_{uv}^t > 0,$$

where $\delta := \min_{m-R \leq k < m} (P^k)_{uu} > 0$. □

We have shown in Theorem 2.40 that aperiodic and irreducible Markov chains on a graph have a unique stationary distribution. Let us now turn our point of view upside down and start with probability distribution $\pi : V \rightarrow \mathbb{R}$ on a connected graph $G = (V, E)$. Can we find a Markov process whose stationary distribution is π ? The answer is yes! And we can use this process to sample from π by simply starting at a vertex $v \in V$ and following the random walk along the graph. This procedure is known as **Metropolis-Hastings algorithm** and it is based on the next theorem.

Theorem 2.43. *Let $G = (V, E)$ be a connected graph with $n = |V|$ and let $\pi : V \rightarrow \mathbb{R}$ be a probability distribution with $\pi(u) > 0$ for all $u \in V$. Let $d > \max_{v \in V} \deg(v)$. Define $P = (p_{uv}) \in \mathbb{R}^{n \times n}$ with*

$$p_{uv} = \begin{cases} \frac{1}{d} \min\{1, \frac{\pi(u)}{\pi(v)}\}, & \text{if } u \neq v \text{ and } \{u, v\} \in E \\ 0, & \text{if } u \neq v \text{ and } \{u, v\} \notin E \\ 1 - \sum_{i \neq v} p_{iv}, & \text{if } u = v. \end{cases}$$

Then, P is the transition matrix of an aperiodic and irreducible Markov process with unique stationary distribution π .

2 Network Analysis

The proof of the theorem is Exercise 2.11.

One interesting property of the transition matrix in Theorem 2.43 is that it only involves ratios of the entries of π . This is useful, when we only know π up to scaling. We give an example of such a situation.

Example 2.44. Let $G = (V, E)$ be a graph. We call a function $F: V \rightarrow \{1, \dots, k\}$ a k -**coloring** of G . Let us call a coloring **admissible**, if $F(u) \neq F(v)$ for all $\{u, v\} \in E$ (the minimal k , such that G has an admissible k -coloring is called the chromatic number of G). Counting the number of admissible k -colorings in G is a hard problem. Nevertheless, we can still sample uniformly a random k -coloring from the set of all admissible k -colorings in G using Theorem 2.43. For this, we assume $k \geq \max_{v \in V} \deg(v) + 2$, and define the graph $\hat{G} := (\hat{V}, \hat{E})$ with

$$\begin{aligned}\hat{V} &:= \{F: V \rightarrow \{1, \dots, k\} \mid F \text{ is admissible}\}, \\ \hat{E} &:= \{\{F, G\} \subset \hat{V} \mid F \text{ and } G \text{ differ in exactly one vertex } u\}.\end{aligned}$$

Then, \hat{G} is connected (see, e.g., [FV07]) and the number of all admissible k -colorings is $|\hat{V}|$. Using Theorem 2.43 we can sample the probability distribution $\pi: \hat{V} \rightarrow \mathbb{R}$ with

$$\pi(u) = \frac{1}{|\hat{V}|} \text{ for all } u \in \hat{V}$$

without knowing $|\hat{V}|$.

Lastly, one thing that we did not discuss in this lecture is the number of steps it takes such that an aperiodic and irreducible Markov process X is somehow close to its stationary distribution π . One way to measure convergence is using the **total variation distance** $d(t) := \max_{u, v \in V} |P(X_t = u \mid X_0 = v) - \pi(u)|$. For a given $\delta > 0$ the minimal t such that $d(t) < \delta$ is called **mixing time** of the process. We refer to [Chu97, Chapter 1] for more details. In practice, however, it often suffices to take a fixed number of steps.

Exercise 2.11. Prove Theorem 2.43. **Hint:** Use Proposition 2.37.

Exercise 2.12. Let G be a graph and X be the uniform Markov process on G . Show:

1. X is aperiodic, if and only if G is not bipartite.
2. X is irreducible, if and only if G is connected.

2 Network Analysis

Exercise 2.13. Prove that the following algorithm implements the Markov process from Example 2.44. Suppose that at the i -th step we have the admissible coloring F_i . Then we do the following:

1. Choose $(u, c) \sim \text{Unif}(V \times \{1, \dots, k\})$.
2. For all $v \in V \setminus \{u\}$ set $F_{i+1}(v) = F_i(v)$, and $F_{i+1}(u) = c$.
3. If F_{i+1} is not an admissible coloring, go back to 1. Otherwise, return F_{i+1} .

Exercise 2.14. Implement the algorithm from Exercise 2.13 for the graph



using 4 colors.

2.4 Centrality Measures

Let $G = (V, E)$ be a graph. In Section 2.2 we interpreted the values $f(u)$, $u \in V$, of a function $f \in \mathcal{F}(V)$ as an indicator of belonging to a class of vertices. For instance, we assigned to a vertex $u \in V$ one of two labels depending on whether $f(u) \geq 0$ or $f(u) < 0$. In Section 2.3 we considered $f \in \mathcal{F}_+(V)$ to be a probability distribution on the set of vertices (recall from Eq. (2.3.1) that we denote by $\mathcal{F}_+(V)$ the vector space of nonnegative functions $V \rightarrow \mathbb{R}$).

In this last section on networks we take a third perspective. We regard a function $f \in \mathcal{F}_+(V)$ as a measure of importance and call f a **centrality measure**. Centrality measures are used to assess the role of a vertex in a network. If $f(u) > f(v)$, we interpret this as u playing a more important role for the network structure than v .

We start with one of the most relevant centrality measures: **Page-Rank** [PBMW98].

Definition 2.45. Let $G = (V, E)$ be a graph. We call a function $c_R \in \mathcal{F}_+(V)$ satisfying

$$c_R(u) = \sum_{v \in V: \{u, v\} \in E} \frac{c_R(v)}{\deg(v)}$$

a Page-Rank of G .

2 Network Analysis

If c_R is a Page-Rank of G , so is $\lambda \cdot c_R$ for every $\lambda \geq 0$. The next result gives conditions when c_R exists and is unique up to such a scaling.

Proposition 2.46. *Let $G = (V, E)$ be a graph. If G is connected and not bipartite, a Page-Rank of G exists and is unique up to scaling.*

Proof. We can write the defining equation in Definition 2.45 as

$$Pc_R = c_R,$$

where P is the transition matrix of the uniform Markov process X on G (see Definition 2.33). It follows from Exercise 2.11 that X is aperiodic and irreducible, and then Theorem 2.40 implies that X has a unique stationary process; i.e., a unique solution $c_R \in \mathcal{F}_+(V)$ with $Pc_R = c_R$ and $\langle e, c_R \rangle = 1$. \square

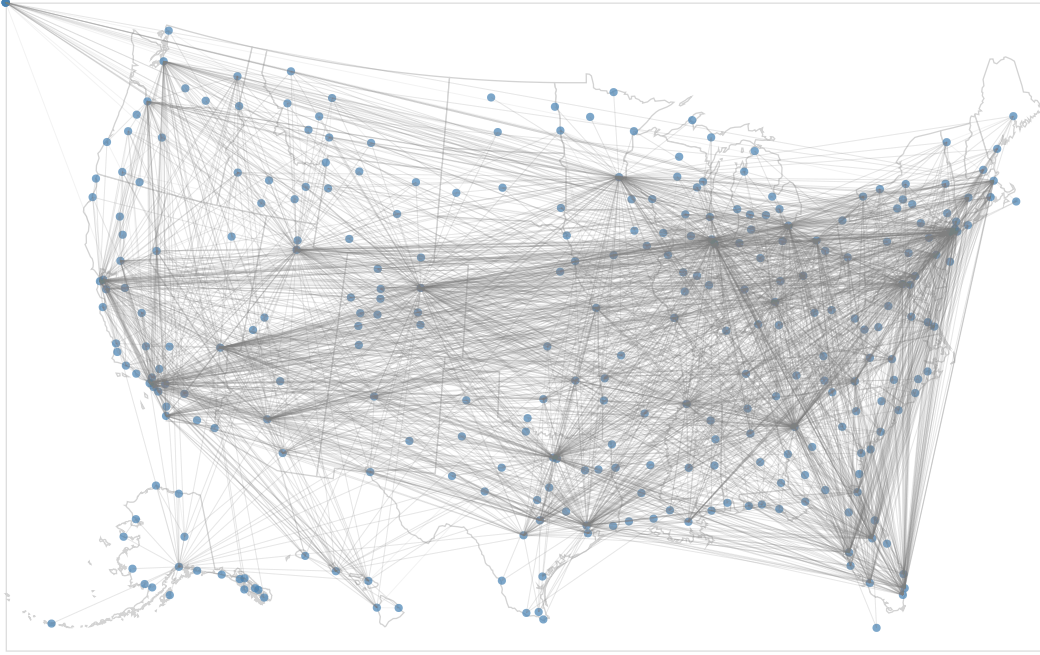


Figure 2.6: The network from Example 2.48. The blue vertices represent 305 airports in the US. There is an edge between two airports if there is a flight from one to the other. The vertex on the top left is Adak Airport.

Remark 2.47. We can use the transition matrix of any irreducible and aperiodic Markov process to define a corresponding Page-Rank.

2 Network Analysis

Example 2.48. This example is based on the Graphs lecture in the Data Science course by Huda Nassar².

We compute Page-Rank for the airport dataset from the `VegaDatasets.jl`³ package. From this dataset we compute a graph G with $n = 305$ vertices and 5668 edges. The vertices represent airports in the US, and there is an edge between two airports if there is a flight from one of the two airports to the other. The graph is shown in Fig. 2.6.



Figure 2.7: Page-Rank (blue) and approximated Page-Rank (brown) for the network from Fig. 2.6.

For computing Page-Rank we set up the transition matrix P of the uniform Markov process on G and then compute an eigendecomposition of P . We also approximate Page-Rank by using the third item in Theorem 2.40. For $1 \leq i \leq N$ we sample a random airport v_0 and then we start a random walk at v_0 using the transition probabilities in P . After m steps we record the locations v_i . We approximate Page-Rank by the empirical distribution function

$$f(u) = \frac{1}{N} |\{i \mid u = v_i\}|.$$

The result of an experiment with $N = 10^4$ and $m = 20$ is shown in Fig. 2.8 and Fig. 2.7.

Next, we introduce several other commonly used centrality measures.

²<https://github.com/JuliaAcademy/DataScience>

³<https://github.com/queryverse/VegaDatasets.jl>

2 Network Analysis



Figure 2.8: Page-Rank (blue) and approximated Page-Rank (brown) for the network from Fig. 2.6. The size of the circles corresponds to the respective centrality measures.

Definition 2.49. Let $G = (V, E)$ be a graph and $u \in V$ be a vertex.

1. The **degree centrality** of u is

$$c_D(u) := \deg(u).$$

2. The **closeness centrality** of u is

$$c_C(u) := \frac{1}{\sum_{v \in V} \text{dist}(u, v)},$$

where $\text{dist}(u, v)$ is the length of a shortest path from v to u .

3. The **harmonic centrality** of u is

$$c_H(u) := \sum_{v \in V} \frac{1}{\text{dist}(u, v)}.$$

4. Let $\sigma_{x,y}$ be the number of shortest paths from x to y , and $\sigma_{x,y}(u)$ the number of shortest paths from x to y passing through u . The **betweenness centrality** of u is

$$c_B(u) := \sum_{x,y \in V \setminus \{u\}, x \neq y} \frac{\sigma_{x,y}(u)}{\sigma_{x,y}}.$$

2 Network Analysis

5. The **Markov centrality** of u is

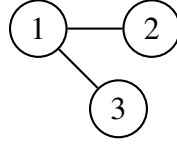
$$c_M(u) := \frac{1}{\sum_{v \in V \setminus \{u\}} \mathbb{E} \tau(u, v)},$$

where $\tau(u, v)$ is the minimal t such that a uniform Markov process X starting at v arrives at u for the first time after t steps:

$$\tau(u, v) = \min\{t \mid (X_t = u \mid X_0 = v)\};$$

i.e., $\mathbb{E} \tau(u, v) = \sum_{t=0}^{\infty} t \cdot \ell_{u,v}^t$ and $\ell_{u,v}^t = P(X_t = u \text{ and } X_i \neq u \text{ for } i < t \mid X_0 = v)$ is as in Eq. (2.3.5).

Example 2.50. Consider the graph from Example 2.3:



We compute the centrality measures from Definition 2.49 for the three vertices. In all cases 1 will have the largest measure. This can be interpreted as 1 taking the most important role in the network.

The degree centralities of the vertices are $c_D(2) = c_D(3) = 1$ and $c_D(1) = 2$. The closeness centralities are

$$C_C(1) = \frac{1}{\text{dist}(1, 2) + \text{dist}(1, 3)} = \frac{1}{1 + 1} = \frac{1}{2}$$

and

$$C_C(2) = \frac{1}{\text{dist}(1, 2) + \text{dist}(2, 3)} = \frac{1}{1 + 2} = \frac{1}{3}.$$

Similarly, $C_C(3) = \frac{1}{3}$. For harmonic centrality we have

$$C_H(1) = \frac{1}{\text{dist}(1, 2)} + \frac{1}{\text{dist}(1, 3)} = 1 + 1 = 2$$

and

$$C_H(2) = \frac{1}{\text{dist}(1, 2)} + \frac{1}{\text{dist}(2, 3)} = 1 + \frac{1}{2} = \frac{3}{2}$$

and $C_H(3) = \frac{3}{2}$ due to symmetry. The betweenness centralities are

$$C_B(1) = \frac{\sigma_{2,3}(1)}{\sigma_{2,3}} = \frac{1}{1} = 1$$

2 Network Analysis

and $C_B(2) = C_B(3) = 0$ as there are no shortest paths from 1 to 3 passing through 2, and no shortest paths from 1 to 2 passing through 3. Finally, we compute the Markov centralities of the three vertices. First, we compute $c_M(1)$:

$$c_M(1) = \frac{1}{\mathbb{E} \tau(1,2) + \mathbb{E} \tau(1,3)}.$$

Starting from either 2 or 3 the next vertex must always be 1, which means that we have $\mathbb{E} \tau(1,2) = \mathbb{E} \tau(1,3) = 1$. Consequently, $c_M(1) = \frac{1}{2}$. Next, we compute $c_M(2)$. We have

$$c_M(2) = \frac{1}{\mathbb{E} \tau(2,1) + \mathbb{E} \tau(2,3)}.$$

We have $\mathbb{E} \tau(2,1) = \sum_{t=0}^{\infty} t \cdot \ell_{2,1}^t$. First, we observe that we can only pass from 1 to 2 in an odd number of steps. This implies $\ell_{2,1}^{2k} = 0$ for all k . If the number of steps is $2k+1$, then this means we have moved from 1 to 3 k times, before moving from 1 to 2. The probability of this event is $\ell_{2,1}^{2k+1} = \frac{1}{2^{k+1}}$. Therefore,

$$\mathbb{E} \tau(2,1) = \sum_{k=0}^{\infty} \frac{2k+1}{2^{k+1}} = 3.$$

Furthermore, for moving from 3 to 2 we always need an even number of steps and we can argue as above to get

$$\mathbb{E} \tau(2,3) = \sum_{k=0}^{\infty} \frac{2k}{2^k} = 4.$$

This shows $c_M(2) = \frac{1}{3+4} = \frac{1}{7}$. Due to symmetry, also $c_M(3) = \frac{1}{7}$.

Exercise 2.15. How would you define Page-Rank for a directed graph? Analyze the data from Example 2.48 using your ideas.

Exercise 2.16. Compute the centrality measures from this section for the following graph:



3 Machine Learning

This chapter is based on [DFO20, Part II].

What is machine learning? Machine learning is a subarea of data analysis. The fundamental motivation is to develop algorithms that *automatically* extract information from datasets. Here we mean “automatic” in the sense that we want to find general methods which can be used in special situations so that we don’t need to create an algorithm for each individual scenario.

The automation of an algorithm happens through the analysis of *training data*. It is also said that one “learns” from the data. The data is understood as numeric vectors.

3.1 Data, Models, and Learning

The mathematical abstraction of a problem in machine learning is:

$$\begin{aligned} &\text{Given } (x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \mathbb{R}^N, \\ &\text{find a function } f : \mathbb{R}^D \rightarrow \mathbb{R}^N \text{ so that} \\ &\quad 1. f(x_i) \approx y_i \text{ for } 1 \leq i \leq n \\ &\quad 2. \text{ for each new data point } (x, y), f(x) \approx y. \end{aligned} \tag{3.1.1}$$

In real-world data we always have to assume the presence of noise, so that modeling $f(x) = y$ is not realistic. Using instead approximations allows for more flexibility. The exact meaning of \approx depends on the problem, and is usually measured with a **loss function**. We will come back to this in Definition 3.5 below.

Definition 3.1. In Eq. (3.1.1):

1. We call the x_i **input data**, or **attributes**.
2. The y_i are called **labels**, **output-variables**, or **response-variables**.

3 Machine Learning

Example 3.2. We are given the following data $(x, y) \in \mathbb{R}^3 \times \mathbb{R}$:

$x^{(1)} = \text{degree}$	$x^{(2)} = \text{current city}$	$x^{(3)} = \text{age}$	$y = \text{yearly income}$
MSc	Osnabrück	36	60.145 €
PhD	Osnabrück	24	72.541 €
BSc	Hannover	31	58.901 €
MSc	Bremen	29	61.005 €

This dataset does not consist of numerical values. We can convert the data into numerical values as follows (though other ways are also possible):

$x^{(1)} = \text{degree}$	$x^{(2)} = \text{latitude}$	$x^{(3)} = \text{longitude}$	$x^{(4)} = \text{age}$	$y = \text{yearly income}$
2	52,28	8,05	36	60.145
3	52,28	8,05	24	72.541
1	52,38	9,73	31	58.901
2	53,1	8,8	29	61.005

Definition 3.3. Variables which have a continuous domain of values are called **continuous variables**. Variables with a discrete domain of values are called **discrete variables** or **categorical variables**.

In Eq. (3.1.1) it is usually not helpful to consider the class of all functions. The selected functions we use are called a **model**. The choice of model is dependent on context, and can be constructed either from an analyst or automatically.

Definition 3.4. Consider given data $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \mathbb{R}^N$.

1. A **deterministic model** is a function $f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^N$ depending on $\theta \in \mathbb{R}^P$.
2. A **statistical model** is a conditional probability distribution for $(y | x) \in \mathbb{R}^N$ that depends on a parameter $\theta \in \mathbb{R}^P$.

In the following, we will always consider statistical models with a density that we denote by $P_\theta(y | x)$. A key assumption for statistical models is that the n pairs in the training data are **chosen independently**.

The goal of machine learning is to determine a parameter θ that accurately describes the data in the model. We also say that we **learn** the parameter θ . Ideally, we have chosen a model and learned a parameter that predicts well on unseen data. For this reason, we split the data into **training data** and **test data**. The training data is used

3 Machine Learning

to learn the parameter. The role of the test data is to simulate unseen data. We assess the quality of prediction of our model by evaluating a **quality function** on the test data. This last step is called **validation**.

Thus, at the core approaching a machine learning problem consists of four steps.

Algorithm 3.1: Core steps of solving a machine learning problem.

- 1 Select a model;
 - 2 Split the data into training and test data;
 - 3 Learn parameters;
 - 4 Validation.
-

Let us first consider the validation step. The common definition of a quality functions that can be used for validation is **empirical risk**.

Definition 3.5. Given data $(x_1, y_1), \dots, (x_n, y_n)$ and the model $f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^N$ (resp. P_θ), let $\ell : \mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{R}$ be a function called the **loss function**. Then the empirical risk depending on ℓ is

$$R(\theta) := \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_\theta(x_i)) \quad \text{resp.} \quad \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\hat{y}_i \sim P_\theta(y|x_i)} \ell(y_i, \hat{y}_i).$$

When there is a discrepancy in the value of the quality function between training and test data, we speak of **overfitting**. Overfitting means that the model fits the training data well, but does not accurately predict the test data.

The simplest way for model selection is to split the data by randomly assigning the data points to either training or test data, learn a parameter θ , and validate θ on the test data using a quality function $Q(\theta)$. The model for which Q is minimized is chosen. A more sophisticated way is **cross-validation**. Here, we randomly split the data into k parts $D_1 \cup \dots \cup D_k$ with $D_i \cap D_j \neq \emptyset$ for $i \neq j$. Then, for every $1 \leq i \leq k$ we learn a parameter θ_i using $\bigcup_{j \neq i} D_j$ as training data and return $Q := \frac{1}{k} \sum_{i=1}^k R(\theta_i)$, where $R(\theta_i)$ is the risk (or any other quality function) for the test data D_i . As before, we choose the model for which Q is minimized. Therefore, cross-validation is a method for model selection, but not for parameter learning (although the process for cross-validation involves the computation of parameters).

For splitting the data into training and test data we can proceed as before and randomly assign training or test labels. Usually between 50%–80% of the data is used for

3 Machine Learning

training. The random choice in this step, however, should be independent of the random choices in the model selection step.

After a model is chosen and data is prepared, we learn parameters. In the deterministic model we utilize a technique called **Empirical Risk Minimization (ERM)**. This means computing a parameter θ^* which minimizes the empirical risk $R(\theta)$ on the training data, namely $\theta^* \in \operatorname{argmin}_{\theta} R(\theta)$. For the statistical model we can use **Maximum-Likelihood Estimation (MLE)** or **Maximum a-Posteriori Estimation (MAP)**. They correspond to maximizing the following functions.

Definition 3.6. Let $(x_1, y_1), \dots, (x_n, y_n)$ be data points and P_{θ} be a statistical model.

1. The **likelihood function** is the probability of observing the response variables given the input data: $L(\theta) := \prod_{i=1}^n P_{\theta}(y_i | x_i)$. The **log-likelihood function** is

$$l(\theta) := \sum_{i=1}^n \log P_{\theta}(y_i | x_i).$$

2. Let us denote $X := \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix}^T \in \mathbb{R}^{n \times D}$ and $Y := \begin{bmatrix} y_1 & \dots & y_n \end{bmatrix}^T \in \mathbb{R}^{n \times N}$; i.e., X and Y have input data and response variables as their rows. Suppose that the parameter θ is a random variable, and that $(\theta | X, Y)$ has a probability density P . The **posteriori function** is

$$\alpha(\theta) = P(\theta | X, Y).$$

The motivation for the definition of l is that in order to maximize L , we can also maximize l . The latter is often simpler.

In particular, in maximum a-posteriori estimation we model the parameter θ as random. Taking the point of view of Bayesian probability allows us to use prior information for modelling the random variable θ . The choice of probability distribution for θ is therefore called the prior distribution. For instance, θ could be a normal distribution around a mean value that we have observed often. Going one step further, we can also keep θ random, so that our model for θ allows fluctuations. The response variable then has a conditional distribution $(y | x, \theta)$, and we can use Eq. (1.2.1) to get $P(y | x, X, Y) = \int_{\theta \in \mathbb{R}^p} P_{\theta}(y | x) \cdot P(\theta | X, Y) d\theta$. These two approaches are summarized under the name **Bayesian machine learning**.

Finally, let us briefly come back to model selection. Suppose that we have to choose among models M_1, \dots, M_r . In Bayesian machine learning we can place a prior $P(M)$

3 Machine Learning

on the choice of model. For instance, we could define $P(M_i) = \frac{1}{r}$ for $1 \leq i \leq r$. This would correspond to choosing a model uniform at random. Then, we have the posteriori function $P(M | X, Y)$, which we use for maximum a-posteriori estimation. It follows from Bayes theorem for densities (Theorem 1.22) that

$$\operatorname{argmax}_{1 \leq i \leq r} P(M_i | X, Y) = \operatorname{argmax}_{1 \leq i \leq r} P(M_i) \cdot P(X, Y | M_i)$$

and $P(X, Y | M_i)$ is the probability of having the data giving the model, also called **evidence** of the model. It is the marginal density, where the random variable θ has been integrated out: $P(X, Y | M_i) = \int_{\theta \in \mathbb{R}^D} P(X, Y | \theta) \cdot P(\theta | M_i) d\theta$, by Eq. (1.2.1). Here, $P(\theta | M_i)$ is the prior distribution of θ for the model M_i and $P(X, Y | \theta)$ the joint density of (X, Y) given θ .

3.2 Nonlinear Regression and Neural Networks

In this lecture we consider a specific machine learning model for continuous input data and continuous response variables, namely linear regression.

Definition 3.7. The **linear model** $f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}$ is given by the following function depending on the parameter $\theta = (\theta_0, \theta_1, \dots, \theta_D)^T \in \mathbb{R}^{D+1}$:

$$f_\theta(x) = x^T \theta' + \theta_0,$$

for $\theta' = (\theta_1, \dots, \theta_D)^T \in \mathbb{R}^D$. The **quadratic loss** is

$$\ell(y, \hat{y}) = (y - \hat{y})^2.$$

In this case the ERM is called the **least squares problem** or **linear regression**.

The empirical risk $R(\theta)$ (see Definition 3.5) for the quadratic loss is also called **mean squared error (MSE)**. Alternatively, one also considers the **root mean squared error (RMSE)** defined by $\operatorname{RMSE}(\theta) := \sqrt{R(\theta)}$.

We again consider training data $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \mathbb{R}$. Recall that we denote

$$X := \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix}^T \in \mathbb{R}^{n \times D} \quad \text{and} \quad Y := \begin{bmatrix} y_1 & \dots & y_n \end{bmatrix}^T \in \mathbb{R}^n. \quad (3.2.1)$$

We also denote

$$\Omega = \begin{bmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{bmatrix}^T \in \mathbb{R}^{n \times (D+1)}, \quad (3.2.2)$$

called **feature matrix**.

3 Machine Learning



Figure 3.1: The picture shows 75% of the data in the dataset `cars` from the `RDatasets` package, and the result of a linear regression for this data.

Theorem 3.8. *Giving the setting from Definition 3.7, let $Y \in \mathbb{R}^n$ be as in Eq. (3.2.1) and Ω be the feature matrix. Recall that $r(\Omega)$ denotes the rank of the matrix Ω .*

1. *When $r(\Omega) < D + 1$, then $\arg\min_{\theta} R(\theta)$ has infinitely many solutions.*
2. *When $r(\Omega) = D + 1$, then $\arg\min_{\theta} R(\theta)$ has a unique solution*

$$\theta^* := \Omega^\dagger Y,$$

where $\Omega^\dagger = (\Omega^T \Omega)^{-1} \Omega^T$ is the pseudoinverse of Ω as in Definition 1.5.

Proof. The empirical risk is

$$R(\theta) = \frac{1}{n} \sum_{i=1}^n (f_{\theta}(x_i) - y_i)^2 = \frac{1}{n} \sum_{i=1}^n (x_i^T \theta' + \theta_0 - y_i)^2 = \frac{1}{n} (\Omega \theta - Y)^T (\Omega \theta - Y).$$

Let $z := \Omega \Omega^\dagger Y$. Recall from Corollary 1.6 that $z = \arg\min_{w \in \text{Im}(\Omega)} (w - Y)^T (w - Y)$. If $r(\Omega) < D + 1$, then $\Omega \theta = z$ has infinitely many solutions. If $r(\Omega) = D + 1$, by Proposition 1.7, $\Omega^\dagger \Omega = \mathbf{1}_{D+1}$, so that $\theta = \Omega^\dagger Y$ is uniquely determined. \square

Example 3.9. The `RDatasets`¹ provides the data set `cars` that features the variables $x = \text{speed}$ and $y = \text{dist}$. The dataset is shown in Fig. 3.1. We expect a linear relation of the form $y = ax + b$ to hold. We estimate a and b using Theorem 3.8 and a subset of 75% of the data for training. The result is shown in Fig. 3.1.

¹<https://github.com/JuliaStats/RDatasets.jl>

3 Machine Learning

Unfortunately, the least-squares parameter often tends to overfitting; i.e., it works well on the training data but not on the test data. The reason is that the quadratic loss function allows too much flexibility for the parameter. To prevent overfitting we can incorporate a **regularizing term**.

Definition 3.10. Let $A \in \mathbb{R}^{(D+1) \times (D+1)}$ be an invertible matrix, and let $\lambda \in \mathbb{R}$. The **regularized loss** for the linear model is given by

$$\ell(\hat{y}, y) = (y - \hat{y})^2 + \lambda \|A\theta\|^2.$$

The special case of Definition 3.10 given by the case $A = \mathbf{1}_{D+1}$ is called **Tikhonov regularization**. The corresponding regression problem is called **ridge regression**. The parameter λ is here not considered a parameter of the model in the sense of step 3. of Algorithm 3.1, but a parameter which is chosen in advance or tuned later. Such parameters are called **hyperparameters**.

Theorem 3.11. Let $\Omega \in \mathbb{R}^{n \times (D+1)}$ be the feature matrix from Eq. (3.2.2) and $Y \in \mathbb{R}^n$ be as in Eq. (3.2.1). Giving the setting from Definition 3.10 and Eq. (3.2.3), for almost every λ there is a unique solution $\theta^* = \operatorname{argmin}_{\theta} R(\theta)$ given by

$$\theta^* = (\Omega^T \Omega + n\lambda A^T A)^{-1} \Omega^T Y.$$

Notice that Theorem 3.11 shows that the regularized linear models not only helps against overfitting, but also in the case that Ω does not have rank $D+1$. This is especially useful in the case when there is too little data with $n \leq D$ or when the data lives in a lower dimensional subspace.

Proof of Theorem 3.11. The risk of the regularized loss is given by

$$R(\theta) = \frac{1}{n} \|\Omega\theta - Y\|^2 + \lambda \|A\theta\|^2 = \frac{1}{n} (\Omega\theta - Y)^T (\Omega\theta - Y) + \lambda \theta^T A^T A \theta. \quad (3.2.3)$$

We compute the derivative of $R(\theta)$ using Exercise 1.1 and set it equal to zero:

$$\nabla R(\theta) = \frac{2}{n} \Omega^T (\Omega\theta - Y) + 2\lambda A^T A \theta = 0.$$

Which implies

$$(\Omega^T \Omega + n\lambda A^T A)\theta = \Omega^T Y.$$

3 Machine Learning

Notice that $\det(\Omega^T \Omega + n\lambda A^T A)$ is a polynomial in λ of degree $D+1$ with at most $D+1$ real zeroes. So for almost every λ the matrix $\Omega^T \Omega + n\lambda A^T A$ is invertible and thus

$$\theta = (\Omega^T \Omega + n\lambda A^T A)^{-1} \Omega^T Y,$$

which also must be the minimum. \square

Let us now consider regression problems using statistical models. Let $f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}$ be a function that depends on a parameter θ . We handle the joint probability density

$$P_\theta(y | x) = \Phi(y | f_\theta(x), \sigma^2),$$

where Φ is the density of the normal distribution as in Eq. (1.2.2). This means that we have $y = f_\theta(x) + \varepsilon$ for $\varepsilon \sim N(0, \sigma^2)$. Note that through this statistical model we model (measurement) errors in the data. The variance σ^2 is not considered a parameter in this model, but it is a hyperparameter. We will consider a model that has σ^2 as a parameter in Proposition 3.14 below.

We now again consider linear regression for $\theta = (\theta_0, \theta_1, \dots, \theta_D)^T \in \mathbb{R}^{D+1}$. That is $f_\theta(x) = x^T \theta' + \theta_0$, where $\theta' = (\theta_1, \dots, \theta_D)$.

Theorem 3.12. *Let $\Omega \in \mathbb{R}^{n \times (D+1)}$ be the feature matrix as defined in Eq. (3.2.2). Let $Y \in \mathbb{R}^n$ be as in Eq. (3.2.1). The maximum likelihood estimation in the above model for linear regression*

1. *is not unique if $r(\Omega) < D+1$,*
2. *is uniquely determined by $\theta_{\text{ML}} = \Omega^\dagger Y$ when $r(\Omega) = D+1$.*

Note that Theorem 3.12 shows that the risk minimization in the deterministic model and in the statistical model give the same answer.

Proof of Theorem 3.12. Following Definition 3.6 the log-likelihood-function is

$$\begin{aligned} l(\theta) &= \sum_{i=1}^n \log P_\theta(y_i | x_i) \\ &= -\frac{n}{2} \log(2\sigma^2) + \sum_{i=1}^n \frac{1}{2\sigma^2} (x_i^T \theta' - \theta_0 - y_i)^2 \\ &= -\frac{n}{2} \log(2\sigma^2) + \frac{1}{2\sigma^2} \|\Omega \theta - Y\|^2. \end{aligned}$$

3 Machine Learning

The first term does not involve θ . As in the proof of Theorem 3.11, we proved that $\operatorname{argmin}_{\theta} \|\Omega\theta - Y\|^2$ has a unique solution exactly when $r(\Omega) = D + 1$, and then the solution is $\theta_{\text{ML}} = \Omega^\dagger Y$. \square

The generalization from linear regression is **nonlinear regression**, which is defined through the following statistical model for $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^P$ and $\theta = (\theta_1, \dots, \theta_P) \in \mathbb{R}^P$:

$$P_{\theta}(y | x) = \Phi(y | \theta^T \phi(x), \sigma^2).$$

For every such model we have a corresponding **feature matrix**, which we denote by

$$\Omega = \begin{bmatrix} \phi(x_1) & \dots & \phi(x_n) \end{bmatrix}^T \in \mathbb{R}^{n \times P}. \quad (3.2.4)$$

This notation is not in conflict with Eq. (3.2.2): in Eq. (3.2.2) we see the special case of Eq. (3.2.4) for linear regression. We note the following.

- In linear regression we have $P = D + 1$ and

$$\phi(x) = \begin{bmatrix} 1 \\ x \end{bmatrix}.$$

- For $D = 1$ and

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ \vdots \\ x^{P-1} \end{bmatrix}$$

we have **polynomial regression**.

The maximum likelihood estimate for nonlinear regression is given next.

Theorem 3.13. *Let Ω be the feature matrix in Eq. (3.2.4) and $Y \in \mathbb{R}^n$ be as in Eq. (3.2.1). The maximum likelihood estimator in nonlinear regression is*

1. *not unique when $r(\Omega) < P$.*
2. *uniquely determined by $\theta_{\text{ML}} = \Omega^\dagger Y$ when $r(\Omega) = P$.*

Proof. See Exercise 3.1 \square

In Theorem 3.12 and Theorem 3.13 we have the variance σ^2 provided. Alternatively we can also model the parameter.

3 Machine Learning

Proposition 3.14. *If we model the variance σ^2 in Theorem 3.13 instead as a parameter, then this gives the MLE $(\theta_{\text{ML}}, \sigma_{\text{ML}}^2)$ with $\theta_{\text{ML}} = \Omega^\dagger Y$, and $\sigma_{\text{ML}}^2 = \frac{1}{n} \|Y - \Omega \theta_{\text{ML}}\|^2$.*

Proof. See Exercise 3.3 □

Similar to linear regression the MLE tends towards **overfitting**. Our solution in the deterministic model was to introduce a penalty term $\lambda \|A\theta\|^2$. In the Bayesian approach θ itself is random, and instead of a penalty term we prescribe a choice of θ called a **prior**. This then leads to a maximum a-posteriori estimator as in Definition 3.6.

The next theorem computes the MAP for the statistical model $(y | x) \sim N(\theta^T \phi(x), \sigma^2)$ using a Gaussian prior $N(\mu, \Sigma)$. The theorem shows that in this case MAP can be understood as the statistical analogue of regularization.

Theorem 3.15. *Let $\mu \in \mathbb{R}^P$ and $\Sigma \in \mathbb{R}^{P \times P}$ be positive. Given the statistical model above, then for the prior $\theta \sim N(\mu, \Sigma)$ we have the MAP*

$$\theta_{\text{MAP}} = (\Omega^T \Omega + \sigma^2 \Sigma^{-1})^{-1} (\Omega^T Y + \sigma^2 \Sigma^{-1} \mu),$$

where Ω is the feature matrix, $Y \in \mathbb{R}^n$ is as in Eq. (3.2.1), and under the assumption that $\Omega^T \Omega + \sigma^2 \Sigma^{-1}$ is invertible.

Proof. Let $\alpha(\theta) = P(\theta | X, Y)$ be defined as in Definition 3.6. By Bayes' Theorem for densities (Theorem 1.22),

$$P(\theta | X, Y) = P(\theta) \frac{P(Y | X, \theta)}{P(Y | X)}.$$

Therefore setting $c = -\log P(Y | X)$ which is constant in θ , we have

$$\log P(\theta | X, Y) = \log P(\theta) + \log P(Y | X, \theta) + c.$$

Recall that

$$P(\theta) = \Phi(\theta | \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^n \det \Sigma}} \exp \left(-\frac{1}{2} (\theta - \mu)^T \Sigma^{-1} (\theta - \mu) \right). \quad (3.2.5)$$

Furthermore, as the n pairs in the training data are assumed to be independent, we have

$$P(Y | X, \theta) = \prod_{i=1}^n \Phi(y_i | \theta^T \phi(x_i), \sigma^2) = \frac{1}{\sqrt{(2\sigma^2)^n}} \exp \left(-\frac{1}{2\sigma^2} \|Y - \Omega \theta\|^2 \right). \quad (3.2.6)$$

3 Machine Learning

Eq. (3.2.5) and Eq. (3.2.6) together yield

$$\log P(\theta | X, Y) = -\frac{1}{2}(\theta - \mu)^T \Sigma^{-1}(\theta - \mu) - \frac{1}{2\sigma^2} \|Y - \Omega\theta\|^2 + c', \quad (3.2.7)$$

where $c' = c + \log \frac{1}{\sqrt{(2\pi)^n \det \Sigma}} + \log \frac{1}{\sqrt{2\sigma^2}^n}$ is independent of θ .

The logarithm is strictly monotone, so we can instead of maximizing $\alpha(\theta)$ also maximize $\log \alpha(\theta) = \log P(\theta | X, Y)$. Taking the derivatives with respect to θ of Eq. (3.2.7):

$$\frac{d}{d\theta} \log \alpha(\theta) = -\Sigma^{-1}(\theta - \mu) + \frac{1}{\sigma^2} \Omega^T (Y - \Omega\theta).$$

Thus, if we set this derivative equal to zero

$$\Sigma^{-1}\theta - \Sigma^{-1}\mu = \frac{1}{\sigma^2} \Omega^T Y - \frac{1}{\sigma^2} \Omega^T \Omega\theta,$$

which implies

$$(\Omega^T \Omega + \sigma^2 \Sigma^{-1})\theta = \frac{1}{\sigma^2} \Omega^T Y + \sigma^2 \Sigma^{-1} \mu.$$

So we conclude

$$\theta = (\Omega^T \Omega + \sigma^2 \Sigma^{-1})^{-1} (\Omega^T Y + \sigma^2 \Sigma^{-1} \mu)$$

and this must also be the minimizer θ_{MAP} . □

Example 3.16. We compute the ML and MAP estimator using Theorem 3.12 and Theorem 3.15 for the data in Example 3.9. A sample from these statistical models is shown in Fig. 3.2. The sample points are connected by lines to plot a piecewise linear function.

In the discussion after Definition 3.6 we observed that instead of finding θ deterministically through MLE or MAP, we can also compute the distribution of θ given the training data. In the context of regression this approach is called **Bayesian regression**. We will assume as before that $\theta \sim N(\mu, \Sigma)$ and $(y | x, \theta) \sim N(\phi(x)^T \theta, \sigma^2)$. The goal of Bayesian regression is to sample from the posteriori distribution with density $\alpha(\theta) = P(\theta | X, Y)$. In this concrete case we can explicitly compute the posteriori distribution.

Theorem 3.17. *In the setting above we have*

$$(\theta | X, Y) \sim N(m, S),$$

where $S = (\sigma^{-2} \Omega^T \Omega + \Sigma^{-1})^{-1}$ and $m = S(\sigma^{-2} \Omega^T Y + \Sigma^{-1} \mu)$.

3 Machine Learning

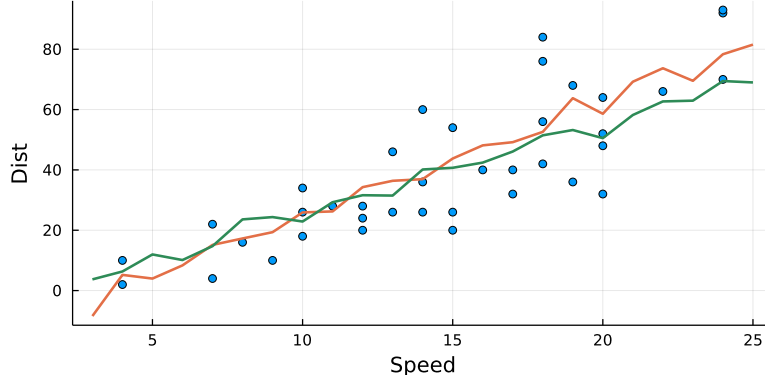


Figure 3.2: The picture shows ML and MAP for the data in Fig. 3.1. For the plot we have sampled points and connected them by lines to get a piecewise linear function illustrating the statistical model. The red line shows the ML estimator, the green line the MAP estimator.

In fact, Theorem 3.17 implies Theorem 3.15, because the density of a Gaussian is maximized at the expected value; i.e., $m = \theta_{\text{MAP}}$.

Proof of Theorem 3.17. Recall from Eq. (3.2.7) that

$$\log P(\theta | X, Y) = -\frac{1}{2}(\theta - \mu)^T \Sigma^{-1}(\theta - \mu) - \frac{1}{2\sigma^2} \|Y - \Omega\theta\|^2 + c,$$

where c is some constant, which is independent of θ . Let us write

$$Q := (\theta - \mu)^T \Sigma^{-1}(\theta - \mu) + \frac{1}{\sigma^2}(Y - \Omega\theta)^T(Y - \Omega\theta),$$

so that $P(\theta | X, Y) = \exp(-\frac{1}{2}Q + c')$. We expand to find

$$Q = \theta^T \Sigma^{-1} \theta - 2\mu^T \Sigma^{-1} \theta + \frac{1}{\sigma^2}(\theta^T \Omega^T \Omega \theta - 2Y^T \Omega \theta) + c',$$

where c' is independent of θ . Setting

$$A = \frac{1}{\sigma^2} \Omega^T \Omega + \Sigma^{-1} \quad \text{and} \quad a = \frac{1}{\sigma^2} \Omega^T Y + \Sigma^{-1} \mu$$

gives

$$Q = \theta^T A \theta - 2a^T \theta + c' = (\theta - b)^T A (\theta - b) + c'',$$

where $Ab = a$ and c'' is independent of θ . Notice

$$b = A^{-1}a = \left(\frac{1}{\sigma^2} \Phi^T \Phi + \Sigma^{-1} \right)^{-1} \left(\frac{1}{\sigma^2} \Phi^T Y + \Sigma^{-1} \mu \right).$$

3 Machine Learning

This shows that the density function of θ given X, Y is

$$P(\theta | X, Y) = \exp\left(-\frac{1}{2}Q + c'\right) = C \exp\left(-\frac{1}{2}(\theta - m)^T S^{-1}(\theta - m)\right)$$

with C independent from $\theta \in \mathbb{R}^P$. Since $P(\theta | X, Y)$ is a density, it must integrate to 1 and we must therefore have $C = 1/\sqrt{(2\pi)^P \det(S)}$. \square

For the validation step we compute the marginal distribution of the response variable $y \in \mathbb{R}$ given an input variable $x \in \mathbb{R}^D$ and the training data $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \mathbb{R}$.

Proposition 3.18. *Giving training data $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \mathbb{R}$ and an additional data point $(x, y) \in \mathbb{R}^D \times \mathbb{R}$, the distribution of $(y | x, X, Y)$ is*

$$(y | x, X, Y) \sim N(\phi^T(x)m, \phi^T(x)S\phi(x) + \sigma^2)$$

where S and m are as in Theorem 3.17, and X, Y are as in Eq. (3.2.1).

Proof. See Exercise 3.4. \square

The discussion in this section centered around models with (non-)linear functions $\mathbb{R}^D \rightarrow \mathbb{R}$. Given a nonlinear function $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^P$, a straightforward generalization of this setting is to multivariate models

$$f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^N, x \mapsto \begin{bmatrix} \phi(x)^T \theta^{(1)} \\ \vdots \\ \phi(x)^T \theta^{(N)} \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta^{(1)} & \dots & \theta^{(N)} \end{bmatrix} \in \mathbb{R}^{P \times N}, \quad (3.2.8)$$

and using quadratic loss $\ell(y, \hat{y}) = \|y - \hat{y}\|^2$ for the deterministic model or the multivariate Gaussian distribution $P_\theta(y | x) = \Phi(y | f_\theta(x), \sigma^2 \mathbf{1}_N)$ for the statistical model. Since the Euclidean norm satisfies $\|y - \hat{y}\|^2 = \sum_{i=1}^N (y_i - \hat{y}_i)^2$ for vectors $y = (y_1, \dots, y_N)^T$ and $\hat{y} = (\hat{y}_1, \dots, \hat{y}_N)^T$, maximization with respect to parameters can be done for each entry of $f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^N$ separately. Therefore, the estimators obtained in this section can be used for each entry of $f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^N$ separately.

Another generalization of linear regression models is given by neural networks. They are obtained by iterating nonlinear regression models.

Definition 3.19. Let $L > 0$ and $N_0, N_1, \dots, N_L > 0$. Let $g_{\theta_i} : \mathbb{R}^{N_{i-1}} \rightarrow \mathbb{R}^{N_i}$ be nonlinear regression models, $1 \leq i \leq L$. Denote $D := N_0, N := N_L$. We call The model

$$f_\theta = (g_{\theta_L} \circ \dots \circ g_{\theta_1}) : \mathbb{R}^D \rightarrow \mathbb{R}^N$$

a **neural network** of depth L . The function g_{θ_i} is called the i -th **layer** of the network.

3 Machine Learning

In this definition, the linear and nonlinear regression models above are neural networks of depth 1. For depth larger than 1, however, we have no analytic expression for MLE or MAP. Instead, we have to use methods from optimization for computing them.

Neural networks often encompass combinations of linear functions with so-called **activation functions**. This means that in Definition 3.19 the i -th layer is given by a function of the form $g_{\theta_i} = \sigma_i(f_{\theta_i}(x))$, where f_{θ_i} is a multivariate nonlinear regression model as in Eq. (3.2.8) and σ_i is a nonlinear activation function. For instance, activation functions for $z \in \mathbb{R}^k$ are the **ReLU function** $\sigma(z) = (\max(0, z_i))_{1 \leq i \leq k}$ the **sigmoid function** $\sigma(z) = (\frac{1}{1+\exp(-z_i)})_{1 \leq i \leq k}$, or the **softmax function**

$$\sigma(z) = \left(\frac{\exp(z_i)}{\sum_{j=1}^k \exp(z_j)} \right)_{1 \leq i \leq k}. \quad (3.2.9)$$

I.e., the softmax function is a smooth version of the max function.

Example 3.20. We use a neural network of depth 2 for the data in Example 3.9. The inner functions should be $\sigma \circ f_1$ and $\sigma \circ f_2$, where σ is the ReLu activation function and $f_1 : \mathbb{R}^1 \rightarrow \mathbb{R}^2$ and $f_2 : \mathbb{R}^2 \rightarrow \mathbb{R}^1$ are linear. Unlike in the previous examples we don't have a closed form for the optimal parameters. Instead, we minimize the empirical risk using optimization methods. The result of the computation is shown in Fig. 3.3. We can see from the figure that the neural network computed an estimator that is piecewise linear. This indicates that there could be a hidden latent variable describing the data with two different linear models depending on whether speed is small or large.

In Example 3.20 we used a neural network to obtain a predictor for a continuous response variable. We can also use a neural network for statistical models of categorical variables. Suppose that we are given data $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \mathbb{R}^N$, where the response variables are categorical and take values in the finite set $\{c_1, \dots, c_k\}$. Let $f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^k$ be a neural network such that $\sum_{j=1}^k (f_\theta(z))_j = 1$ and $(f_\theta(z))_i \geq 0$ for all $1 \leq i \leq k$ and $z \in \mathbb{R}^D$. We get a **statistical model for categorical variables** $P_\theta(y | x)$ by setting

$$P_\theta((y | x) = c_i) := (f_\theta(x))_i.$$

For instance, we can choose in the last layer the soft-max from Eq. (3.2.9) as activation function. One should interpret soft-max as smooth version of the max function (hence the name). The soft-max function is often used for **classification problems**. The idea is to choose the label, which has the highest probability predicted by the model.

3 Machine Learning



Figure 3.3: The picture shows a function describing the data in Fig. 3.1 that was computed using a neural network.

Exercise 3.1. Prove Theorem 3.13. **Hint:** Adapt the proof of Theorem 3.12.

Exercise 3.2. Reformulate and prove Theorem 3.8 for linear models $\mathbb{R}^D \rightarrow \mathbb{R}^N$ and the quadratic loss $\ell(y, \hat{y}) = \|y - \hat{y}\|^2$.

Exercise 3.3. Prove Proposition 3.14.

Exercise 3.4. Prove Proposition 3.18. **Hint:** By Eq. (1.2.1) the marginal density is given by $P(y | x, X, Y) = \int_{\mathbb{R}^P} P(y | x, \theta) \cdot P(\theta | X, Y) d\theta$.

Exercise 3.5. From the RDatasets package in Julia load the pressure data set. This data set contains the variables temperature and pressure, which give the values of pressure of mercury depending on temperature. The **Antoine Equation** is a simple model for this dependency:

$$\log(\text{pressure}) = a - \frac{b}{\text{temperature}}.$$

Set up and solve a regression problem to estimate a and b .

3.3 Support Vector Machines

At the end of the last section we discussed how to use neural networks for classification problems. Another machine learning method for classification are **support vector machines (SVM)**.

3 Machine Learning

We discuss SVMs in the context of **binary classification**. In this setting the labels y_1, \dots, y_n are elements in $\{1, -1\}$ and we want to find a model that fits the data $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \{-1, 1\}$. We will work here with finding a deterministic model as in Definition 3.4 given by

$$f_\theta : \mathbb{R}^D \rightarrow \{-1, 1\}.$$

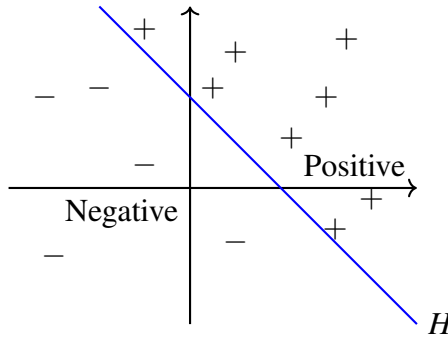
The generalization to $k \geq 2$ categories is then given by taking several of such functions, so that we map into $\{-1, 1\}^\ell$ for $2^\ell \geq k$.

SVMs provide the following family of parametrized functions:

$$f_\theta(x) = \text{sgn}(\langle a, x \rangle + b) \quad \text{for} \quad \theta = (a, b) \in \mathbb{R}^D \times \mathbb{R}, a \neq 0. \quad (3.3.1)$$

Thus, the idea behind SVMs is to find a hyperplane $H = \{x \in \mathbb{R}^D \mid \langle a, x \rangle + b = 0\}$ and split the data into on each side of H . Notice that we can define the same plane H if we take $\langle a, a \rangle = 1$, since $(\lambda a, \lambda b)$ and (a, b) define the same hyperplane.

Example 3.21. Suppose that our input data is contained in the plane \mathbb{R}^2 . Consider the line $H = \{x = (x_1, x_2) \in \mathbb{R}^2 \mid x_1 + x_2 - 2 = 0\}$ shown in blue below. Then H splits the plane in the following two regions.



We classify data in \mathbb{R}^2 by assigning them a plus or a minus sign. The data points on the upper right of the plane are labeled with a plus, because for them $x_1 + x_2 - 2 > 0$. The other data points are labeled with a minus sign, because here $x_1 + x_2 - 2 < 0$.

The plan for the rest of this lecture is now to find a suitable way computing parameters for the model Eq. (3.3.1). To do so, we first make two small observations.

Lemma 3.22. Let $x \in \mathbb{R}^D$ and $\theta = (a, b) \in \mathbb{R}^D \times \mathbb{R}$. For $y \in \{-1, 1\}$, we have $y = f_\theta(x)$ if and only if $y(\langle a, x \rangle + b) > 0$.

3 Machine Learning



Figure 3.4: Geometrical setup in the proof for Lemma 3.23.

Proof. In the case that $f_\theta(x) = -1$, we have $\langle a, x \rangle + b < 0$ and $y = -1$, which implies $y(\langle a, x \rangle + b) > 0$. Similarly when $f_\theta(x) = 1$, then $\langle a, x \rangle + b > 0$ and $y = 1$ so that in this case we also have $y(\langle a, x \rangle + b) > 0$. \square

We can see from Example 3.21 that a hyperplane H which separates the data doesn't need to be unique. SVMs select the Hyperplane which maximizes the distance to the data. The following lemma helps us to formulate the right optimization problem.

Lemma 3.23. *Let $x \in \mathbb{R}^D$ and $H = \{z \in \mathbb{R}^D \mid \langle a, z \rangle + b = 0\}$ for $\langle a, a \rangle = 1$. The Euclidean distance from x to H is $y(\langle a, x \rangle + b)$ where $y = f_\theta(x)$.*

Proof. The geometrical setup of this proof is depicted in Fig. 3.4. Let $z \in H$ be the point which minimizes the distances to x . So we can write $x = z + \varepsilon r a$ where $r = \|x - z\| \geq 0$ and $\varepsilon \in \{-1, 1\}$. Since $z \in H$ we compute

$$\langle a, x \rangle + b = \langle a, z + \varepsilon r a \rangle + b = \langle a, z \rangle + \varepsilon r \langle a, a \rangle + b = 0 + \varepsilon r.$$

Note if $\varepsilon = -1$, then $y = -1$ and therefore $y(\langle a, x \rangle + b) = y\varepsilon r = r$. Similarly if $\varepsilon = 1$ then $y = 1$ and therefore $y(\langle a, x \rangle + b) = y\varepsilon r = r$. \square

In order to compute the hyperplane that maximizes the distance to the data, we can now write it in terms of solving the following optimization problem:

$$\begin{aligned} \max_{\theta=(a,b)} \quad & r \\ \text{s.t.} \quad & y_k(\langle a, x_k \rangle + b) \geq r, \text{ for } k = 1, \dots, n, \\ & \text{and } \langle a, a \rangle = 1, r \geq 0. \end{aligned} \tag{3.3.2}$$

3 Machine Learning

The optimization problem of Eq. (3.3.2) is often solved via an equivalent problem by normalizing the value of r . Namely if $a' = \frac{a}{r}$ and $b' = \frac{b}{r}$, then the constraints now ask for $y_k(\langle a', x_k \rangle + b') \geq 1$. Then in this case $\|a'\| = \frac{1}{r}$, so we can either maximize $\frac{1}{\|a'\|}$, or minimize $\|a'\|$. This leads to the following parameter finding problem.

Definition 3.24. The **Hard Margin SVM** is given by the optimization problem

$$\begin{aligned} \min_{\theta=(a,b)} \|a\|^2 \\ \text{s.t. } y_k(\langle a, x_k \rangle + b) \geq 1, \text{ for } k = 1, \dots, n. \end{aligned}$$

An issue with Hard Margin SVM when working with noisy data is that it does not allow for outliers. To compensate for this, we will introduce a **slack variable** ξ_k .

Definition 3.25. The **Soft Margin SVM** is given by the optimization problem

$$\begin{aligned} \min_{a,b,\xi} \|a\|^2 + C \sum_{k=1}^n \xi_k \\ \text{s.t. } y_k(\langle a, x_k \rangle + b) \geq 1 - \xi_k, \text{ for } \xi_k \geq 0, k = 1, \dots, n. \end{aligned}$$

The parameter C is called the **regularization parameter**.

Here, the regularization parameter is not taken as a parameter of the model but as a hyperparameter. Moreover, soft Margin SVM can be understood as empirical risk minimization. This is proved in the next proposition.

Proposition 3.26. Consider the **Hinge loss** $\ell(y, \hat{y}) = \max\{0, 1 - y \cdot \hat{y}\}$. The Soft Margin SVM can be derived from ERM with respect to a regularized Hinge loss.

Proof. Exercise 3.6. □

The SVMs in Definition 3.24 and Definition 3.25 both fall in the category of **Primal SVMs**. Another formulation is the **Dual SVM** which we now work toward defining. We first define the Lagrange function for Soft Margin SVM:

$$\mathcal{L}(a, b, \xi, \alpha, \beta) = \|a\|^2 + C \sum_{k=1}^n \xi_k - \sum_{k=1}^n \alpha_k (y_k(\langle a, x_k \rangle + b) - (1 - \xi_k)) - \sum_{k=1}^n \beta_k \xi_k, \quad (3.3.3)$$

where $\alpha_k, \beta_k \geq 0$. The KKT conditions say that the optimum occurs when

$$\frac{\partial \mathcal{L}}{\partial a} = 0, \quad \frac{\partial \mathcal{L}}{\partial b} = 0, \quad \frac{\partial \mathcal{L}}{\partial \xi} = 0.$$

3 Machine Learning

We will write these equations by rewriting \mathcal{L} using

$$u := (\frac{1}{2}\alpha_k y_k)_{k=1}^n, \quad v := (\langle a, x_k \rangle + b)_{k=1}^n, \quad e = (1, \dots, 1)^T \in \mathbb{R}^n,$$

which gives

$$\mathcal{L} = \langle a, a \rangle + C \langle \xi, e \rangle - 2 \langle u, v \rangle - \langle \alpha + \beta, \xi \rangle + \langle \alpha, e \rangle.$$

Therefore

$$\begin{aligned} 0 = \frac{\partial \mathcal{L}}{\partial a} &= 2a - 2 \sum_{k=1}^n u_k x_k \quad \Rightarrow \quad a = \sum_{k=1}^n u_k x_k \\ 0 = \frac{\partial \mathcal{L}}{\partial b} &= -2 \langle u, e \rangle \quad \Rightarrow \quad 0 = \sum_{k=1}^n u_k \\ 0 = \frac{\partial \mathcal{L}}{\partial \xi} &= Ce - (\alpha + \beta) \quad \Rightarrow \quad Ce = \alpha + \beta \quad \Rightarrow \quad \alpha_i \leq C \end{aligned} \tag{3.3.4}$$

Remark 3.27. The equation $\frac{\partial \mathcal{L}}{\partial a} = 0$ is the meaning behind the name Support Vector Machine. Namely, the x_i with $u_k \neq 0$ and equivalently $\alpha_k \neq 0$ are the **support** of the vector a .

First, observe that from Eq. (3.3.4) it follows that.

$$\langle u, v \rangle = \sum_{k=1}^n u_k \langle a, x_k \rangle + b \sum_{k=1}^n u_k = \sum_{k=1}^n u_k \langle a, x_k \rangle = \langle \sum_{k=1}^n u_k x_k, \sum_{k=1}^n u_k x_k \rangle.$$

Now, we put the equations in Eq. (3.3.4) into \mathcal{L} and get

$$\begin{aligned} \mathcal{L} &= \langle a, a \rangle + C \langle \xi, e \rangle - 2 \langle u, v \rangle - \langle \alpha + \beta, \xi \rangle + \langle \alpha, e \rangle \\ &= \langle \sum_{k=1}^n u_k x_k, \sum_{k=1}^n u_k x_k \rangle - 2 \langle u, v \rangle + \langle \alpha, e \rangle \\ &= -u^T G u + \langle \alpha, e \rangle \quad \text{where} \quad G = (\langle x_k, x_\ell \rangle)_{k, \ell=1}^n, \text{ and } 0 \leq \alpha_i \leq C \end{aligned}$$

This all leads to the following definition of the Dual SVM.

Definition 3.28. The **Dual SVM** is given by the optimization problem

$$\begin{aligned} \max_{\alpha} \quad & -u^T G u + \sum_{k=1}^n \alpha_k \\ \text{s.t.} \quad & \sum_{k=1}^n \alpha_k y_k = 0, \text{ and } 0 \leq \alpha_i \leq C, \end{aligned}$$

where $u := (\frac{1}{2}\alpha_k y_k)_{k=1}^n$ and $G = (\langle x_k, x_\ell \rangle)_{k, \ell=1}^n$.

3 Machine Learning

Remark 3.29. *Why a maximum?* The solution for the optimization problem for the Primal SVM is given by

$$\min_{a,b,\xi} \max_{\alpha,\beta} \mathcal{L}(a,b,\xi,\alpha,\beta). \quad (3.3.5)$$

But the function in Definition 3.28 is independent of a, b, ξ , and β , so that we can remove the minimization step.

When the Dual SVM problem has been solved, we can use the following result to compute direction of the hyperplane a and the offset b .

Proposition 3.30. *Let α be an optimal solution from the Dual SVM. Then, we have optimal values for the Soft Margin SVM from Definition 3.25 by setting*

1. $a^* := \sum_{k=1}^n u_k x_k$, where $u_k = y_k \alpha_k$;
2. b^* is the median value of $y_k - \langle a^*, x_k \rangle$ for all k with $\alpha_k \neq 0$.

Proof. The formula for a^* follows from Eq. (3.3.4) Using that we have $\alpha + \beta = C$ in the optimal value, the Lagrange function from Eq. (3.3.3) becomes

$$\mathcal{L} = \|a\|^2 - \sum_{k=1}^n \alpha_k (y_k (\langle a, x_k \rangle + b) - 1).$$

By Eq. (3.3.5) we are maximizing the Lagrangian over α . Therefore,

$$y_k (\langle a, x_k \rangle + b) - 1 \leq 0 \quad \Rightarrow \quad \alpha_k = 0. \quad (3.3.6)$$

This implies

$$\mathcal{L} = \|a\|^2 - \sum_{k: \alpha_k > 0} \alpha_k (y_k (\langle a, x_k \rangle + b) - 1) = \|a\|^2 - C \sum_{k: \alpha_k > 0} (y_k (\langle a, x_k \rangle + b) - 1),$$

as the term in the middle is maximized for setting all $\alpha_k = C$. Now, by Eq. (3.3.6) the summands on the right are all nonnegative. Using that $|y_k| = 1$ for all k we get that b^* is a point on the real line that minimizes $\sum_{k: \alpha_k > 0} |y_k - \langle a^*, x_k \rangle - b|$. The median of $y_k - \langle a^*, x_k \rangle$ for $\alpha_k \neq 0$ is a minimizer (see Exercise 3.7). \square

It is not always possible to find a suitable hyperplane using Soft-Margin SVM. We give an example of a situation where this occurs, and potential solutions.

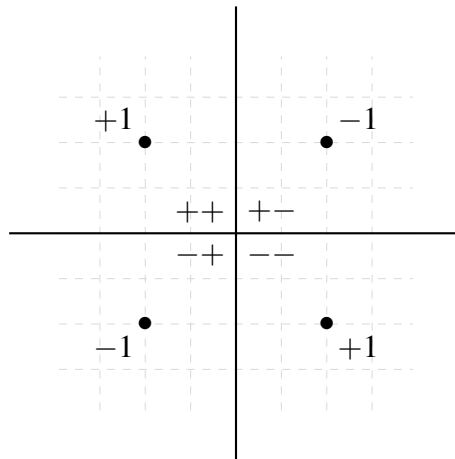
Example 3.31. There is no hyperplane separating the following 4 points into 2 classes. Even using Soft-Margin SVM the 2 classes can't be well separated.

3 Machine Learning

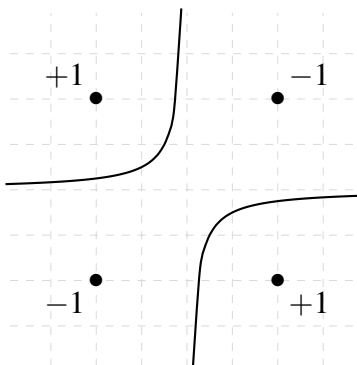


We can think of two solutions:

Solution 1: Combine multiple hyperplanes. The problem with this first solution is that we have multiple labels for each class.



Solution 2: Separate through curved hypersurfaces. In the below example we separate the data via a hyperbola given by a polynomial of degree 2.



3 Machine Learning

For solution 2 above the idea is to combine an SVM with a **feature map**

$$\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$$

(as we did in nonlinear Regression), and compute an SVM for the modified data points $(\phi(x_1), y_1), \dots, (\phi(x_n), y_n)$. In this situation the Dual SVM is well suited. In particular, in the formulation of Definition 3.28 only the inner product between input variables occurs. So it is enough to know the value of the inner products $\langle \phi(x_k), \phi(x_\ell) \rangle$. These are parametrized through **positive semi-definite matrices**.

Lemma 3.32. *Let $G \in \mathbb{R}^{n \times n}$. Then G is positive semi-definite if and only if there exists M and vectors $z_1, \dots, z_n \in \mathbb{R}^M$ with $G = (\langle z_i, z_j \rangle)_{i,j=1}^n$.*

Proof. Let $Z \in \mathbb{R}^{M \times n}$ be the matrix whose columns are the z_i . If $G = Z^T Z$, for every $w \in \mathbb{R}^n \setminus \{0\}$ we then have $w^T G w = (Zw)^T Zw \geq 0$. On the other hand, if G is positive semi-definite, we can find a Cholesky-decomposition $G = Z^T Z$. \square

In the context of SVMs, the positive definite matrices are also called **kernels**. The function

$$\kappa(x, y) := \langle \phi(x_k), \phi(x_\ell) \rangle$$

is called **kernel map**. By Proposition 3.30 we have the optimal value $a^* = \frac{1}{2} \sum_{i=1}^n u_i \phi(x_i)$. Let us define

$$\psi(x) := \frac{1}{2} \sum_{i=1}^n u_i \kappa(x_i, x).$$

Then, we have $\psi(x) = \langle a^*, \phi(x) \rangle$ by linearity. The optimal value for b is then the median of $y_k - \psi(x_k)$ for $\alpha_k \neq 0$. Moreover, we can evaluate the model from Eq. (3.3.1) as $f_\theta(x) = \text{sgn}(\psi(x) + b)$. All this leads to the following algorithm.

Algorithm 3.2: Binary classification by Dual SVM.

- 1 **Input:** Training data $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \{-1, 1\}$, a kernel map $\kappa(x, y)$, and a regularization parameter C .
 - 2 **Output:** A function $f : \mathbb{R}^D \rightarrow \{-1, 1\}$ of the form $f(x) = \text{sgn}(\langle a, \phi(x) \rangle + b)$.
 - 3 Compute the kernel matrix $G = (\kappa(x_i, x_j))_{1 \leq i, j \leq n}$;
 - 4 Using C and G solve the Dual SVM problem from Definition 3.28 for $\alpha \in \mathbb{R}^n$;
 - 5 Define the function $\psi(x) := \frac{1}{2} \sum_{i=1}^n y_i \alpha_i \kappa(x_i, x)$;
 - 6 Take b as the median of $\{y_k - \psi(x_k) \mid \alpha_k \neq 0\}$;
 - 7 Return $f(x) = \text{sgn}(\psi(x) + b)$.
-

3 Machine Learning

An important example of a kernel is the **polynomial kernel**.

Lemma 3.33. Consider the feature map $\phi(x) = (x_1^{i_1} \cdots x_D^{i_D})_{(i_1, \dots, i_D) \in I_D}$, where the index set is $I_D := \{(i_1, \dots, i_D) \in \mathbb{N}^D \mid 0 \leq i_1 + \dots + i_D \leq d\}$. I.e., $\phi(x)$ gives all monomials in x of degree at most d . Then for $u, v \in \mathbb{R}^D$

$$\langle \phi(u), \phi(v) \rangle = (\langle u, v \rangle + 1)^d.$$

Proof. We write $u_0 := 1$ and $v_0 := 1$. Then

$$\begin{aligned} \langle \phi(u), \phi(v) \rangle &= \phi(u)^T \phi(v) = \sum_{0 \leq i_0 + \dots + i_D = d} u_0^{i_0} u_1^{i_1} \cdots u_D^{i_D} v_0^{i_0} v_1^{i_1} \cdots v_D^{i_D} \\ &= \sum_{0 \leq i_0 + \dots + i_D = d} (u_0 v_0)^{i_0} (u_1 v_1)^{i_1} \cdots (u_D v_D)^{i_D} \\ &= (u_0 v_0 + u_1 v_1 + \dots + u_D v_D)^d = (\langle u, v \rangle + 1)^d. \end{aligned}$$

□

Exercise 3.6. Prove Proposition 3.26.

Exercise 3.7. Let $w_1, \dots, w_n \in \mathbb{R}$. Prove that the median of the w_i minimizes the aggregated distances $d(v) = \sum_{i=1}^n |w_i - v|$.

Exercise 3.8. The `MLDatasets.jl`² package provides the MNIST dataset³. This dataset contains pictures of handwritten digits. Load the training data for pictures of zeros and ones and implement an algorithm that learns to separate these two classes. After the learning step let your algorithm predict the labels of test data points.

3.4 Principal Component Analysis

In this lecture we consider **Principal Component Analysis (PCA)**, which is one of the principal methods for **dimensionality reduction**. This means the following:

Given data pairs $x_1, \dots, x_n \in \mathbb{R}^D$, with the help of PCA we model the data so that we reduce the number of parameters that the data describes. This can have several motivations. For instance, we could be interested in **data compression**, and reducing

²<https://github.com/JuliaML/MLDatasets.jl>

³<http://yann.lecun.com/exdb/mnist/>

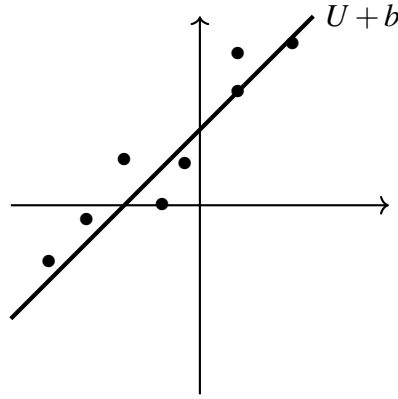
3 Machine Learning

parameters would reduce the memory for storing the data. Another motivation is to interpret the parameters as geometric information, so that the goal here would be to **learn the shape of the data**.

For PCA we will not be taking into account response variables $y_1, \dots, y_n \in \mathbb{R}^N$. Learning without using response variables is called **unsupervised learning**. By contrast, the settings from the previous sections are summarized as **supervised learning**.

The basic idea of PCA is to find a linear space $U \subseteq \mathbb{R}^D$ of dimension $d \ll D$ and a vector $b \in \mathbb{R}^D$ so that x_1, \dots, x_n lay “nearby” $U + b$.

Example 3.34. Here is a small graphic when $D = 2$ and $d = 1$. The line is represented by $U + b$, and the points are the points x_1, \dots, x_8 .



For the moment, we take d as a fixed input parameter. We will discuss later how this parameter can be chosen.

As depicted in Example 3.34, the assumption in PCA is that the data centers around a low-dimensional linear subspace, and the goal is to determine this subspace. This assumption, however, is not always fulfilled. Data points can also lie on a **nonlinear** subspace. For this reason, we again work with a nonlinear feature map $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$, and we set

$$z_i := \phi(x_i), \quad 1 \leq i \leq n.$$

For instance, if we take the polynomial feature map from Lemma 3.33, then finding a linear subspace for the z_i means finding a low-dimensional **algebraic variety** describing the data; i.e., the vanishing set of a system of polynomial equations.

3 Machine Learning

In the following, we assume that the data x_1, \dots, x_D are independent samples from an (unknown) random variable $x \in \mathbb{R}^D$, and we set $z := \phi(x)$, so that

$$z_1, \dots, z_n \stackrel{\text{i.i.d.}}{\sim} z.$$

We also denote the **expected value** $\mu := \mathbb{E}z \in \mathbb{R}^M$ and the **covariance matrix** of z :

$$\Sigma := \begin{bmatrix} \text{Cov}(z_1, z_1) & \cdots & \text{Cov}(z_1, z_D) \\ & \ddots & \\ \text{Cov}(z_D, z_1) & \cdots & \text{Cov}(z_D, z_D) \end{bmatrix} \in \mathbb{R}^{M \times M}.$$

By Exercise 3.9, Σ is positive semi-definite, and thus diagonalizable with orthogonal matrices and has eigenvalues $\lambda_1 \geq \dots \geq \lambda_M \geq 0$.

In practice we rarely know Σ exactly, because we only have a finite data sample available. Thus we need to approximate Σ through the **empirical covariance matrix**.

Definition 3.35. Let $z_1, \dots, z_n \in \mathbb{R}^M$.

1. The **empirical average** is given by the arithmetic mean $\bar{z} = \frac{1}{n}(z_1 + \dots + z_n)$.
2. The **empirical covariance matrix** is $S = (s_{ij}) \in \mathbb{R}^{M \times M}$ with

$$s_{ij} = \frac{1}{n} \sum_{k=1}^n ((z_k)_i - \bar{z}_i)((z_k)_j - \bar{z}_j).$$

Remark 3.36. A common approach is to **standardize the data**. This means to replace for every k the i -th entry $(z_k)_i$ by $(z_k)_i / \sqrt{s_{ii}}$, and then using the modified data for PCA.

As before, let us denote the feature matrix (see Eq. (3.2.4))

$$\Omega = \begin{bmatrix} \phi(x_1) & \cdots & \phi(x_n) \end{bmatrix}^T = \begin{bmatrix} z_1 & \cdots & z_n \end{bmatrix}^T \in \mathbb{R}^{n \times P}.$$

Writing also $e = (1, \dots, 1)^T \in \mathbb{R}^n$, then we can write

$$\bar{z} = \frac{1}{n} \Omega^T e, \quad \text{and} \quad S = \frac{1}{n} (\Omega^T - \bar{z} e^T) (\Omega^T - \bar{z} e^T)^T. \quad (3.4.1)$$

In particular, the second equation together with Lemma 3.32 shows that the empirical covariance S , like Σ , is also positive semi-definite.

We will consider two interpretations of “nearby”. The first is to select U as the space which maximizes the variance in the sense of the following definition.

3 Machine Learning

Definition 3.37. Let $z \in \mathbb{R}^M$ be a random variable with expectation $\mu := \mathbb{E}z \in \mathbb{R}^M$. Let $1 \leq d < D$. A **space of maximal variance** of dimension d is a linear subspace $U \subset \mathbb{R}^M$ of dimension $d = \dim(U)$ such that $U \in \operatorname{argmax}_{U: \dim(U)=d} \mathbb{E} \|P_U(z - \mu)\|^2$, where P_U is the orthogonal projection onto U .

We have the following theorem.

Theorem 3.38. Let $\lambda_1 \geq \dots \geq \lambda_M \geq 0$ be the eigenvalues of the covariance matrix $\Sigma \in \mathbb{R}^{M \times M}$, and let $u_i \in \mathbb{R}^M$ be the eigenvectors associated to λ_i so that $\langle u_i, u_j \rangle = \delta_{ij}$. Then

$$U = \operatorname{span}\{u_1, \dots, u_d\}$$

is a space of maximal variance of dimension d . Furthermore, U is uniquely determined in the case that $\lambda_d > \lambda_{d+1}$.

As mentioned above, in practice we do not have Σ available and instead use the eigenvectors of S for computing U .

Before we prove this theorem, let us discuss for a brief moment the choice of d . The statement of Theorem 3.38 implies that d should be chosen, such that the data spreads out in the direction of the eigenvectors for $\lambda_1, \dots, \lambda_d$, but not in the directions of the eigenvectors for $\lambda_{d+1}, \dots, \lambda_M$. One choice for such a d is to take $\lambda_d > 0$ and $\lambda_{d+1} \approx 0$. Alternatively, we can choose d such that λ_d/λ_{d+1} or $\lambda_d - \lambda_{d+1}$ is maximized.

Proof of Theorem 3.38. Let U be a subspace of dimension d and u_1, \dots, u_d an orthonormal basis of U (that is $\langle u_i, u_j \rangle = \delta_{ij}$). Also let $A = \begin{bmatrix} u_1 & \dots & u_d \end{bmatrix} \in \mathbb{R}^{M \times d}$. By Corollary 1.6, the orthogonal projection onto U is given by $P_U = AA^\dagger$ and, by Proposition 1.7, we have $A^\dagger = (A^T A)^{-1} A^T = A^T$. Therefore,

$$P_U = AA^T,$$

and so

$$\begin{aligned} \|P_U(z - \mu)\|^2 &= \|AA^T(z - \mu)\|^2 \\ &= (AA^T(z - \mu))^T (AA^T(z - \mu)) \\ &= (z - \mu)^T AA^T AA^T (z - \mu) \\ &= (z - \mu)^T AA^T (z - \mu); \end{aligned}$$

3 Machine Learning

the last line because $A^T A = \mathbf{1}_d$. Now, using that $A^T A = \sum_{i=1}^d u_i u_i^T$ this gives

$$\|P_U(z - \mu)\|^2 = \sum_{i=1}^d (z - \mu)^T u_i u_i^T (z - \mu) = \sum_{i=1}^d u_i^T (z - \mu)(z - \mu)^T u_i.$$

Using Exercise 3.9 and that the expected value is linear we get

$$\mathbb{E} \|P_U(z - \mu)\|^2 = \sum_{i=1}^d u_i^T \mathbb{E}[(z - \mu)(z - \mu)^T] u_i = \sum_{i=1}^d u_i^T \Sigma u_i.$$

Thus, a space of maximal variance $U = \text{span}\{u_1, \dots, u_d\}$ maximizes the scalar function $\sum_{i=1}^d u_i^T \Sigma u_i$. This function is maximized by taking the eigenvectors of the largest eigenvalues of Σ . We prove this.

We maximize with Lagrange Multipliers. Let

$$\mathcal{L}(u_1, \dots, u_d, \ell_{ij}) = \sum_{i=1}^d u_i^T \Sigma u_i - \sum_{1 \leq i \leq j \leq d} (u_i^T u_j - \delta_{ij}) \ell_{ij}.$$

We compute

$$\frac{\partial \mathcal{L}}{\partial u_j} = 2\Sigma u_j - 2u_j \ell_{jj} - \sum_{i < j} u_j \ell_{ij} = 0 \quad 1 \leq i \leq d. \quad (3.4.2)$$

$$\frac{\partial \mathcal{L}}{\partial \ell_{ij}} = u_i^T u_j - \delta_{ij} = 0 \quad 1 \leq i \leq j \leq d. \quad (3.4.3)$$

From Eq. (3.4.2) for $i = 1$ it follows that $\Sigma u_1 = \ell_{11} u_1$ thus $\ell_{11} = \lambda_1$. Similarly,

$$2u_1^T \Sigma u_2 - \ell_{12} = 0 \Rightarrow \ell_{12} = 0 \Rightarrow \Sigma u_2 = \ell_{22} u_2 \Rightarrow \ell_{22} = \lambda_2.$$

By continuing this process, the proof follows. Finally, if $\lambda_d > \lambda_{d+1}$, then U is uniquely determined as the sum of the eigenspaces for $\lambda_1, \dots, \lambda_d$. \square

Definition 3.39. The u_1, \dots, u_d in Theorem 3.38 are called **principal components** of Σ .

In the first approach we computed the principle components of S , thus obtaining $U \subset \mathbb{R}^D$. The data would then be compressed by applying the map

$$z_i \mapsto P_U(z_i - \bar{z}) + \bar{z}.$$

In the second approach, instead of maximizing the variance, we want to minimize the squared distance the data points.

$$\sum_{i=1}^n \|(z_i - \bar{z}) - P_U(z_i - \bar{z})\|^2. \quad (3.4.4)$$

3 Machine Learning

The next theorem shows that, although this approach is conceptually different from maximizing the variance, we get the same minimizer as in Theorem 3.38 (when replacing the covariance matrix Σ by the empirical covariance matrix S).

Theorem 3.40. *Let $\lambda_1, \dots, \lambda_M \geq 0$ be the eigenvalues of the empirical covariance matrix S . Let u_i be eigenvectors of λ_i so that $\langle u_i, u_j \rangle = \delta_{ij}$. Then $U = \text{span}\{u_1, \dots, u_d\}$ minimizes the squared distance in Eq. (3.4.4). If $\lambda_d > \lambda_{d+1}$, the subspace U is uniquely determined.*

Proof. Let $w_i = z_i - \bar{z}$. As in the proof of Theorem 3.38 let $A = [u_1, \dots, u_d] \in \mathbb{R}^{M \times d}$ so that $P_U = AA^T$. Then, w_i is the i -th column of $W := \Omega^T - \bar{z}e^T$. By Eq. (3.4.1), we have

$$WW^T = nS.$$

We note that $AA^T = \sum_{i=1}^d u_i u_i^T$ and $\mathbf{1}_D = \sum_{i=1}^D u_i u_i^T$. Therefore, we also have

$$W - AA^T W = \left(\sum_{i=d+1}^D u_i u_i^T \right) W.$$

Observe that

$$W - P_U W = W - AA^T W = \begin{bmatrix} w_1 - P_U(w_1) & \dots & w_n - P_U(w_n) \end{bmatrix}.$$

Therefore, Eq. (3.4.4) can be rewritten as $\text{Trace}((W - P_U W)(W - P_U W)^T)$. We get

$$\text{Trace}((W - P_U W)(W - P_U W)^T) = \sum_{i=d+1}^D \text{Trace}(u_i u_i^T W W^T) = n \sum_{i=d+1}^D u_i^T S u_i.$$

As in the proof from Theorem 3.38, one can show that the u_i must be eigenvectors of S so that $U = \text{span}\{u_1, \dots, u_d\}$ for $S u_i = \lambda_i u_i$. The uniqueness statement follows as in the proof of Theorem 3.38. \square

In both cases we need to compute eigenvectors of the covariance matrix S . In the case that $n \ll M$, so that the number of data points is significantly smaller than the dimension of the data, then we can proceed as follows. As in the proof of Theorem 3.40, let $W := \Omega^T - \bar{x}e^T$, so that the empirical covariance matrix is $S = \frac{1}{n} W W^T$ by Eq. (3.4.1). We assume noisy data, so that the rank of W can be assumed to n . Let also

$$W = U D V^T$$

3 Machine Learning

be the **singular value decomposition** (see Theorem 1.8) of W , where

$$U \in \mathbb{R}^{M \times n}, \quad V \in \mathbb{R}^{n \times n},$$

and $D = \text{diag}(\sigma_1, \dots, \sigma_n)$ with singular values $\sigma_1 \geq \dots \geq \sigma_n \geq 0$. Then, we have

$$S = \frac{1}{n} U D^2 U^T \in \mathbb{R}^{M \times M}.$$

This shows that, using PCA, we actually compute a low-dimensional representation of the data using SVD. We can exploit this: Consider the eigendecomposition

$$W^T W = V D^2 V^T \in \mathbb{R}^{n \times n}.$$

Since the i -th columns of U and V are related by $u_i = W v_i / \|W v_i\|$, this shows that it suffices to compute an eigendecomposition of an $n \times n$ -matrix, rather than a decomposition of an $M \times M$ -matrix. Moreover, $W^T W$ can be computed using only the kernel map.

Lemma 3.41. *Let $\kappa(x_1, x_2) = \langle \phi(x_1), \phi(x_2) \rangle$ be the kernel map for ϕ , and let the kernel matrix be $G = (\kappa(x_k, x_\ell))_{1 \leq k, \ell \leq n}$. Then, for $W = \Omega^T - \bar{x} e^T$ we have*

$$W^T W = (\mathbf{1}_n - \frac{1}{n} e e^T) G (\mathbf{1}_n - \frac{1}{n} e e^T),$$

where, as before, $e = (1, \dots, 1)^T \in \mathbb{R}^n$.

Proof. We get from Eq. (3.4.1) that $W = \Omega^T (\mathbf{1}_n - \frac{1}{n} e e^T)$. This implies,

$$W^T W = (\mathbf{1}_n - \frac{1}{n} e e^T) \Omega \Omega^T (\mathbf{1}_n - \frac{1}{n} e e^T),$$

and we have $\Omega \Omega^T = G$. □

In particular, this lemma shows that $W^T W$ will always have rank at most $n - 1$.

Let us now assume a statistical setting for PCA. Recall that we have data points $x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} x$, where $x \in \mathbb{R}^D$ is a random variable. Let $d \leq D$. We assume that there is a Gaussian **latent variable** $y \sim N(0, \mathbf{1}_d)$ such that

$$x = Ay + b + \varepsilon,$$

where $A \in \mathbb{R}^{D \times d}$ is a matrix of rank d , and we have a Gaussian noise $\varepsilon \sim N(0, \sigma^2 \mathbf{1}_D)$. That is,

$$(x \mid y) \sim N(Ay + b, \sigma^2 \mathbf{1}_D).$$

In particular, we do not work with nonlinear feature maps here.

3 Machine Learning

Proposition 3.42. *In the statistical setting above we have*

$$x \sim N(b, AA^T + \sigma^2 \mathbf{1}_D).$$

Before we prove this result, let us state an immediate corollary.

Corollary 3.43. *If above we assume the prior $y \sim N(v, B)$ with $v \in \mathbb{R}^D$ and $B \in \mathbb{R}^{D \times D}$ positive definite, then we have $x \sim N(Av + b, ABA^T + \sigma^2 \mathbf{1}_D)$.*

Proof. If $y \sim N(v, B)$, then by Lemma 1.29 $y := Ry' + v$ with $y' \sim N(0, \mathbf{1}_D)$, where $RR^T = B$. This gives $x = Ay + b = ARy' + Av + b$. \square

Now, we prove Proposition 3.42.

Proof of Proposition 3.42. By Eq. (1.2.1), we have $P(x) = \int_{\mathbb{R}} P(x | y)P(y) dy$. By assumption, we have

$$\log P(x | y) + \log P(y) = -\frac{1}{2\sigma^2} \|x - (Ay + b)\|^2 - \frac{1}{2} \|y\|^2 + c,$$

where c is independent of x and y . This shows that we can find vectors $v \in \mathbb{R}^d$ and $\mu \in \mathbb{R}^D$, and positive definite matrices Σ and B (where v and B are allowed to depend on x), such that

$$\log P(x | y) + \log P(y) = -\frac{1}{2}(x - \mu)\Sigma^{-1}(x - \mu) - \frac{1}{2}(y - v)B^{-1}(y - v) + c'$$

where c' is independent of x and y . Consequently, if we integrate out y , we obtain a Gaussian random variable for x . We use Lemma 1.26 to determine the parameters

$$\mu = \mathbb{E}x = \mathbb{E}(Ay + b) = b$$

and

$$\begin{aligned} \Sigma &= \mathbb{E}(x - b)(x - b)^T \\ &= \mathbb{E}(Ay + \varepsilon)(Ay + \varepsilon)^T \\ &= \mathbb{E}(Ay)(Ay)^T + \mathbb{E}\varepsilon\varepsilon^T \\ &= \mathbb{E}Ayy^T A^T + \sigma^2 \mathbf{1}_D. \end{aligned}$$

We have, by linearity, $\mathbb{E}Ayy^T A^T = A(\mathbb{E}yy^T)A^T = AA^T$. \square

3 Machine Learning

We can use Proposition 3.42 for maximum likelihood, maximum a-posteriori estimation (see Definition 3.6), or for taking the Bayesian perspective where we assume priors for the parameters $(A, b) \in \mathbb{R}^{D \times d} \times \mathbb{R}^D$. Alternatively, we can also compute the distribution of y given a data point x . We can use this last approach to use data for updating the distribution of y , and then generating synthetic data points by sampling $(x | y)$.

Theorem 3.44. *Suppose that we have a prior $y \sim N(\mathbf{v}, B)$ and $(x | y) \sim N(Ax + b, \sigma^2 \mathbf{1}_D)$. Then, the posterior distribution of y given x is*

$$(y | x) \sim N(m, C)$$

with covariance matrix $C = (\sigma^{-2} A^T A + B^{-1})^{-1}$ and $m = C(\sigma^{-2} A^T (x - b) + B^{-1} \mathbf{v})$.

Proof. By Bayes' theorem for densities (Theorem 1.22) we have

$$\log P(y | x) = \log P(x | y) + \log P(y) + c,$$

where c does not depend on y . By assumption, we have $(x | y) \sim N(Ay + b, \sigma^2 \mathbf{1}_D)$ and $y \sim N(\mathbf{v}, B)$, so that

$$\log P(y | x) = -\frac{1}{2\sigma^2} \|Ay - (x - b)\|^2 - \frac{1}{2} (y - \mathbf{v})^T B^{-1} (y - \mathbf{v}) + c'$$

where c' does not depend on y . This expression has the same form as in Eq. (3.2.7), and we have computed the conditional density in Eq. (3.2.7) in the proof of Theorem 3.17. Following the proof gives us $(y | x) \sim N(m, C)$ with covariance matrix and expected value as stated. \square

Exercise 3.9. Let $z \in \mathbb{R}^M$ be a random variable with $\mu := \mathbb{E}z \in \mathbb{R}^M$. Show that the covariance matrix of z is given by $\Sigma = \mathbb{E}(z - \mu)(z - \mu)^T$. Use this to show that Σ is positive semi-definite.

Exercise 3.10. Take again the MNIST dataset from the `MLDatasets.jl` package and load the training data for pictures of ones and zeros. Use PCA to reduce the number of parameters representing these pictures. Then, load a point x from the test data set and compute the posterior distribution for $(y | x)$ in Theorem 3.44 using x . Use the posterior distribution to generate synthetic data.

Bibliography

- [Ash70] Robert B. Ash. *Basic probability theory*. John Wiley & Sons, Inc., New York-London-Sydney, 1970.
- [AZ18] Martin Aigner and Günter M. Ziegler. *Proofs from THE BOOK*. Springer Publishing Company, Incorporated, 6th edition, 2018.
- [BEKS17] Jeff Bezanson, Alan Edelman, Stefan Karpinski, and Viral Shah. Julia: A fresh approach to numerical computing. *SIAM Review*, 59(1):65–98, 2017.
- [Chu97] Fan R. K. Chung. *Spectral graph theory*, volume 92 of *CBMS Regional Conference Series in Mathematics*. Published for the Conference Board of the Mathematical Sciences, Washington, DC; by the American Mathematical Society, Providence, RI. Available at <https://mathweb.ucsd.edu/~fan/research/revised.html>, 1997.
- [Chu10] Fan Chung. Graph theory in the information age. *Notices Amer. Math. Soc.*, 57(6):726–732, 2010.
- [DFO20] Marc Peter Deisenroth, A. Aldo Faisal, and Cheng Soon Ong. *Mathematics for Machine Learning*. Cambridge University Press, 2020.
- [FV07] Alan Frieze and Eric Vigoda. A survey on the use of Markov chains to randomly sample colourings. *Combinatorica*, 01 2007.
- [GK12] Venkatesan Guruswami and Ravi Kannan. Computer science theory for the information age, 2012. Available at <https://www.cs.cmu.edu/~venkatg/teaching/CStheory-infoage/>.
- [HJ92] Richard A. Horn and Charles R. Johnson. *Matrix analysis*, volume 349. Cambridge University Press, Cambridge, 1992.
- [PBMW98] Larry Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The PageRank citation ranking: Bringing order to the Web. In *Proceedings of the 7th International World Wide Web Conference*, pages 161–172, Brisbane, Australia, 1998.

Bibliography

- [RGS09] José C. Rosales and Pedro A. García-Sánchez. *Numerical Semigroups*. Developments in Mathematics. Springer New York, 2009.
- [SS11] Thomas Sauerwald and He Sun. Spectral graph theory, 2011. Available at <https://resources.mpi-inf.mpg.de/departments/d1/teaching/ws11/SGT/>.
- [Str93] Gilbert Strang. The fundamental theorem of linear algebra. *Amer. Math. Monthly*, 100(9):848–855, 1993.