

Mathematics for Machine Learning

NYU Paris, [CSCI-UA 9473](#)

Augustin Cosse

Draft, May 26, 2019

These notes provide a summary of the mathematics needed for an introductory class in machine learning. The version is temporary. Please direct any comments or questions to acosse@nyu.edu or acosse@ens.fr

1 Notations

$\mathbb{R}, \mathbb{Z}, \mathbb{N}$,	Real, integer and natural numbers
$\langle \cdot, \cdot \rangle$	Inner product, $\langle \mathbf{v}, \mathbf{w} \rangle = \sum_{i=1}^n v_i w_i$, $\langle \mathbf{A}, \mathbf{B} \rangle = \sum_{i=1}^n \sum_{j=1}^m A_{ij} B_{ij}$
$\text{Tr}(\mathbf{A})$	Trace of a matrix $\text{Tr}(\mathbf{A}) = \sum_{i=1}^n A_{ii}$
$\mathcal{N}(\mu, \sigma^2)$	Normal distribution with mean μ , variance σ^2 (see below)
$\boldsymbol{\theta}$	vector encoding the parameters of a distribution, e.g. for the Gaussian, $\boldsymbol{\theta} = (\mu, \sigma^2)$
$A \propto B$	A Proportional to B . We will use this symbol often when discussing inference. For example we will often write $p(\theta x) \propto P(x \theta)p(\theta)$ where we neglect the normalizing constant $p(x)$.
$\delta(x = k)$	Indicator function, $\delta(x = k) = 1$ if $x = k$ and 0 otherwise
\mathcal{D}	Usually used to represent a dataset, $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$
$\phi(\mathbf{x})$	For a given prototype \mathbf{x} , we usually reserve the notation $\phi(\mathbf{x})$ to denote the corresponding feature vector generated from \mathbf{x}
$\frac{\partial f}{\partial x_i}$	partial derivative of $f(x_1, \dots, x_N)$ with respect to the variable x_i
$\nabla f(\mathbf{x})$	gradient, $\nabla f(\mathbf{x}) = (\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n})$
$\mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$	Eigenvalue decomposition, $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_N)$
$\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$	Singular value decomposition $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_N)$
$\mathcal{H}(p)$	Shannon Entropy

Table 1: Notations

1.1 Notions from optimization

When discussing optimization, we will often make use of the notations max, min as well as argmax and argmin. When looking for the minimum *value*, f^* , of a function, we write

$$f^* = \min_x f(x) \quad (1)$$

When looking for *the point* x^* at which the function achieves a minimal value (minimal argument), we write

$$x^* = \underset{x}{\text{argmin}} f(x) \quad (2)$$

When looking for the solution of an optimization problem (such as the computation of the MLE or MAP) below, a common approach is to use a gradient descent iteration,

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k) \quad (3)$$

By moving in the direction of $-\nabla f(\mathbf{x}_k)$, provided that we use a sufficiently small step length (or learning rate) α , we are guaranteed to decrease the value of the function at each step. Whenever the Hessian is available and not too expensive, one can turn to second order methods, the most popular of which is Newton's method, whose iterations are defined as

$$\mathbf{x}_{k+1} = \mathbf{x}_k - (\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k) \quad (4)$$

The use of the Hessian enables faster convergence (we can think of second order methods as methods that are using a more global view of the optimization landscape than gradient descent steps) but induces a higher computational cost. In practice we therefore turn to approximations of this Hessian. The corresponding optimization methods are then called quasi-Newton methods. The iterations corresponding to those methods are defined as

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{B}_k^{-1} \nabla f(\mathbf{x}_k) \quad (5)$$

where \mathbf{B}_k is an approximation of the Hessian which is updated at each step. The most popular quasi Newton method is the BFGS algorithm (see for example [17], chapter 8).

1.2 Statistical Inference

For a given probability model p parametrized by the parameter θ which takes values in $\{\theta_1, \theta_2, \dots, \theta_N\}$, Bayes's Theorem (see below) gives

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})} \quad (6)$$

$p(\theta|\mathbf{x})$ is called the [posterior](#), $p(\mathbf{x}|\theta)$ is called [likelihood](#). Finally $p(\theta)$ is called the [prior](#).

Contents

1	Notations	1
2	Linear algebra	3
3	Differential Calculus	4
4	Statistics and probability	5
5	Probabilistic classifiers	11
6	Optimization	14
7	Halfspaces and hyperplanes	16
8	Regression and regularization	17
9	Bias Variance tradeoff	20
10	Kernels	26
11	Support Vector Machines (SVM)	28
12	Neural Networks	30
13	Clustering	34
14	Linear latent variable models	38
15	Manifold Learning	46
16	Reinforcement learning	52

2 Linear algebra

Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$,

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \ddots & & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix} \quad (7)$$

we define its transpose as the matrix $\mathbf{A}^T \in \mathbb{R}^{n \times m}$ whose j^{th} column is defined as the j^{th} row of \mathbf{A} , i.e.

$$\mathbf{A}^T = \begin{pmatrix} a_{11} & a_{21} & \dots & a_{m1} \\ a_{12} & a_{22} & \dots & a_{n2} \\ \vdots & \ddots & & \vdots \\ a_{1n} & \dots & a_{m-1,m} & a_{mn} \end{pmatrix} \quad (8)$$

If \mathbf{A} and \mathbf{B} are $m \times n$ matrices, we have the following properties

- $(\mathbf{A}^T)^T = \mathbf{A}$
- $(\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T$
- $(\alpha \mathbf{A})^T = \alpha \mathbf{A}^T$
- $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$

Definition 1 ([PseudoInverse](#)). Given a matrix \mathbf{A} , the pseudo inverse is defined as the matrix \mathbf{A}^+ that satisfies the following conditions

$$\mathbf{AA}^+ \mathbf{A} = \mathbf{A} \quad (9)$$

$$\mathbf{A}^+ \mathbf{AA}^+ = \mathbf{A}^+ \quad (10)$$

$$(\mathbf{AA}^+)^T = \mathbf{AA}^+ \quad (11)$$

$$(\mathbf{A}^+ \mathbf{A})^T = \mathbf{A}^+ \mathbf{A} \quad (12)$$

When the matrix $\mathbf{A}^T \mathbf{A}$ is invertible, the pseudo inverse can be defined as $(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$

2.1 Norms and inner products

Given two vectors $\mathbf{v} \in \mathbb{R}^n$ and $\mathbf{w} \in \mathbb{R}^n$, one can define an inner product as

$$\mathbf{v} \cdot \mathbf{w} = \langle \mathbf{v}, \mathbf{w} \rangle = \sum_{i=1}^n v_i w_i \quad (13)$$

Generally speaking, a function $\langle \cdot, \cdot \rangle$ is an inner product if it satisfies the following properties

- $\langle x, x \rangle \geq 0$, $\langle x, x \rangle = 0 \iff x = 0$ (positivity)
- $\langle x, y \rangle = \langle y, x \rangle$ (symmetry)
- $\langle x + y, z \rangle = \langle x, z \rangle + \langle y, z \rangle$ (additivity)
- $\langle \alpha x, y \rangle = \alpha \langle x, y \rangle$, for all $\alpha \in \mathbb{R}$ (homogeneity)

So far, we have discussed inner products for vectors, but it is also possible to define an inner product on matrices. A standard inner product can be defined in a similar way on matrices $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{m \times n}$,

$$\langle \mathbf{X}, \mathbf{Y} \rangle = \text{Tr}(\mathbf{X}^T \mathbf{Y}) = \sum_{i=1}^m \sum_{j=1}^n X_{ij} Y_{ij} \quad (14)$$

Here we used the notation $\text{Tr}(\mathbf{A})$ to denote the trace of a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, i.e. $\text{Tr}(\mathbf{X}) = \sum_{i=1}^n A_{ii}$.

Another important function is the norm. A norm is a function $f : \mathbb{R}^n \mapsto \mathbb{R}$ and which has the following properties

- $f(\mathbf{x}) \geq 0$, $f(\mathbf{x}) = 0 \iff \mathbf{x} = 0$ (positivity)
- $f(\alpha \mathbf{x}) = |\alpha| f(\mathbf{x})$, $\forall \alpha \in \mathbb{R}$ (homogeneity)
- $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$ (triangle inequality)

Examples of norms include

- The ℓ_2 norm: $\|\mathbf{v}\|_2 = \sqrt{\sum_{i=1}^n v_i^2}$
- The ℓ_1 norm: $\|\mathbf{v}\|_1 = \sum_{i=1}^n |v_i|$
- The ℓ_∞ norm: $\|\mathbf{v}\|_\infty = \max_i |v_i|$
- More generally, the ℓ_p norms are defined for $p \geq 1$ as $\|\mathbf{v}\|_p = (\sum_{i=1}^n |v_i|^p)^{1/p}$

Just as for vectors, one can define norms on matrices. Those norms will satisfy the same properties. The most popular one is the Frobenius norm

$$\|\mathbf{A}\|_F = \|\mathbf{A}\|_2 = \sqrt{\sum_{i=1}^n \sum_{j=1}^n |A_{i,j}|^2} = \sqrt{\text{Tr}(\mathbf{A}^T \mathbf{A})} \quad (15)$$

In fact given any inner product (on vectors or matrices), one can define an "induced" norm by letting $f(\mathbf{v}) = \|\mathbf{v}\| = \sqrt{\langle \mathbf{v}, \mathbf{v} \rangle}$

3 Differential Calculus

Given a multivariate scalar function $f(\mathbf{x})$, $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$, one defines the gradient of the function as the vector encoding the derivatives of the function f with respect to each of the components of \mathbf{x} , i.e.,

$$\nabla f(\mathbf{x}) = \frac{\partial f}{\partial \mathbf{x}} = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right) \quad (16)$$

Note that when we differentiate a field or vector function, we get a matrix, whose entries correspond to the partial derivatives of each of the components of the field with respect to each of the variables. This matrix is known as the [Jacobian matrix](#). For $\vec{f}(\mathbf{x}) = (f_1(x_1, \dots, x_n), \dots, f_m(x_1, \dots, x_n))$, we have

$$\frac{\partial \vec{f}}{\partial \mathbf{x}} \equiv \left(\frac{\partial f_i(\mathbf{x})}{\partial x_j} \right)_{ij} = J = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \ddots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix} \quad (17)$$

The Hessian matrix of a multivariate function $f : \mathbb{R}^n \mapsto \mathbb{R}$, is the matrix encoding all second order partial derivatives of f ,

$$\nabla^2 f = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \ddots & & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix} \quad (18)$$

Compactly we can write

$$H = \nabla^2 f = \left(\frac{\partial^2 f}{\partial x_i \partial x_j} \right)_{ij} \quad (19)$$

3.1 Matrix differentiation

When considering the derivative of a function $f(\mathbf{A})$, $\mathbf{A} \in \mathbb{R}^{m \times n}$ with respect to any matrix \mathbf{A} , we will use the abridged notation $\frac{\partial f(\mathbf{A})}{\partial \mathbf{A}}$ to denote the $m \times n$ matrix encoding the first order (partial) derivatives of the function f with respect to each of the entries of the matrix \mathbf{A} , i.e.

$$\frac{\partial f(\mathbf{A})}{\partial \mathbf{A}} = \begin{pmatrix} \frac{\partial f(\mathbf{A})}{\partial a_{11}} & \frac{\partial f(\mathbf{A})}{\partial a_{12}} & \cdots & \frac{\partial f(\mathbf{A})}{\partial a_{1m}} \\ \frac{\partial f(\mathbf{A})}{\partial a_{21}} & \frac{\partial f(\mathbf{A})}{\partial a_{22}} & \cdots & \frac{\partial f(\mathbf{A})}{\partial a_{2m}} \\ \vdots & \ddots & & \vdots \\ \frac{\partial f(\mathbf{A})}{\partial a_{m1}} & \cdots & \frac{\partial f(\mathbf{A})}{\partial a_{m,n-1}} & \frac{\partial f(\mathbf{A})}{\partial a_{mn}} \end{pmatrix} \quad (20)$$

4 Statistics and probability

4.1 Probability

Definition 2. Given a sample space \mathcal{S} and an associated algebra \mathcal{B} , a probability function is a function P with a domain \mathcal{B} that satisfies

- $P(A) \geq 0$ for all $A \in \mathcal{B}$
- $P(\mathcal{S}) = 1$
- If $A_1, A_2, \dots \in \mathcal{B}$ are pairwise disjoint, then $P(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$

Theorem 1. If P is a probability function and A is any set in \mathcal{B} , then

- $P(\emptyset) = 0$, where \emptyset is the empty set
- $P(A) \leq 1$
- $P(A^c) = 1 - P(A)$

Theorem 2. if P is a probability function and A and B are any sets in \mathcal{B} , then

- $P(B \cap A^c) = P(B) - P(A \cap B)$
- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- If $A \subset B$, then $P(A) \leq P(B)$.

Definition 3. If A and B are events in S , and $P(B) > 0$, then the conditional probability of A given B , written $P(A|B)$ is

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad (21)$$

Given definition (3), it is easy to see that $P(A \cap B) = P(A|B)P(B) = P(B \cap A) = P(B|A)P(A)$ (this idea is sometimes known as the [chain rule](#)). The generalization of this gives Bayes' rule:

Theorem 3 (Bayes' rule). Let A_1, A_2, \dots be a partition of a sample space, and let B be any set. Then, for each $i = 1, 2, \dots$,

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{\sum_{j=1}^{\infty} P(B|A_j)P(A_j)} \quad (22)$$

(the denominator is simply the expansion $P(B) = \sum_{j=1}^{\infty} P(B|A_j)P(A_j)$)

Sometimes, it could happen that the occurrence of an event B has no influence on the probability of another event A , i.e.

$$P(A|B) = P(A) \quad (23)$$

If this condition holds, using Bayes's Theorem, it immediately implies $P(B|A) = \frac{P(A|B)P(B)}{P(A)} = P(B)$. From this, we also have $P(A \cap B) = P(A)P(B)$. This last relation is the precise definition of [statistical independence](#)

Definition 4 (Independence). Two events A and B are statistically independent if

$$P(A \cap B) = P(A)P(B) \quad (24)$$

Definition 5. A collection of events A_1, \dots, A_n are mutually independent if for any subcollection A_{i_1}, \dots, A_{i_k} , we have

$$P\left(\bigcap_{j=1}^k A_{i_j}\right) = \prod_{j=1}^k P(A_{i_j}) \quad (25)$$

Theorem 4. If two events A and B are independent, the following pairs of events are also independent

- A and B^c
- A^c and B
- A^c and B^c

4.2 Random variables

Definition 6 (Random variable). A random variable is a function from a sample space into the real numbers

It is important to understand that the set of random variables includes the samples themselves as one can of course choose the random variable to match a sample onto the value of this sample. However, the probability is only defined on the sample space. For example, a simple Gaussian random variable X is a function from the sample space $\mathcal{S} \equiv (-\infty, \infty)$ into this same sample space (or if you prefer the value of the sample). I.e. $X(s_i) = s_i$ in this case. But the definition is more general and one could for example define another random variable $Y = X^2$. Now the random variable is not the identity anymore. It a function which takes samples from the sample space \mathcal{S} and raise them to the power 2. The connection to the sample space is important because the

notion of probability P is defined only on this sample space. I.e we only know how likely an event s_i is, we don't know about $X(s_i)$. It is however possible to define an induced probability on the random variable X , by using P ,

$$P_X(X \in A) = P(\{s \in \mathcal{S} : X(s) \in A\}) \quad (26)$$

In words, we look at the samples s for which $X(s) \in A$ and we define the probability of the event $X \in A$ from the probability of the corresponding sample $\{s \in \mathcal{S} : X(s) \in A\}$.

We are now ready to introduce the notions of cumulative distribution function (cdf), probability mass function (pmf) and probability density function (pdf).

Definition 7. The *cumulative distribution function* or *cdf* of a random variable X which we denote $F_X(x)$, is defined by

$$F_X(x) = P(X \leq x), \quad \text{for all } x \quad (27)$$

Definition 8. The *probability mass function* (pmf) of a discrete random variable X defined as

$$f_X(x) = P(X = x), \quad \text{for all } x \quad (28)$$

Definition 9. The *probability density function* or *pdf*, $f_X(x)$ of a continuous random variable X is the function defined as

$$F_X(x) = \int_{-\infty}^x f_X(t) dt, \quad \text{for all } x \quad (29)$$

We now introduce the notion of expectation or expected value (a.k.a mean or average) of a random variable. The expected value $\mathbb{E}\{X\}$ of a random variable X can be thought of as an indication of a typical value that the random variable X will take, hence a measure of the "center" of the distribution.

Definition 10. The *expected value* of a random variable $h(X)$, which we denote $\mathbb{E}\{h(X)\}$ is defined as

$$\mathbb{E}h(X) = \begin{cases} \int_{-\infty}^{\infty} h(x)f_X(x) dx & \text{if } X \text{ is continuous} \\ \sum_{x \in \mathcal{X}} h(x)f_X(x) = \sum_{x \in \mathcal{X}} h(x)P(X = x) & \text{if } X \text{ is discrete} \end{cases} \quad (30)$$

In the definition above, \mathcal{X} is used to denote the set of all possible outcomes of the random variable X .

Theorem 5. If X is a random variable and α, β, γ are constants, then for any functions h_1 and h_2 with well defined expectations, we have

- $\mathbb{E}\{\alpha h_1(X) + \beta h_2(X) + \gamma\} = \alpha \mathbb{E}\{h_1(X)\} + \beta \mathbb{E}\{h_2(X)\} + \gamma$
- If $h_1(x) \geq 0$ for all x , then $\mathbb{E}\{h_1\} \geq 0$
- If $h_1(x) \geq h_2(x)$ for all x , then $\mathbb{E}\{h_1(X)\} \geq \mathbb{E}\{h_2(X)\}$
- If $\alpha \leq h_1(x) \leq \beta$ for all x , then $\alpha \leq \mathbb{E}\{h_1(X)\} \leq \beta$.

A particular class of expectations can be defined for a random variable X by choosing the function $h_1(X)$ to encode the monomials X^α , $\alpha \in \mathbb{N}$. Those expectations are then called *moments* of the random variable X .

Definition 11. For each integer $n \in \mathbb{N}$, the *n^{th} moment* of X , μ'_n is defined as

$$\mu'_n = \mathbb{E}\{X^n\} \quad (31)$$

Accordingly, the *n^{th} central moment* of X , μ_n is defined as

$$\mu_n = \mathbb{E}\{(X - \mu)^n\} \quad (32)$$

In particular, note that we have $\mu = \mu'_1 = \mathbb{E}X$.

Definition 12. The *variance* of a random variable X is its second central moment, i.e. $\sigma^2 = \text{Var}(X) = \mathbb{E}\{(X - \mathbb{E}\{X\})^2\}$. The (positive) square root of σ^2 , σ is called *standard deviation*.

For two random variables X and Y , one can also define the correlation and covariance of X and Y ,

Definition 13. The *covariance* of X and Y is the number defined by

$$\text{Cov}(X, Y) = \mathbb{E}\{(X - \mu_X)(Y - \mu_Y)\} \quad (33)$$

Definition 14. The *correlation* or *correlation coefficient* of X and Y is the number defined by

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \quad (34)$$

If similar values of X and Y occur at the same time, the covariance will be large and positive. If on the contrary, large positive values of X occur with large negative values for Y , the covariance $\text{Cov}(X, Y)$ will be large and negative. Hence we can already say that the sign of $\text{Cov}(X, Y)$ gives information on the type of relation that exists between X and Y .

The correlation coefficient $\rho_{X,Y}$ between X and Y takes value between -1 and 1 with a value of 1 indicating a perfect linear relationship between the two variables. Another important aspect of the correlation coefficient is its connection to statistical independence.

Theorem 6. If two variables X and Y are independent, then $\text{Cov}(X, Y) = 0$ and $\rho_{X,Y} = 0$.

Proof. The proof follows from the definition of independence and the fact that independence implies $\mathbb{E}\{XY\} = \mathbb{E}X\mathbb{E}Y$ \square

4.3 Classical PDFs

We start by listing some important distributions:

- The *Gaussian/Normal Distribution* is the most popular distribution in statistics. The univariate distribution reads as

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}|x - \mu|^2\right) \quad (35)$$

The multivariate distribution, for a d -dimensional vector \mathbf{x} , reads as

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\} \quad (36)$$

- The *Binomial distribution* (which we write $\text{Bin}(k|n, \theta)$) is used to represent the probability of observing the positive/successful outcomes (the probability of success is given by θ) that repeats itself k times over n successive trials.

$$\text{Bin}(k|n, \theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k} \quad (37)$$

The first factor encodes the number of possibilities of choosing k elements among the n without replacement.

$$\binom{n}{k} = \frac{n!}{(n-k)!k!} \quad (38)$$

The mean and variance for this distribution are given by $\mu = n\theta$ and $\sigma^2 = n\theta(1 - \theta)$

- The **Bernoulli distribution**. The Bernoulli distribution is defined on a binary variable. I.e if we consider a single experiment (tossing a coin) with only two possible outcomes which are encoded by the variable $y = \{1, 0\}$ and appear with probability p and $1 - p$ respectively. to encode the probability of an outcome 0 or 1 we can thus just raise the probability θ to the power y , as show below. This gives the Bernoulli pdf

$$\text{Ber}(x|\theta) = \theta^y(1 - \theta)^{1-y} \quad (39)$$

$$(40)$$

This is also sometimes written

$$\text{Ber}(x|\theta) = \theta^{\delta(y=0)}(1 - \theta)^{\delta(y=1)} \quad (41)$$

$$(42)$$

Where $\delta(\mathcal{C}) = 1$ if the predicate \mathcal{C} is satisfied.

- **Categorical distribution** or **Multinoulli distribution**. When considering a variable with K possible states (e.g. a dice with 6 faces), we will often encodes those states through dummy variables. This means that for example for a variable that can take three different values, $y = 1$, $y = 2$ or $y = 3$ (we will see later that this applies to features or classes in classification problems), we would encode each of these feature through the 3-tuples $\mathbf{t} = (1, 0, 0)$ ($y = 1$), $\mathbf{t} = (0, 1, 0)$ ($y = 2$) and $\mathbf{t} = (0, 0, 1)$ ($y = 3$). If we use $\theta_j = p_j$, $j = 1, \dots, 3$ to denote the probability to get each feature, the total probability density function can now read compactly as

$$\text{Cat}(x|\boldsymbol{\theta}) = \text{Mu}(x|\boldsymbol{\theta}) = p(x|\boldsymbol{\theta}) = \prod_{j=1}^K \theta_j^{t_j} = \theta_1^{t_1} \dots \theta_K^{t_K} = \boldsymbol{\theta}^{\mathbf{t}} \quad (43)$$

- The **Laplace distribution** is characterized by "fatter" tails than the Gaussian distribution (i.e the probability of getting values that are very different from the mean will be higher in the Laplace distribution). The pdf of the Laplace distribution reads as

$$p(x|\mu, \beta) = \frac{1}{2\beta} \exp\left(-\frac{|x - \mu|}{\beta}\right) \quad (44)$$

one way to see this is that in the Laplace distribution the argument of the exponential is the absolute value whereas in the Normal distribution, the argument is squared which implies a faster decrease in the pdf for large values of the deviation to the mean, $|x - \mu|$.

4.4 A few important Theorems

Theorem 7 (Weak Law of large numbers). Let X_1, X_2, \dots be identically and independently distributed (i.i.d) random variables with $\mathbb{E}X_i = \mu$ and $\text{Var}X_i = \sigma^2$. Define $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$. Then, for every $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| < \varepsilon) = 1 \quad (45)$$

That is, \bar{X}_n converges in probability to μ .

Theorem 8 (Strong Law of large numbers). Let X_1, X_2, \dots be i.i.d. random variables with $\mathbb{E}X_i = \mu$ and $\text{Var}X_i = \sigma^2 < \infty$ and define $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$. Then for every $\epsilon > 0$,

$$P\left(\lim_{n \rightarrow \infty} |\bar{X}_n - \mu| < \epsilon\right) = 1 \quad (46)$$

that is \bar{X}_n converges almost surely to μ

Theorem 9 (Central Limit Theorem). Let X_1, X_2, \dots be a sequence of i.i.d random variables whose mgf's exist in a neighborhood of 0 (that is $M_{X_i}(t)$ exists for $t < |h|$, for some positive h). Let $\mathbb{E}X_i = \mu$ and $\text{Var}X_i = \sigma^2$ (both μ and σ^2 are finite since the mgf exists). Define the empirical mean \bar{X}_n as $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$. Let $G_n(x)$ denote the cdf of $\sqrt{n}(\bar{X}_n - \mu)/\sigma$. Then for any x , $-\infty < x < \infty$,

$$\lim_{n \rightarrow \infty} G_n(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy \quad (47)$$

That is, $\sqrt{n}(\bar{X}_n - \mu)/\sigma$ has a limiting standard normal distribution.

4.5 MAP and MLE

This part gives a very basic introduction to the difference between the Frequentist and Bayesian notions of estimators. For more details see [4]

Definition 15. A *point estimator* is any function of the sample $W(x_1, \dots, x_n)$

The most popular method to derive an estimator on the parameters of a distribution is the method of *maximum likelihood*. For a given probability distribution $p(x_i|\theta_1, \dots, \theta_k)$ parametrized by θ , the likelihood function is defined as

$$L(\theta|\mathbf{x}) = p(\mathbf{x}|\theta) = \prod_{i=1}^n p(x_i|\theta_1, \dots, \theta_n) \quad (48)$$

This function encodes the "likelihood" of observing the samples x_i for a particular choice of parameters θ (we write $L(\theta|\mathbf{x})$ because we view it as a function of θ). Given this function, it seems reasonable to look for the parameters $\theta_1, \dots, \theta_n$ which give the highest probability of observing the sample (i.e the distribution which gives the highest probability of observing the x_i is likely to be the distribution of those x_i). This idea leads to the *Maximum likelihood estimator* (MLE) whose formal definition is given below

Definition 16. For each sample point \mathbf{x} , let $\hat{\theta}(\mathbf{x})$ be a parameter value at which $L(\theta|\mathbf{x})$ attains its maximum as a function of θ , with \mathbf{x} held fixed. A *maximum likelihood estimator* (MLE) of the parameter θ based on a sample X is $\hat{\theta}(X)$. Mathematically, we write

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} p(\mathbf{x}|\theta) \quad (49)$$

The MLE is considered a Frequentist estimator as it only relies on the sample itself without making any prior assumption on the parameters.

The Bayesian approach to statistics is fundamentally different. So far, we assumed that θ was an unknown but fixed quantity. I.e. A random sample X_1, \dots, X_n was drawn from a population and based on those observations, some knowledge on θ was obtained. In the Bayesian approach, θ is described by a probability distribution (the prior distribution). We thus assume some level of uncertainty on θ . The prior is a subjective distribution which is set by the observer.

The idea is that the observer updates his original prior based on the sample that he draws from the population. The update is done by means of Bayes' rule,

$$p(\theta|\mathbf{x}) \propto p(\mathbf{x}|\theta)p(\theta) \quad (50)$$

We can then estimate the parameter θ by finding the value which maximizes the new prior, $p(\theta|\mathbf{x})$. This idea is known as *Maximum a posteriori* (MAP). The MAP estimator is thus defined as

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} p(\mathbf{x}|\theta)p(\theta) \propto p(\theta|\mathbf{x}) \quad (51)$$

4.5.1 Linear Gaussian system

By linear Gaussian system, we mean here that we consider both a Gaussian likelihood and a Gaussian prior, and connect the two by means of any linear transformation,

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} | \boldsymbol{\mu}_\theta, \boldsymbol{\Sigma}_\theta) \quad (52)$$

$$p(\mathbf{x} | \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x} | \mathbf{A}\boldsymbol{\theta} + \mathbf{b}, \boldsymbol{\Sigma}_x) \quad (53)$$

If we write down the probability of \mathbf{x} , from such a model, we get

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{A}\boldsymbol{\mu}_\theta + \mathbf{b}, \boldsymbol{\Sigma}_x + \mathbf{A}\boldsymbol{\Sigma}_\theta\mathbf{A}^T) \quad (54)$$

4.6 Exponential family

Many of the probability density functions discussed above are part of a general family of distributions known as the [exponential family](#).

A probability density function, $f(\mathbf{x} | \boldsymbol{\theta})$ is said to be in the exponential family if it is of the form

$$f(\mathbf{x} | \boldsymbol{\theta}) = h(\mathbf{x}) c(\boldsymbol{\theta}) \exp \left(\sum_{i=1}^k w_i(\boldsymbol{\theta}) t_i(\mathbf{x}) \right) \quad (55)$$

$$= \frac{1}{Z(\boldsymbol{\theta})} h(\mathbf{x}) \exp(\boldsymbol{\theta}^T \phi(\mathbf{x})) \quad (56)$$

$$= h(\mathbf{x}) \exp(\boldsymbol{\theta}^T \phi(\mathbf{x}) - A(\boldsymbol{\theta})) \quad (57)$$

where

$$Z(\boldsymbol{\theta}) = \int_{\mathcal{X}^m} h(\mathbf{x}) \exp(\boldsymbol{\theta}^T \phi(\mathbf{x})) d\mathbf{x} \quad (58)$$

$$A(\boldsymbol{\theta}) = \log Z(\boldsymbol{\theta}) \quad (59)$$

$Z(\boldsymbol{\theta})$ is called the partition function and $A(\boldsymbol{\theta})$ is called the log partition function. $\boldsymbol{\theta}$ are called the natural parameters or canonical parameters.

5 Probabilistic classifiers

When discussing classification, we often either look at classifiers as geometric objects (i.e separating planes, maximal margin hyperplanes,...) or as probabilistic objects based on Bayes Theorem. In the last framework, we consider two main groups of classifiers

- [Generative](#) classifiers
- [Discriminative](#) classifiers

The first class defines a model for the joint probability distribution $p(\mathbf{x}, y)$ (or equivalently for the class conditional density $p(\mathbf{x} | y)$) and hence provides a way to generate new samples \mathbf{x} . The second class defines a model for the probability distribution $p(y | \mathbf{x})$ (i.e the class posterior) and hence only provides a way to discriminate between classes. Below we discuss the two most popular discriminative models: The [Naive Bayes classifier](#) (generative) and the [Logistic regression classifier](#) (discriminative)

5.1 Naive Bayes classifier

The Naive Bayes classifier is a generative classifier (i.e we learn a model for $p(\mathbf{x} | y)$ or $p(\mathbf{x}, y)$). We consider a classification problem in which the feature vectors are D dimensional. Moreover,

each of the D features can take K values. I.e, $\mathbf{x} \in \{1, \dots, K\}^D$. If we assume that the features are independent within a class c (note that this is usually not the case), we can write the “class conditional density” as

$$p(\mathbf{x}|y = \mathcal{C}_\ell, \boldsymbol{\theta}) = \prod_{j=1}^D p(x_j|y = \mathcal{C}_\ell, \boldsymbol{\theta}_{j\ell}) \quad (60)$$

In the equation above $\boldsymbol{\theta}_{j\ell}$ is the parameter associated to feature j in class \mathcal{C}_ℓ (see below).

Possible choices for the probability distribution include

- the Gaussian distribution

$$p(\mathbf{x}|y = \mathcal{C}_\ell; \boldsymbol{\theta}) = \prod_{j=1}^D \mathcal{N}(x_j|\mu_{j\ell}, \sigma_{j,\ell}^2) \quad (61)$$

Here $\mu_{j,\ell}$ and $\sigma_{j,\ell}^2$ are the mean and variance for the j^{th} feature in class \mathcal{C}_ℓ .

- When the features are encoded through binary variables $\beta = \{0, 1\}$, the multivariate Bernoulli distribution (or Bernoulli product model) is often used. In this case, we have

$$p(\mathbf{x}|y = \mathcal{C}_\ell, \boldsymbol{\theta}) = \prod_{j=1}^D \text{Ber}(x_j|\mu_{j,\ell}) \quad (62)$$

Again, here $\mu_{j,\ell}$ is the probability of getting feature j in class ℓ .

5.1.1 Deriving the parameters through Maximum Likelihood

We want to train a classification algorithm from a set of pairs (\mathbf{x}_i, c_i) . The prototypes \mathbf{x}_i are represented by D -dimensional feature vectors $x_i = (x_{i1}, \dots, x_{iD})$. From this, if we assume that the probabilities of observing the features are independent, the probability of observing a given pair (\mathbf{x}_i, c_i) reads as

$$p(\mathbf{x}_i, c_i|\boldsymbol{\theta}) = \prod_{j=1}^D p(x_{ij}|c_i, \boldsymbol{\theta}_j) p(c_i|\boldsymbol{\pi}) \quad (63)$$

Here $p(c_i|\boldsymbol{\pi})$ encodes the probability of observing the class i , $p(c_i|\boldsymbol{\pi}) = (\pi_1^{\delta(c_i=1)}, \dots, \pi_N^{\delta(c_i=N)})$ encode the probabilities of getting a particular class \mathcal{C}_1 to \mathcal{C}_C (we use C to denote the number of classes). Using the δ notation, we can thus write (63) as

$$p(\mathbf{x}_i, c_i|\boldsymbol{\theta}) = \prod_{c=1}^C \pi_c^{\delta(c_i=c)} \prod_{j=1}^D \prod_{c=1}^C p(x_{ij}|\boldsymbol{\theta}_{c,j})^{\delta(y_i=c)} \quad (64)$$

Taking the product over all pairs $(x_i, c_i) \in \mathcal{D}$ and then the log, the log likelihood can read as

$$\log(\mathcal{D}|\boldsymbol{\theta}) = \sum_{c=1}^C N_c \log(\pi_c) + \sum_{j=1}^D \sum_{c=1}^C \sum_{i \text{ s.t. } y_i=c} \log p(x_{ij}|\boldsymbol{\theta}_{j,c}) \quad (65)$$

Maximizing the log likelihood with respect to the parameters of the model then gives

$$\hat{\pi}_c = \frac{N_c}{N} \quad (66)$$

as an estimate for the class prior (the probability of being in Class c). This estimate is thus the ratio between the number of samples in class c and the total number of training samples.

If we assume that all N features x_{1i}, \dots, x_{1N} are binary and that all the vectors within a class follow the same distribution, then a good choice for $p(x_{ij}|\theta_{jc})$ is to take a Bernoulli distribution with parameter θ_{jc} . In this case we can show the MLE for the parameter θ_{jc} is given by the ratio of the number of vectors from class c which have their j^{th} feature equal to 1 over the total number of points from class c ,

$$\hat{\theta}_{jc} = \hat{p}(x_{ij} = 1|y = c) = \frac{N_{jc}}{N_c} \quad (67)$$

Once we have computed the parameters of the model, this model can be used to classify new points by relying again on Bayes (neglecting the normalizing constant), we have

$$p(y = c|\mathbf{x}, \mathcal{D}) \propto p(y = c|\mathcal{D}) \prod_{j=1}^D p(x_j|y = c, \mathcal{D}) \quad (68)$$

$$\propto p(y = c|\mathcal{D})p(\mathbf{x}|y = c, \mathcal{D}) \quad (69)$$

5.2 Logistic regression

Logistic regression is perhaps the most popular [discriminative classifier](#). For data pairs $\{\mathbf{x}, y\}$ of prototypes and their associated labels, the probabilistic model is defined as

$$p(y|\mathbf{x}, \mathbf{w}) = \text{Ber}(y|\text{sigm}(\mathbf{w}^T \mathbf{x})) \quad (70)$$

where the sigmoid $\text{sigm}()$ is defined as

$$\text{sigm}(x) = \frac{1}{1 + \exp(-x)} \quad (71)$$

We thus have

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \frac{\exp(\mathbf{w}^T \mathbf{x})}{1 + \exp(\mathbf{w}^T \mathbf{x})} \quad (72)$$

$$p(y = 0|\mathbf{x}, \mathbf{w}) = 1 - \frac{\exp(\mathbf{w}^T \mathbf{x})}{1 + \exp(\mathbf{w}^T \mathbf{x})} \quad (73)$$

The interesting aspect of logistic regression is that the ratio between the probability of each class (which can be used to do classification) is linear in the parameters of the model. You can indeed verify that

$$\log\left(\frac{p(y = 1|\mathbf{x}, \mathbf{w})}{p(y = 0|\mathbf{x}, \mathbf{w})}\right) = \mathbf{w}^T \mathbf{x} \quad (74)$$

This model in fact falls in the class of “generalized” linear models for that precise reason. This idea easily extends to multiple classes by introducing “log-odd ratios” between the probability of each class and the probability of the last class. For K classes, we have $K - 1$ such ratios (the last probability is fixed through $\sum_i p(\mathcal{C} = i|x_n) = 1$)

$$\log\left(\frac{P(\mathcal{C} = 1|x)}{P(\mathcal{C} = K|x)}\right) = \mathbf{w}_1^T \mathbf{x} \quad (75)$$

$$\log\left(\frac{P(\mathcal{C} = 1|x)}{P(\mathcal{C} = K|x)}\right) = \mathbf{w}_1^T \mathbf{x} \quad (76)$$

$$\vdots \quad (77)$$

$$\log\left(\frac{P(\mathcal{C} = K - 1|x)}{P(\mathcal{C} = K|x)}\right) = \mathbf{w}_{K-1}^T \mathbf{x} \quad (78)$$

The posterior class probabilities are then defined as

$$P(\mathcal{C} = k|\mathbf{x}) = \frac{\exp(\mathbf{w}_k^T \mathbf{x})}{1 + \sum_{\ell=1}^{K-1} \exp(\mathbf{w}_\ell^T \mathbf{x})}, \quad k = 1, \dots, K-1 \quad (79)$$

$$P(\mathcal{C} = K|\mathbf{x}) = \frac{1}{1 + \sum_{\ell=1}^{K-1} \exp(\mathbf{w}_\ell^T \mathbf{x})} \quad (80)$$

The model (and the associated parameters) can be learned through Maximum likelihood, by writing the likelihood function for the N observation (noting c_n the class of \mathbf{x}_n),

$$\ell(\mathbf{w}_\ell, \ell = 1, \dots, K-1) = \sum_{n=1}^N \log(p(\mathcal{C} = c_n|\mathbf{x})) \quad (81)$$

And then minimizing the function (for example through gradient descent).

6 Optimization

In this section, we discuss some additional details in optimization. Finding the global minimizer of sufficiently complex functions is usually hard because iterative algorithms will get trapped in "local minima". When solving optimization problems, we will often be interested in understanding whether the particular point to which our algorithm converges is a local or a global minimum (resp maximum) of the function. Under some conditions (smoothness of the function), the general characterization of the local extremas of a function is relatively easy as shown by the following two propositions

Proposition 1 (First order necessary conditions). *If \mathbf{x}^* is a local minimizer and f is continuously differentiable in an open neighborhood of \mathbf{x}^* , then $\nabla f(\mathbf{x}^*) = 0$.*

Proposition 2 (Second Order Necessary conditions). *If \mathbf{x}^* is a local minimizer of f and $\nabla^2 f$ is continuous in an open neighborhood of \mathbf{x}^* , then $\nabla f(\mathbf{x}^*) = 0$ and $\nabla^2 f(\mathbf{x}^*)$ is positive semidefinite.*

Given a function f , there also exist sufficient conditions under which one can guarantee that a point is a minimizer. This idea is summarized by proposition 3 below

Proposition 3 (Second Order Sufficient conditions). *Suppose that $\nabla^2 f$ is continuous in an open neighborhood of \mathbf{x}^* and that $\nabla f(\mathbf{x}^*) = 0$ and $\nabla^2 f(\mathbf{x}^*)$ is positive definite. Then \mathbf{x}^* is a strict local minimizer of f .*

6.1 Convexity

Definition 17 (Convex set). *A set of points, \mathcal{C} is convex if the line segment between any two points in the set is included in the set, that is to say for any $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{C}$, and any scalar $\alpha \in \mathbb{R}$ with $0 \leq \alpha \leq 1$, we have*

$$\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2 \in \mathcal{C} \quad (82)$$

When a set is not convex, it is always possible to define the smallest convex set that contains \mathcal{C} . This idea in particular plays an important role when replacing non convex function (which therefore cannot be optimized efficiently) by convex approximation whose solution can be found and studied by iterative algorithms.

Definition 18 (convex hull). *Given a set of points \mathcal{S} , the convex hull of \mathcal{S} is defined as*

$$\text{conv}\mathcal{S} = \{\alpha_1 \mathbf{x}_1 + \dots, \alpha_k \mathbf{x}_k | \mathbf{x}_i \in \mathcal{S}, \alpha_i \geq 0, i = 1, \dots, k, \alpha_1 + \dots + \alpha_k = 1\} \quad (83)$$

I.e., $\text{conv}\mathcal{S}$ is the set defined from all possible convex combinations of points from \mathcal{S} .

6.2 Constrained problems and Lagrangian

When considering a basic mathematical programming problem of the form

$$\begin{aligned} \min \quad & f_0(\mathbf{x}) \\ \text{subject to} \quad & f_\ell(\mathbf{x}) \leq 0, \ell = 1, \dots, L. \end{aligned} \quad (84)$$

We will use K to denote the set of feasible points for the program (84) (i.e. the set of points \mathbf{x} satisfying the constraints in (84)), i.e. $K = \{\mathbf{x} \mid f_\ell(\mathbf{x}) \leq 0\}$

We can write the Lagrangian by introducing positive multipliers $\lambda_\ell \in \mathbb{R}$ for each of the constraints $f_\ell(x) \leq 0$. As those constraints are negative, the Lagrangian then reads as

$$L = f_0(x) - \sum_{\ell=1}^L \lambda_\ell f_\ell(\mathbf{x}) \quad (85)$$

The following Theorem then gives necessary conditions (on the Lagrangian) for a point \mathbf{x} to be a minimizer from the feasible set, K .

Theorem 10 (Kuhn-Tucker conditions). *Assume that $f_\ell(\mathbf{x})$ ($\ell = 1, \dots, L$) are all differentiable. If the function $f_0(\mathbf{x})$ attains a local minimum at a point \mathbf{x}^* which belongs the feasible set K , then there exists a vector of multipliers, λ_ℓ^* , $\ell = 1, \dots, L$, such that the following conditions are satisfied*

$$\frac{\partial f_0(\mathbf{x}^*)}{\partial x_j} + \sum_{\ell=1}^L \lambda_\ell^* \frac{\partial f_\ell(\mathbf{x}^*)}{\partial x_j} = 0, \quad (j = 1, \dots, J) \quad (86)$$

$$f_\ell(\mathbf{x}^*) \leq 0, \quad (\ell = 1, \dots, L) \quad (87)$$

$$\lambda_\ell^* f_\ell(\mathbf{x}^*) = 0, \quad (\ell = 1, \dots, L) \quad (88)$$

$$\lambda_\ell^* \geq 0, \quad (\ell = 1, \dots, L) \quad (89)$$

Proof. We can always replace the inequality constraints in (84), by introducing slack variables y_ℓ . I.e., we always have the equivalence

$$f_\ell(\mathbf{x}) \leq 0 \iff f_\ell(\mathbf{x}) + y_\ell^2 = 0 \quad (90)$$

for some variables y_ℓ .

For those slack variables, the Lagrangian reads as

$$L = f_0(\mathbf{x}) + \sum_{\ell=1}^L \lambda_\ell (f_\ell(\mathbf{x}) + y_\ell^2) \quad (91)$$

Any minimum of the Lagrangian must satisfy the zero gradient condition, so we have

$$\frac{\partial L}{\partial x_j} = \frac{\partial f_0(\mathbf{x}^*)}{\partial x_j} + \sum_{\ell=1}^L \lambda_\ell \frac{\partial (f_\ell(\mathbf{x}^*) + (y_\ell^*)^2)}{\partial x_j} = 0, \quad (\ell = 1, \dots, L) \quad (92)$$

$$\frac{\partial L}{\partial y_\ell} = 2\lambda_\ell y_\ell = 0 \quad (93)$$

$$\frac{\partial L}{\partial \lambda_\ell} = f_\ell(\mathbf{x}) + y_\ell^2 = 0 \quad (94)$$

Condition (93) above is equivalent to the "complementary slackness" condition (88). To see this, simply note that $\lambda_\ell y_\ell = 0$ implies either $\lambda_\ell = 0$ or $y_\ell = 0$.

In the second case, we have $y_i = 0$ and hence $f_\ell(x) + y_\ell^2 = 0$ which implies $f_\ell(x) = 0$. For the reverse implication, simply note that $\lambda_\ell \neq 0$ implies $f_\ell = 0$ and hence $y_\ell^2 = -f_\ell = 0$ which finally gives $y_\ell \lambda_\ell = 0$

If $y_\ell \neq 0$ (equivalently, $\lambda_\ell = 0$), on the other hand, then (94) implies $y_\ell \lambda_\ell = -f_\ell \lambda_\ell = 0$ and the equivalence is straightforward.

Condition (94) is equivalent to (86) (i.e the slack variables can be eliminated from the derivative)

We are thus left with showing that conditions (92) to (94) imply the non negativity of the multipliers λ_ℓ . For a general minimization problem,

$$\min \quad f_0(\mathbf{x}) \quad (95)$$

$$\text{subject to} \quad f_i(\mathbf{x}) \leq b_i, \quad (i = 1, 2, \dots, m) \quad (96)$$

if we denote the optimal solution corresponding to a given value of the bounds b_i , as $\mathbf{x}_0(\mathbf{b})$, we have

$$\frac{\partial f_0(\mathbf{x}_0(\mathbf{b}))}{\partial b_i} = -\lambda_i \quad (97)$$

But on the other any increase in the b_i leads to a larger feasible set and hence the variation in the value of f_0 for such a small increase in b_i must always be negative (i.e we can only do better when we increase the feasible set), following from this, we must have

$$\frac{\partial f_0(\mathbf{x}_0(\mathbf{b}))}{\partial b_i} \leq 0 \quad (98)$$

which implies $\lambda_i \geq 0$

□

7 Halfspaces and hyperplanes

Given a vector $\mathbf{w} \in \mathbb{R}^N$, one can define a(n) (affine) hyperplane as the set of all points $\mathbf{x} \in \mathbb{R}^N$ satisfying the relation $\mathbf{w}^T \mathbf{x} + b = 0$.

The vector \mathbf{w} encodes the normal to the hyperplane, which also means that for any two points $\mathbf{x}_1, \mathbf{x}_2$ that lie on the hyperplane, we necessarily have

$$\mathbf{w}^T \mathbf{x}_1 + b = \mathbf{w}^T \mathbf{x}_2 + b = 0 \quad (99)$$

(i.e the two points belong to the hyperplane) and hence,

$$\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0 \quad (100)$$

Every hyperplane naturally defines a splitting of the space into two halfspaces. The set of points which are lying above the plane (i.e formally, the points for which $\mathbf{w}^T \mathbf{x} + b > 0$), and the set of points which are lying under the plane ($\mathbf{w}^T \mathbf{x} + b < 0$). Every hyperplane can thus be used as a natural classifier,

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b \quad (101)$$

where we put the point \mathbf{x} in class \mathcal{C}_0 if $y(\mathbf{x}) > 0$ and in class \mathcal{C}_1 if $y(\mathbf{x}) < 0$.

7.1 Distance of a point to a hyperplane

When discussing the robustness of a classifier (and in particular when introducing the notion of margin), we will need the notion of distance of a point to a hyperplane.

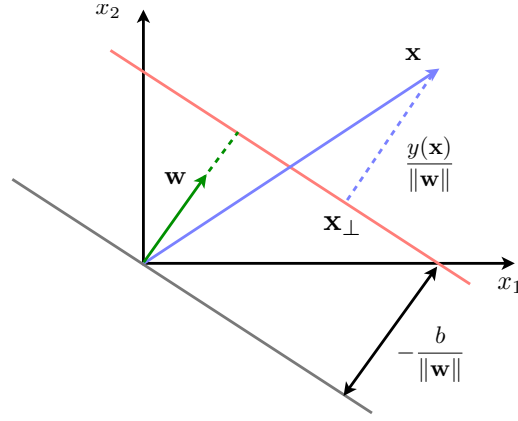


Figure 1: Distance of a point to a hyperplane

Consider Fig 1 below. We consider the plane in red defined as $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$. For any point \mathbf{x} (shown in blue), one can consider the decomposition

$$\mathbf{x} = \mathbf{x}_\perp + d \frac{\mathbf{w}}{\|\mathbf{w}\|} \quad (102)$$

That is we write \mathbf{x} as the combination of its projection onto the plane and a contribution of length d along the normal vector to the plane, \mathbf{w} . d thus encodes the distance of \mathbf{x} to the hyperplane. Now if we multiply (102) by \mathbf{w}^T ,

$$\mathbf{w}^T \mathbf{x} = \mathbf{w}^T \mathbf{x}_\perp + \mathbf{w}^T d \frac{\mathbf{w}}{\|\mathbf{w}\|} \quad (103)$$

and add the bias, we get

$$\mathbf{w}^T \mathbf{x} + b = \mathbf{w}^T \mathbf{x}_\perp + b + \mathbf{w}^T d \frac{\mathbf{w}}{\|\mathbf{w}\|} \quad (104)$$

The left handside is simply $y(\mathbf{x})$. The first term on the right handside is 0 as \mathbf{x}_\perp lies along the plane (hence is orthogonal to \mathbf{w}). We are thus left with

$$y(\mathbf{x}) = \|\mathbf{w}\| \mathbf{d} \quad (105)$$

from which follows

$$d = \frac{y(\mathbf{x})}{\|\mathbf{w}\|} \quad (106)$$

So the signed (because \mathbf{d} can be both positive and negative here) distance, is given by the ratio of the prediction $y(\mathbf{x})$ over the norm of the normal vector \mathbf{w} . We will use this when discussing Maximal margin classifiers.

8 Regression and regularization

The simplest regression model, and one of the most widely used, is the linear regression model. In linear regression we want to learn a plane (\mathbf{w}, b) that describes the data as well as possible. The idea here is that if we can show that such a plane is a good representation of the relation between $y(\mathbf{x}_n)$ and \mathbf{x}_n for the training data we have, then it might give reliable predictions on new data \mathbf{x}_n

To do so, for a given set of training points, (\mathbf{x}_n, y_n) , we minimize the sum of squared errors between the predictions $y(\mathbf{x}_n)$ from the plane and the exact target values t_n ,

$$\underset{\mathbf{w}, b}{\operatorname{argmin}} \sum_{n=1}^N |\hat{y}(\mathbf{x}^n; \mathbf{w}, b) - t^n|^2 \quad (107)$$

$$= \underset{\mathbf{w}, b}{\operatorname{argmin}} \sum_{n=1}^N |\mathbf{w}^T \mathbf{x}^n + b - t^n|^2 \quad (108)$$

$$= \underset{\mathbf{w}, b}{\operatorname{argmin}} \sum_{n=1}^N \left| \sum_{i=1}^D \mathbf{x}_i^n \mathbf{w}_i + b - t_n \right|^2 \quad (109)$$

The model remains linear if we replace the original points by a representation in some feature space, $\phi(\mathbf{x})$. In this case, the model reads as

$$y(\mathbf{x}^n; \mathbf{w}; b) = b + \sum_{j=1}^D w_j \phi_j(\mathbf{x}^n) \quad (110)$$

and, including the bias b in the weight vector \mathbf{w} , $\tilde{\mathbf{w}} = [1, \mathbf{w}]$ and letting $\tilde{\phi}(\mathbf{x}) = [1, \phi(\mathbf{x})]$, we have

$$\underset{\mathbf{w}, b}{\operatorname{argmin}} \sum_{n=1}^N |y(\mathbf{x}^n; \tilde{\mathbf{w}}) - t^n|^2 = \sum_{n=1}^N |\mathbf{w}^T \phi(\mathbf{x}_n) - t_n|^2 \quad (111)$$

The solution to problem (111) can be computed in closed form by setting the derivatives to 0 and then solving for \mathbf{w} and b . Setting the derivative of (107) (on the generic feature formulation) to 0, we get

$$\sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n)^T = 0 \quad (112)$$

(Verify this using the matrix derivatives). When solving for \mathbf{w} , we thus have

$$\sum_{n=1}^N t_n \phi(\mathbf{x}_n)^T - \mathbf{w}^T \left(\sum_{n=1}^N \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_n) \right) = 0 \quad (113)$$

which gives

$$\mathbf{w} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad (114)$$

Here \mathbf{t} is the vector concatenating the targets $\mathbf{t} = (t_1, \dots, t_n)$, and Φ is the matrix whose columns encode the feature vectors $\phi(\mathbf{x}_n)$, i.e.

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_0) & \phi_1(\mathbf{x}_0) & \dots & \phi_M(\mathbf{x}_0) \\ \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_M(\mathbf{x}_1) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \dots & \phi_M(\mathbf{x}_N) \end{pmatrix} \quad (115)$$

The matrix $\Phi^+ = (\Phi^T \Phi)^{-1} \Phi^T$ is precisely the [Moore Penrose Pseudo inverse](#) which was introduced in Definition 1.

This solution highlights an important aspect of linear regression: For the solution of problem (111) to be well defined, we need the Gram matrix $\Phi^T \Phi$ to be invertible.

When this matrix is not invertible, it means that columns or rows, of that matrix are linearly dependent and hence, that a given subset of the feature vectors $\phi(\mathbf{x}_n)$ can be written from the

knowledge of all the remaining features. The issue with such redundancy in the representation of the data is that one can define a perfectly valid classifier for the dataset, but for which the weights can vary arbitrarily and in particular, be arbitrarily high as they can cancel each other. Indeed, assume that for all training points, features $\phi_1(\mathbf{x}_n)$ can be obtained as the combination $\phi_2(\mathbf{x}_n) + 2\phi_3(\mathbf{x}_n)$. Then for any given classifier,

$$y(\mathbf{x}) = b + w_1\phi_1 + w_2\phi_2 + w_3\phi_3 \quad (116)$$

adding a term

$$\alpha w\phi_1 - w(\phi_2 + 2\phi_3) \quad (117)$$

will not change the regression of the training points, for any $w \in \mathbb{R}$, since at all those training points we have

$$\phi_1(\mathbf{x}_n) = \phi_2(\mathbf{x}_n) + 2\phi_3(\mathbf{x}_n) \quad (118)$$

In particular, we could take w arbitrarily large and get a classifier that "fits" the training data perfectly well. The issue however is that for new data (or test data), the two classifiers will give very different results. In practice, we therefore want to avoid such situations and we will add a penalty on the coefficients (w_1, \dots, w_N) and bias b .

8.1 Regularizers

There are three popular approaches at regularizing the linear regression problem (111).

- ℓ_2 (Ridge regression). Here we simply minimize the sum of squared weights (squared ℓ_2 norm),

$$\operatorname{argmin}_{\mathbf{w}, b} \sum_{n=1}^N |\mathbf{w}^T \phi(\mathbf{x}_n) - t_n|^2 + \lambda \sum_{i=1}^D |\mathbf{w}_i|^2 \quad (119)$$

- ℓ_1 (LASSO). Here we minimize the sum of the absolute value of the weights,

$$\operatorname{argmin}_{\mathbf{w}, b} \sum_{n=1}^N |\mathbf{w}^T \phi(\mathbf{x}_n) - t_n|^2 + \lambda \sum_{i=1}^D |w_i| \quad (120)$$

- Finally in Best Subset Selection (which can be thought of as a minimization of the number of non zero weights), we found for each K , the best size- K subset of regression coefficients w_1, \dots, w_N .

In all of those approaches, the optimal choice of λ or for best subset selection, the optimal size of the subset of regression coefficients, is fixed through cross validation, by computing the regression coefficients associated to a particular value of λ for a subset \mathcal{S}_1 of the whole dataset, then testing the model on the remaining $\mathcal{S} \setminus \mathcal{S}_1$ data. And then selecting the λ that leads to the smallest generalization error. The two formulations (119) and (120) both admit constrained variations,

$$\operatorname{argmin}_{\mathbf{w}, b} \sum_{n=1}^N |\mathbf{w}^T \phi(\mathbf{x}_n) - t_n|^2 \quad (121)$$

$$\text{subject to } \sum_{i=1}^D |\mathbf{w}_i|^2 \leq t \quad (122)$$

and

$$\underset{\mathbf{w}, b}{\operatorname{argmin}} \quad \sum_{n=1}^N |\mathbf{w}^T \phi(\mathbf{x}_n) - t_n|^2 \quad (123)$$

$$\text{subject to} \quad \sum_{i=1}^D |\mathbf{w}_i| \leq t \quad (124)$$

respectively.

Finally, note that other regularization terms are possible. In particular, any ℓ_p norm, $p \geq 1$

$$\|\mathbf{x}\|_p = \left(\sum_{i=1}^D w_i^p \right)^{1/p} \quad (125)$$

can be used as a regularizer.

In the case of the ridge regression formulation (119), it remains possible to derive a closed form solution as follows. We consider the unconstrained formulation, take the derivatives with respect to (\mathbf{w}, b) and set those derivatives to 0, we get

$$\mathbf{w}^{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{t} \quad (126)$$

or in the case of feature vectors $\phi(\mathbf{x})$,

$$\mathbf{w}^{\text{ridge}} = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{t} \quad (127)$$

Hence you see that we have replaced the matrix $(\Phi^T \Phi)$ by a better conditioned version of this matrix as it is now combined to the identity which is invertible.

9 Bias Variance tradeoff

9.1 Expected Loss

Suppose our data set is given by $\mathcal{D} = \{(x, t)\}$ where t are the labels associated to each point \mathbf{x} . We will assume that the data is distributed according to an ideal function $y(\mathbf{x})$ plus some perturbations.

$$t = y(\mathbf{x}) + \varepsilon \quad (128)$$

t is thus the observed label. In order to study how well a given regression model is fitting the data, one can look at the expected loss

$$\mathbb{E}[L] = \int \int (h(\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt \quad (129)$$

Note that given the objective (129), the optimal model we can choose for $h(\mathbf{x})$ is to solve the minimization problem

$$\min_{h(\mathbf{x})} \int \int (h(\mathbf{x}) - t)^2 d\mathbf{x} dt \quad (130)$$

To do this, we set the first order derivatives to 0, and get

$$\frac{\delta \mathbb{E} L}{\delta h(\mathbf{x})} = 2 \int (h(\mathbf{x}) - t) p(\mathbf{x}, t) d\mathbf{x} dt \quad (131)$$

from this we get

$$\int h(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} = \int t p(t, \mathbf{x}) d\mathbf{x} dt \quad (132)$$

The t disappeared in the left handside because $h(\mathbf{x})$ does not depend on t so we can simply integrate over all values of t .

This last equation tells us that in general we should at every value \mathbf{x} choose the model $h(\mathbf{x})$ corresponding to

$$h(\mathbf{x}) = \frac{\int t p(\mathbf{x}, t) d\mathbf{x} dt}{p(\mathbf{x})} = \int t p(t|\mathbf{x}) dt = \mathbb{E}_t \{t|\mathbf{x}\} \quad (133)$$

the last line follows from $p(\mathbf{x}, t) = p(t|\mathbf{x})p(\mathbf{x})$. This is however not always possible.

On the other hand, we always have

$$(h(\mathbf{x}) - t)^2 = (h(\mathbf{x}) - \mathbb{E} \{t|\mathbf{x}\} + \mathbb{E} \{t|\mathbf{x}\} - t)^2 \quad (134)$$

$$= (h(\mathbf{x}) - \mathbb{E} \{t|\mathbf{x}\})^2 + (\mathbb{E} \{t|\mathbf{x}\} - t)^2 + 2 (h(\mathbf{x}) - \mathbb{E} \{t|\mathbf{x}\}) (\mathbb{E} \{t|\mathbf{x}\} - t) \quad (135)$$

If you substitute this last expression into the expected Loss (129) you can see that the cross terms will disappear as it leads to the following four terms

$$2 \int h(\mathbf{x}) \mathbb{E} \{t|\mathbf{x}\} p(\mathbf{x}, t) d\mathbf{x} dt = 2 \int h(\mathbf{x}) \mathbb{E} \{t|\mathbf{x}\} p(\mathbf{x}) d\mathbf{x} \quad (136)$$

$$- 2 \int h(\mathbf{x}) t p(\mathbf{x}, t) d\mathbf{x} dt = -2 \int \mathbb{E} \{t|\mathbf{x}\} h(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \quad (137)$$

$$- 2 \int \mathbb{E} \{t|\mathbf{x}\}^2 p(\mathbf{x}, t) d\mathbf{x} dt = -2 \int \mathbb{E} \{t|\mathbf{x}\}^2 p(\mathbf{x}) d\mathbf{x} \quad (138)$$

$$2 \int \mathbb{E} \{t|\mathbf{x}\} t p(\mathbf{x}, t) d\mathbf{x} dt = 2 \int \mathbb{E} \{t|\mathbf{x}\}^2 p(\mathbf{x}) d\mathbf{x} \quad (139)$$

So that the average loss reduces to

$$\mathbb{E}[L] = \int (t - h(\mathbf{x}))^2 d\mathbf{x} dt = \int (h(\mathbf{x}) - \mathbb{E} \{t|\mathbf{x}\})^2 p(\mathbf{x}) d\mathbf{x} \quad (140)$$

$$+ \int (\mathbb{E} \{t|\mathbf{x}\} - t)^2 dt p(\mathbf{x}) d\mathbf{x} \quad (141)$$

$$= \int (h(\mathbf{x}) - \mathbb{E} \{t|\mathbf{x}\})^2 p(\mathbf{x}) d\mathbf{x} \quad (142)$$

$$+ \int \sigma_x^2 p(\mathbf{x}) d\mathbf{x} \quad (143)$$

Where I used

$$\int (\mathbb{E} \{t|\mathbf{x}\} - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt = \int \text{Var}[t|\mathbf{x}] p(\mathbf{x}) d\mathbf{x} \quad (144)$$

Once again, the t is integrated over in the first integral because none of the functions $h(\mathbf{x})$ and $\mathbb{E}\{t|\mathbf{x}\}$ depend on t anymore.

The final expression of the average loss is then

$$\mathbb{E}[L] = \int (h(\mathbf{x}) - \mathbb{E}\{t|\mathbf{x}\})^2 p(\mathbf{x}) d\mathbf{x} + \int \sigma_x^2 p(\mathbf{x}) d\mathbf{x} \quad (145)$$

You see that there is a first term (MSE) that depends on the model and a second term that depends only on the data. Now we will show that the first term can be decomposed into a *Variance* and a *bias* contributions.

9.2 Bias Variance Tradeoff

In the derivations above, we haven't made any assumption on the regression model learned $h(\mathbf{x})$. In general, when we learn a regression model, we learn it from a subset of the data. We will denote this subset \mathcal{D}_i so that $\mathcal{D}_i \subset \mathcal{D}$ and I will now denote our regression model as $h(\mathbf{x}|\mathcal{D}_i)$ to emphasize the dependence on the particular choice of \mathcal{D}_i . Now that being said, if we use $\bar{t}(x)$ to denote $\mathbb{E}\{t|\mathbf{x}\}$, we have the decomposition

$$(h(\mathbf{x}; \mathcal{D}_i) - \bar{t}(x))^2 \quad (146)$$

$$= \left(h(\mathbf{x}; \mathcal{D}_i) - \mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} + \mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} - \bar{t}(x) \right)^2 \quad (147)$$

$$= \left(h(\mathbf{x}; \mathcal{D}_i) - \mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} \right)^2 + \left(\mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} - \bar{t}(x) \right)^2 \quad (148)$$

$$+ 2 \left(h(\mathbf{x}; \mathcal{D}_i) - \mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} \right) \left(\mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} - \bar{t}(x) \right) \quad (149)$$

If we take the expectation over all possible datasets \mathcal{D}_i , we will get a measure on how well the model $h(\mathbf{x})$ performs on average for the different datasets. Below I use $\mathbb{E}_{\mathcal{D}_i}$ to denote the average when sample random datasets \mathcal{D}_i from the full set of samples \mathcal{D} .

$$\mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i) - \bar{t}(x)\}^2 \quad (150)$$

$$= \mathbb{E}_{\mathcal{D}_i} \left(h(\mathbf{x}; \mathcal{D}_i) - \mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} \right)^2 + \mathbb{E}_{\mathcal{D}_i} \left(\mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} - \bar{t}(x) \right)^2 \quad (151)$$

$$+ 2 \mathbb{E}_{\mathcal{D}_i} \left(h(\mathbf{x}; \mathcal{D}_i) - \mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} \right) \left(\mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} - \bar{t}(x) \right) \quad (152)$$

$$= \underbrace{\left(\mathbb{E}_{\mathcal{D}} \{y(\mathbf{x}; \mathcal{D})\} - \bar{t}(x) \right)^2}_{\text{bias}^2} + \underbrace{\mathbb{E}_{\mathcal{D}} \left[(h(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}} \{h(\mathbf{x}; \mathcal{D})\})^2 \right]}_{\text{variance}} \quad (153)$$

What this tells you in terms of the complexity of the regression models is that complex models will usually have a large variance because within the class of complex models, the models will usually *vary* a lot. I.e if you have a very complicated model with many many parameters that perfectly fits the data in \mathcal{D}_i , then when you will take another dataset \mathcal{D}_j , the model that you build for \mathcal{D}_j will be completely different from the model you built with \mathcal{D}_i . So complex models usually have larger variance.

The bias, on the other hand is a measure of how well the model "captures" the average behavior of the data $\mathbb{E}\{t|\mathbf{x}\}$. In general (except if the data is assumed to be linear which is the framework of the Gauss Markov Theorem discussed below), the simpler models will not be able to capture such behavior. As an example, imagine we have data of the form

$$t = y(x) + \varepsilon = \sin(x) + \varepsilon \quad (154)$$

	simpler models	complex models
bias	large	small
variance	small	large

Table 2: The bias variance trade-off and the model complexity

Then without using non linear terms (and defining our new variables as $x' = \phi(x)$ where $\phi(x)$ encodes the non linearity), it is impossible to design a model of the form $h(x) = \beta_0 + \langle \beta_1, x \rangle$ such that

$$h(\mathbf{x}) = \sin(x) \quad (155)$$

This should illustrate the fact that simpler models have a large bias (When we don't make assumptions on the data distribution)

In general we thus have the relation summarized in table 2.

9.3 The particular case of linearly generated data and the Gauss Markov Theorem

So far we haven't made any assumptions on the data. When we don't make assumptions on the data, the relation of table 2 holds in general.

However, when we know that the data has a particular distribution, we can better describe the quality of approximation of given models. This is the framework of the *Gauss Markov Theorem* (GMT).

When the data is distributed according to a linear model, so that

$$t = y(x) + \varepsilon = \alpha_0 + \alpha_1 x + \varepsilon, \quad (156)$$

then the linear model has no bias (note that this might seem counterintuitive compared to table 2 but this is because here we assumed that the data is linear).

Moreover, the *Gauss Markov theorem* states that among all possible linear estimators of the α_0 and α_1 appearing in (for linear data), the linear regression model, which is computed by minimizing the residual sum of squares,

$$\underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N (t_i - \beta_0 - \beta_1 x_i)^2 \quad (157)$$

is the best estimator. In the Gauss Markov framework, when we say data, we mean the t_i . What we now call data are the noisy measurements $t_i = \beta_0 + \beta_1 x_i + \varepsilon$. Because we've put ourselves in a framework where we now know how the data is distributed, $t_i(x) = \beta_0 + \beta_1 x_i + \varepsilon$, t_i is enough to recover the data when there is no noise and we thus call t_i the data. In this particular framework, a linear estimator for the coefficients β is an estimator of the form

$$\tilde{\beta} = M\mathbf{t} \quad (158)$$

where $\mathbf{t} = [t_1, t_2, \dots, t_N]$ encode the labels from \mathcal{D}_i . The estimator provided by the linear regression approach (157) is a particular such estimator as it can read as

$$\hat{\beta}_{\text{RSS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \quad (159)$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\beta + \varepsilon) \quad (160)$$

To prove the GMT, we consider any other linear estimator of the coefficients $\beta = \beta_0, \beta_1$, $\tilde{\beta} = \hat{\beta}_{\text{RSS}} + \Delta \mathbf{t}$ where Δ is a perturbation.

We first require the linear estimator $\tilde{\beta}$ to be unbiased, which gives

$$\mathbb{E} \{ \tilde{\beta} \} = \mathbb{E} \{ \hat{\beta} \} + \mathbb{E} \{ \Delta \mathbf{t} \} \quad (161)$$

$$= \mathbb{E} \{ \Delta \mathbf{X} \beta + \varepsilon \Delta \} \quad (162)$$

$$= \Delta \mathbf{X} \beta + 0 \quad (163)$$

$$= 0 \quad (164)$$

The only source of randomness here is the noise ε and we assumed that $\mathbb{E} \varepsilon = 0$. As we haven't made any assumption on β , the last line implies that for the linear estimator to be unbiased, we must have

$$\Delta \mathbf{X} = 0 \quad (165)$$

Now, for the variance, first note that

$$\text{Var} \{ \hat{\beta} \} = \mathbb{E} \{ (\hat{\beta} - \mathbb{E} \hat{\beta})^2 \} \quad (166)$$

$$= \mathbb{E} \{ \hat{\beta}^2 \} - \beta^2 \quad (167)$$

as the estimator is unbiased. For any other unbiased estimator, $\tilde{\beta}$ of the form $\tilde{\beta} = \mathbb{E} \{$

$$\text{Var} \{ \tilde{\beta} \} \quad (168)$$

$$= \mathbb{E} \{ (\tilde{\beta} - \mathbb{E} \tilde{\beta})(\tilde{\beta} - \mathbb{E} \tilde{\beta})^T \} \quad (169)$$

$$= \mathbb{E} \left\{ \left(\hat{\beta}_{\text{RSS}} + \Delta \mathbf{t} - \beta - \mathbb{E} \{ \Delta \mathbf{t} \} \right) \left(\hat{\beta}_{\text{RSS}} + \Delta \mathbf{t} - \beta - \mathbb{E} \{ \Delta \mathbf{t} \} \right)^T \right\} \quad (170)$$

$$= \mathbb{E} \left\{ \left(\hat{\beta}_{\text{RSS}} - \beta \right) \left(\hat{\beta}_{\text{RSS}} - \beta \right)^T \right\} \quad (171)$$

$$+ \mathbb{E} \{ (\Delta \mathbf{t} - \mathbb{E} \{ \Delta \mathbf{t} \})(\Delta \mathbf{t} - \mathbb{E} \{ \Delta \mathbf{t} \})^T \} \quad (172)$$

$$- \mathbb{E} \{ (\hat{\beta}_{\text{RSS}} - \beta)(\Delta \mathbf{t} - \mathbb{E} \{ \Delta \mathbf{t} \})^T \} \quad (173)$$

$$- \mathbb{E} \{ (\Delta \mathbf{t} - \mathbb{E} \{ \Delta \mathbf{t} \})(\hat{\beta}_{\text{RSS}} - \beta)^T \} \quad (174)$$

As $\hat{\beta}_{\text{RSS}}$ is unbiased, the last two (cross) terms disappear and we are left with

$$\text{Var} \{ \tilde{\beta} \} = \text{Var} \{ \hat{\beta}_{\text{RSS}} \} + \mathbb{E} \{ (\Delta \mathbf{t})(\Delta \mathbf{t})^T \} \quad (175)$$

$$- \mathbb{E} \{ (\Delta \mathbf{t}) \} \mathbb{E} \{ (\Delta \mathbf{t}) \}^T \quad (176)$$

$$= \text{Var} \{ \hat{\beta}_{\text{RSS}} \} + \Delta \mathbf{X} \beta \beta^T \mathbf{X}^T \Delta^T \quad (177)$$

$$+ \Delta \varepsilon \varepsilon^T \Delta^T \quad (178)$$

$$= \text{Var} \{ \hat{\beta}_{\text{RSS}} \} + \sigma^2 \Delta \Delta^T \quad (179)$$

As soon as $\Delta \neq 0$, the diagonal of $\Delta \Delta^T$ is non zero as well. This in particular means that there is one component of $\tilde{\beta}$ for which we always have

$$\mathbb{E} \{ (\tilde{\beta}_k - \beta_k)^2 \} > \mathbb{E} \{ (\hat{\beta}_k - \beta_k)^2 \} \quad (180)$$

Both estimators have zero bias (we will overestimate and underestimate the truth equally) but the errors we make will often be larger when we don't use the RSS estimator.

In conclusion, you can thus see that any linear estimator will only increase the variance with respect to $\hat{\beta}_{\text{RSS}}$ (This is because we are in the particular framework of linearly generated data)

For this reason, the estimator $\hat{\beta}_{\text{RSS}}$ is sometimes called the BLUE (Best Linear Unbiased Estimator)

9.3.1 Increasing the bias to reduce generalization under the assumption of linearly generated data

That being said, the fact that the BLUE is the best estimator among all unbiased estimators with respect to linearly generated data, does not mean that it is not possible to do better in terms of future predictions. And this is because a similar bias variance tradeoff as the one discussed in section 9.2 applies.

Imagine that we take a new point x_0 and we apply our estimator $\tilde{\beta}$ (*now we don't assume that this estimator is unbiased anymore*, this is just a general linear estimator) to this point to get an estimated label y , we can then write the expected mean square prediction error as before

$$\mathbb{E}_{\varepsilon} \left\{ \left(\langle \tilde{\beta}, x_0 \rangle - \langle \beta, x_0 \rangle \right)^2 \right\} \quad (181)$$

$$= \mathbb{E}_{\varepsilon} \left\{ \left(\langle \tilde{\beta}, x \rangle - \mathbb{E} \left\{ \langle \tilde{\beta}, x \rangle \right\} + \mathbb{E} \left\{ \langle \tilde{\beta}, x \rangle \right\} - \langle \beta, x \rangle \right)^2 \right\} \quad (182)$$

$$= \mathbb{E}_{\varepsilon} \left\{ \left(\langle \tilde{\beta}, x \rangle - \mathbb{E} \left\{ \langle \tilde{\beta}, x \rangle \right\} \right)^2 \right\} \quad (183)$$

$$+ \mathbb{E}_{\varepsilon} \left\{ \left(\mathbb{E} \left\{ \langle \tilde{\beta}, x \rangle \right\} - \langle \beta, x \rangle \right)^2 \right\} \quad (184)$$

$$+ 2\mathbb{E}_{\varepsilon} \left\{ \left(\langle \tilde{\beta}, x_0 \rangle - \mathbb{E}_{\varepsilon} \left\{ \tilde{\beta}, x_0 \right\} \right) \left(\mathbb{E} \left\{ \langle \tilde{\beta}, x_0 \rangle \right\} - \langle \beta, x_0 \rangle \right) \right\} \quad (185)$$

The last term vanishes as the second factor is deterministic and the first factor has mean 0. We thus once again have

$$\mathbb{E}_{\varepsilon} \left\{ \left(\langle \tilde{\beta}, x_0 \rangle - \langle \beta, x_0 \rangle \right)^2 \right\} = \underbrace{\mathbb{E}_{\varepsilon} \left\{ \left(\langle \tilde{\beta}, x \rangle - \mathbb{E} \left\{ \langle \tilde{\beta}, x \rangle \right\} \right)^2 \right\}}_{\text{variance}} \quad (186)$$

$$+ \underbrace{\mathbb{E}_{\varepsilon} \left\{ \left(\mathbb{E} \left\{ \langle \tilde{\beta}, x \rangle \right\} - \langle \beta, x \rangle \right)^2 \right\}}_{\text{bias}^2} \quad (187)$$

This shows that even in the framework of linearly generated data, despite the fact that $\hat{\beta}$ is the BLUE estimator, it does not mean that this estimator will be the best (on average over the noise) at predicting a new value of $t = \langle x, \beta \rangle$. There might be an estimator $\tilde{\beta}$ with a non zero bias (on linear data) that will have smaller variance.

That is to say, if I generate many values $t_k = \langle \beta, x_k \rangle + \varepsilon_k$, learn my estimator $\hat{\beta}$ from those values, then this estimator will be unbiased but when I will want to get the prediction for a new point x_0 , the average prediction might be better if I lower the variance and increase the bias a little. Again this might sound counter intuitive with respect to table 2 but remember that we assumed linearly distributed data here.

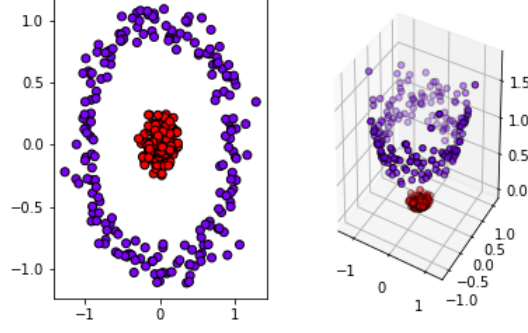


Figure 2: From non linearly separable data to linearly separable data via higher dimensional feature space

10 Kernels

Let $\beta_0 \in \mathbb{R}$ to denote the bias and $\beta_1 \in \mathbb{R}^D$ to denote the vector of weights. In practice, we often have to deal with prototypes, $\{\mathbf{x}_\mu\}_{\mu=1}^N$ which are not linearly separable in their original space. One approach then consists in introducing features and to "map" the original prototypes \mathbf{x}_μ into a space (the feature space) in which those prototypes become linearly separable. We use ϕ to denote the underlying transformation so that for a given datapoint \mathbf{x} , $\phi(\mathbf{x})$ denotes the **feature vector** of \mathbf{x} .

Example 1. Consider the dataset shown in Fig. 2 (left). In the original \mathbb{R}^2 space, the data given by pairs (x_1, x_2) is clearly not linearly separable. However, it is possible to introduce a transformation $\phi : \mathbf{x} \mapsto \phi(\mathbf{x})$ defined as $\phi(\mathbf{x}) = \phi(x_1, x_2) = (x_1, x_2, x_1^2 + x_2^2)$.

Here, introducing an additional dimension, and setting this dimension to be equal to the radius of the points in the original space, results in the purple points being placed above the red points, thus leaving some space for a separating plane.

Relying on feature vectors to learn classifiers is at the core of machine learning. However, such an approach often requires to map the data to a much higher dimensional space (i.e higher than the number of prototypes \mathbf{x}_μ). In this case, it is more interesting to avoid computing the feature vectors explicitly and to rely instead on a measure of similarity between the points in the feature space. After all, the important information here is really how the points compare to each other in the feature space. This information is encoded through **kernels**

Mathematically, a kernel is nothing else than a function that measure similarity between points. This is summarized by the definition below

Definition 19 (Kernel). We define a kernel to be the real valued function of two arguments \mathbf{x} and \mathbf{x}' from a space \mathcal{X} . Typically this function is taken to be symmetric $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$ and non negative, $k(\mathbf{x}, \mathbf{x}') \geq 0$. Those two properties enable us to interpret the kernel $k(\mathbf{x}, \mathbf{x}')$ as a measure of similarity between points.

Examples of popular kernels are listed below

- The Gaussian kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^\top \Sigma^{-1}(\mathbf{x} - \mathbf{x}')\right)$$

- When the covariance matrix Σ is isotropic/spherical, we get the isotropic kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$$

- The cosine similarity:

$$k(\mathbf{x}, \mathbf{x}') = \frac{\mathbf{x}^T \mathbf{x}'}{\|\mathbf{x}\|_2 \|\mathbf{x}'\|_2}$$

- The (non stationnary) polynonial kernel:

$$k(\mathbf{x}, \mathbf{x}') = (\gamma \mathbf{x}^T \mathbf{x}' + r)^M$$

where $r > 0$. This kernel corresponds to a feature vector $\phi(\mathbf{x})$ that contains all the monomials up to degree M .

The isotropic kernel is an example of **Radial basis function** (a function whose distance decreases or increases monotonically with respect to a central point). Those functions are multivariate (i.e they are defined on $\mathbf{x} \in \mathbb{R}^D$) but they reduce to a scalar function of the Euclidean norm $\|\mathbf{x}\|^2$ of their argument \mathbf{x} , i.e.,

$$F(\mathbf{x}) = f(\|\mathbf{x}\|_2) = \phi(r), \quad \mathbf{x} \in \mathbb{R}^D \quad (188)$$

Examples of radial basis functions include

- The Gaussian RBF: $F(\|\mathbf{x} - \mathbf{x}'\|) = \exp(-\alpha^2 \|\mathbf{x} - \mathbf{x}'\|^2)$
- The Multiquadratic RBF: $F(\|\mathbf{x} - \mathbf{x}'\|) = \sqrt{1 + \alpha^2 \|\mathbf{x} - \mathbf{x}'\|^2}$

As soon as one has access to the feature vector $\phi(\mathbf{x})$ (for example when this feature is finite dimensional), one can build a valid kernel by defining $k(\mathbf{x}, \mathbf{x}')$ as the inner product of the feature vectors of \mathbf{x} and \mathbf{x}' , i.e.,

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') \quad (189)$$

10.1 From linear regression to kernel regression

When the dataset is not linearly separable in the original space, we turn to a formulation on the feature space and try to find the plane defined by (β_0, β_1) that minimizes $J(\beta)$, i.e.,

$$\min J(\beta) = \frac{1}{2} \sum_{n=1}^N \left\{ \beta_1^T \phi(\mathbf{x}_n) + \beta_0 - t_n \right\}^2 + \frac{\lambda}{2} \|\beta\|^2 \quad (190)$$

Here the separating plane is defined through the vector of parameters $\beta = (\beta_0, \beta_1)$, $\beta_0 \in \mathbb{R}$, $\beta_1 \in \mathbb{R}^D$. The solution for β can be obtained by computing the derivative of $J(\beta)$ with respect to β and setting it to zero. For pairs $\{(\mathbf{x}_n, t_n)\}_{n=1}^N$, we have

$$\frac{\partial J(\beta)}{\partial \beta} = 0 \Rightarrow \beta = -\frac{1}{\lambda} \sum_{n=1}^N \left\{ \beta^T \tilde{\phi}(\mathbf{x}_n) - t_n \right\} \tilde{\phi}(\mathbf{x}_n) \quad (191)$$

Now introducing the notation a_n for the quantities

$$a_n = -\frac{1}{\lambda} \left\{ \beta^T \tilde{\phi}(\mathbf{x}_n) - t_n \right\} \quad (192)$$

Then the solution for β can read as

$$\beta = \sum_{n=1}^N a_n \phi(\mathbf{x}_n) = \Phi^T \mathbf{a} \quad (193)$$

Where Φ is the matrix whose i^{th} row is given by $\phi^T(x_i)$. Now substituting this expression into (191), we get

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a} \quad (194)$$

From the discussion above, you see that the product $\Phi \Phi^T$ is actually encoding all inner products $\langle \phi(\mathbf{x}_n), \phi(\mathbf{x}_m) \rangle$ and we can thus replace the product $\Phi \Phi^T$ with the Kernel $K(\mathbf{x}, \mathbf{x}')$. Doing this yields an objective, or energy function that only depends on the kernel $K(\mathbf{x}, \mathbf{x}')$,

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a} \quad (195)$$

We can then solve for \mathbf{a} or simply substitute the expression (193) we found for β in (192), to find

$$\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{t} \quad (196)$$

Once we have this expression (note that so far everything can be expressed with the Kernel), the classifier simply reads as

$$y(\mathbf{x}) = \beta^T \phi(\mathbf{x}) = \mathbf{a}^T \Phi \phi(\mathbf{x}) = \sum_{n=1}^N a_n K(\mathbf{x}_n, \mathbf{x}) \quad (197)$$

11 Support Vector Machines (SVM)

Support vector machine, which are also known as sparse kernel machines, or Maximal margin classifiers, extend the notion of separating hyperplane to find a separating plane that reduces as much as possible the "risk" of misclassifying new points. To achieve this, the separating plane is defined as the plane maximizing the distance to its closest point(s).

Using the discussion from section 7.1, we can then write the derivation of such maximal margin hyperplane as

$$\operatorname{argmax}_{\mathbf{w}, \mathbf{b}} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + \mathbf{b})] \right\} \quad (198)$$

That is we look for the \mathbf{w} and b (which define the plane) such that the distance of the closest point (minimization on \mathbf{x}_n) to the plane is maximized.

Formulation (198) is invariant under any rescaling of the pair (\mathbf{w}, \mathbf{b}) . I.e as \mathbf{w} and/or b both appear at the numerator and at the denominator, replacing those quantities by $(\alpha \mathbf{w}, \alpha \mathbf{b})$ does not affect the solution. In fact any such pair, for any α will be a valid solution.

In order to simplify the formulation, we can therefore fix the value of α and decide to take this value equal to the value satisfying

$$t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + \mathbf{b}) = 1 \quad (199)$$

for the closest point \mathbf{x}_n . By assumption, any other point, must satisfy $t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + \mathbf{b}) \geq 1$, and the problem reads as

$$\begin{aligned} & \operatorname{argmin}_{\mathbf{w}, \mathbf{b}} \quad \frac{1}{2} \|\mathbf{w}\|^2 \\ & \text{subject to} \quad t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + \mathbf{b}) \geq 1 \end{aligned} \quad (200)$$

Note that we have used $\operatorname{argmax}_x \|\mathbf{x}\|^{-1} = \operatorname{argmin}_x \|\mathbf{x}\|$. Problem (200) is a quadratic program that is convex, hence, has a single basin of attraction and can be solved efficiently.

To solve this problem, we write it as an unconstrained formulation (Lagrangian) by introducing multipliers, λ_i for each constraints (see section 6.2)

$$L(\mathbf{w}, \mathbf{b}, \lambda) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N \lambda_n \{t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + \mathbf{b}) - 1\} \quad (201)$$

To find the solutions for \mathbf{w} and b , we compute the derivatives of L with respect to those variables and set the derivatives to 0. From this, we get

$$\mathbf{w} = \sum_{n=1}^N \lambda_n t_n \phi(\mathbf{x}_n) \quad (202)$$

$$0 = \sum_{n=1}^N \lambda_n t_n \quad (203)$$

Now recall that the general form of the classifier we are after is given by

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + \mathbf{b} \quad (204)$$

If we substitute the expression for \mathbf{w} that we derived from the zero derivatives in the expression of the classifier, we get

$$\begin{aligned} y(\mathbf{x}) &= \sum_{n=1}^N \lambda_n t_n \langle \phi(\mathbf{x}_n), \phi(\mathbf{x}) \rangle + b \\ &= \sum_{n=1}^N \lambda_n t_n k(\mathbf{x}, \mathbf{x}_n) + b \end{aligned} \quad (205)$$

Now besides the maximization of the margins, the other key idea of SVM is the fact that we don't need to keep track of all the training samples in (205). Indeed from Theorem 10, we know that the solution (\mathbf{w}, \mathbf{b}) of problem (200) has to satisfy the conditions

$$\lambda_n \geq 0 \quad (206)$$

$$t_n y(\mathbf{x}_n) - 1 \geq 0 \quad (207)$$

$$\lambda_n (t_n y(\mathbf{x}_n) - 1) = 0 \quad (208)$$

From the last condition in particular, we can see that for every training sample, $n = 1, \dots, N$, either we have $\lambda_n = 0$ or we have $t_n y(\mathbf{x}) - 1 = 0$. The second case corresponds to points that are located the closest to the maximal margin hyperplane as the distance to each points to the hyperplane is always lower bounded by

$$\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} \geq \frac{1}{\|\mathbf{w}\|} \quad (209)$$

In other words, for every training point that appears in (205):

- Either that point does not contribute to the expression (205) of the classifier (as $\lambda_n = 0$)
- Or it is lying on the margins.

Using this idea to simplify (205), if we introduce the notation \mathcal{S} to denote the set of “support vectors”, that is the points that are lying on the margins, or located the closest to the plane, the expression of the classifier reduces to

$$y(\mathbf{x}) = \sum_{n \in \mathcal{S}} \lambda_n t_n k(\mathbf{x}, \mathbf{x}_n) + b \quad (210)$$

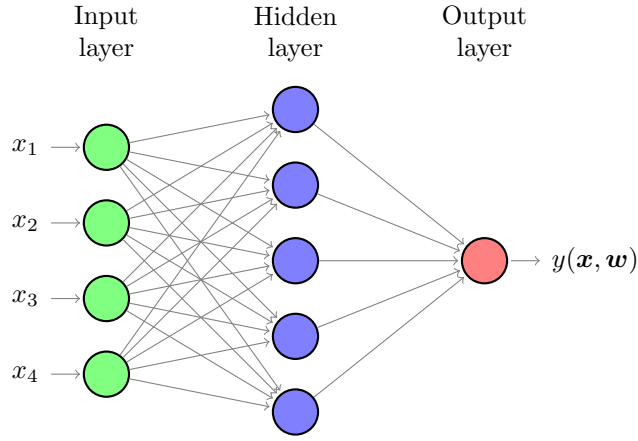


Figure 3: Graphical representation of a one layer neural network. The arrows represent multiplication by the weights w_{ij} . The hidden units (shown in blue) represent the application of the activation function.

Furthermore, we can use condition (208) for the support vectors, (i.e. those for which we have $t_n(y(\mathbf{x}_n) - 1) = 0$) to write b . indeed note that substituting (210) into condition (208) gives

$$t_n \left(\sum_{n \in \mathcal{S}} \lambda_n t_n k(\mathbf{x}, \mathbf{x}_n) + b \right) = 1, \quad \text{for every } \mathbf{x}_n \in \mathcal{S} \quad (211)$$

To obtain a robust estimate of b , we can then just sum all those equation and divide by N_S . This gives

$$b = \frac{1}{N_S} \sum_{n \in \mathcal{S}} \left(t_n - \sum_{n' \in \mathcal{S}} \lambda_{n'} t_{n'} k(\mathbf{x}', \mathbf{x}_n) \right) \quad (212)$$

12 Neural Networks

The general formulation of a (two hidden layers) neural network is as follows,

$$y_k = y_k(\mathbf{x}, \mathbf{w}) = \sigma^{(2)} \left(\sum_{j=1}^M w_{kj}^{(2)} \sigma^{(1)} \left(\sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right) \quad (213)$$

When training neural networks, we usually don't take all the training samples into account but rather rely on a small subset of those training samples at each gradient iteration. These subsets are known as [minibatches](#). There are three main approaches at training a network

- [Batch gradient descent](#). Here all the training samples are taken into account and the weight are updated only after all the error has been evaluated at all the training samples (Batch gradient descent is just the regular gradient descent)
- [Stochastic gradient descent \(SGD\)](#). Here, the algorithm computes the error for a given training sample and update the weights immediately after (i.e the weights are updated for each training sample)
- Finally, the tradeoff between batch gradient descent and SGD, [MiniBatch gradient descent](#) updates the weight after evaluating the error/loss on a subset of the training samples.

Among the possible regularization approaches for neural networks, the most popular are the following

- **Activity regularization.** For a given network, we call activations, the outputs of the activation functions, i.e $\sigma^{(\ell)}(\cdot)$. Placing a regularizer on the activation will result in training a network in which only a fraction of the neurons fire at the same time (at least for data similar to the training data).
- **Weight regularization.** Here we add a penalty (e.g. ℓ_2 or ℓ_1 norm) on the weights.
- **Early stopping.** When training neural networks, one often observe a decrease in both the training and validation error for the first iterations. Then, after a sufficient number of gradient iterations, the validation error starts increasing [6]. One way to avoid overfitting is then to get the values of the parameters at each iteration and return the values of the parameters that give the lowest validation error (we stop the iterations when no improvement on the validation error has been obtained for some iterations)
- **Dropout,** which shares some similarities with activity regularization in fact it can be considered as some sort of subset selection alternative to the ℓ_1 or ℓ_2 penalties used in activity regularization, can be considered as optimizing over a family of "sparse" networks. When training a network with dropout, we select some random subset of the neurons and only update the weights associated to those neurons. Mathematically, when performing SGD, we thus compute the forward pass (i.e the output of the network) by considering a sparse version of the original network [21]. We choose a probability p . Recall that for a traditional neural network, the activity of each unit can be defined as

$$z_i^{(\ell+1)} = \langle \mathbf{w}_i^{(\ell+1)}, \mathbf{y} \rangle + b_i^{(\ell+1)} \quad (214)$$

$$y_i^{(\ell+1)} = \sigma(z_i^{(\ell+1)}) \quad (215)$$

With dropout propagation through the network now reads as

$$r_j^{(\ell)} \sim \text{Bernoulli}(p) \quad (216)$$

$$\tilde{\mathbf{y}}^{(\ell)} = \mathbf{r}^{(\ell)} \odot \mathbf{y}^{(\ell)} \quad (217)$$

$$z_i^{(\ell+1)} = \mathbf{w}_i^{(\ell+1)} \tilde{\mathbf{y}}^{(\ell)} + b_i^{(\ell+1)} \quad (218)$$

$$y_i^{(\ell+1)} = f(z_i^{(\ell+1)}) \quad (219)$$

So that we select only some of the activations y_j and only update the weights associated to those activations.

- **Weight Sharing.** The idea of weight sharing is especially meaningful in convolutional neural networks (CNN) where filters are applied to an image to extract information from this image. As the name indicates, weights sharing requires all the neurons within a particular subgroup to have identical weights. Weights sharing is implicitly present in convolutional neural. In those network, each neuron can be thought as computing a local average across the neighborhood of a pixel. The point of such averages is to extract features from a part of the image. Usually, when using convolutional neural networks, we are interested in detecting a particular object in the input images. Hence the features we are looking for should be the same regardless of the part of the image on which the filter is applied. As a consequence, the weights are constants across each layers. Equality of the weights can be enforced in a hard (exact equality such as in CNNs) or in a soft way. Soft weight sharing is done by assuming a Gaussian distribution of the weights within a given group of neurons.

In the simplest framework, we assume that all the weights w_i in the group \mathcal{G} are all relatively close to some mean value μ (with variance σ^2) and we for example assume the distribution of those weights to be Gaussian

$$p(w_i) = \mathcal{N}(w_i | \mu, \sigma^2) \quad (220)$$

We can make this model slightly more general by assuming that the weights have a high probability to take a few distinct values. In this case, we assume for example that the

distribution of those weights is given by a mixture of Gaussian, i.e.,

$$p(w_i) = \sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2) \quad (221)$$

We then use this distribution as a regularization in the objective used to train the network. I.e, just as we did for the MAP and MLE estimators, we want the training step to select weights that have a high probability in the model (221). This corresponds to maximizing (221) over the weights, or (as we did for the MAP and MLE) to minimizing the negative log likelihood

$$\Omega(w) = -\log\left(\sum_{j=1}^M \mathcal{N}(w | \mu_j, \sigma_j^2)\right) \quad (222)$$

When there are several ($N_{\mathcal{G}}$) weights in the group \mathcal{G} , we simply sum over the contribution of each weight,

$$\Omega(\mathbf{w}) = -\sum_{i=1}^{N_{\mathcal{G}}} \log\left(\sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2)\right) \quad (223)$$

The loss function that is minimized when training the network is then given by

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \lambda \Omega(\mathbf{w}) \quad (224)$$

12.1 Backpropagation

Consider the two hidden layers NN of (213). The extension of this simpler model to L layers can read as

$$y_k = y_k(\mathbf{x}, \mathbf{w}) = \sigma^{(L)}\left(\sum_{j=1}^M w_{kj}^{(L)} \sigma^{(L-1)}\left(\dots \sigma^{(1)}\left(\sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)}\right) \dots\right) + w_{k0}^{(L)}\right) \quad (225)$$

In order to train such a network, we minimize the [loss](#)

$$\sum_{n=1}^N \sum_{k=1}^K (y_k(\mathbf{x}_n) - t_k^n)^2 \quad (226)$$

if we put ourselves in a stochastic gradient descent framework, where we only consider a single sample at each iteration, we can drop the sum over n in (226) and consider the minimization

$$L = \sum_{k=1}^K (y_k(\mathbf{x}_n) - t_k^n)^2 \quad (227)$$

In order to apply one SGD iteration on this objective, and to update the weights

$$w_{ij}^\ell \leftarrow w_{ij}^\ell - \eta \frac{\partial L}{\partial w_{ij}^\ell} \quad (228)$$

we need to compute the expression of the derivatives $\frac{\partial L}{\partial w_{ij}^\ell}$. To do so, we will rely extensively on the chain rule.

Computing the derivative with the outer weights (the weights that are the closest to the output of the network) is easy. But computing the derivatives with respect to the most inner weights, i.e. w_{ij}^1 , can be quite involved.

Backpropagation is an efficient way to compute the gradient of neural networks. First note that

$$\frac{\partial L}{\partial y_k} = 2(y_k - t_k) \quad (229)$$

We then want to compute all derivatives of the form

$$\frac{\partial L}{\partial w_{kj}^\ell} \quad (230)$$

For this, we use the chain rule a first time, splitting (230) into

$$\frac{\partial L}{\partial w_{kj}^\ell} = \frac{\partial L}{\partial a_k^\ell} \frac{\partial a_{kj}^\ell}{\partial w_{kj}^\ell} \quad (231)$$

that is we start from the weight we want to compute and we go up one step towards the output of the network by considering the quantities

$$a_k^\ell = \sum_{j=1}^M w_{k,j}^\ell z_j^{\ell-1} \quad (232)$$

where $z_j^{\ell-1} = \sigma^{\ell-1}(\dots)$.

Note that the second factor in (231) is simply $z_j^{\ell-1} = \sigma^{\ell-1}(\dots)$ which is computed when we pass the data through the network. To get the first factor, we apply the chain rule again, going one more level towards the output of the network, y_k . We have

$$\frac{\partial L}{\partial a_k^\ell} = \sum_j \frac{\partial L}{\partial a_j^{\ell+1}} \frac{\partial a_j^{\ell+1}}{\partial a_k^\ell} \quad (233)$$

This last equation is the equation that really encodes backpropagation. To see this, first note that from (225), we have

$$a_j^{\ell+1} = \sum_m w_{jm}^\ell z_m \quad (234)$$

$$= \sum_m w_{jm}^\ell \sigma^\ell(a_m) \quad (235)$$

and we can thus write

$$\frac{\partial a_j^{\ell+1}}{\partial a_k^\ell} = (\sigma^\ell)' w_{jm}^\ell \quad (236)$$

Putting this back into (233), we get

$$\frac{\partial L}{\partial a_k^\ell} = \sum_j \frac{\partial L}{\partial a_j^{\ell+1}} (\sigma^\ell)' w_{jm}^\ell \quad (237)$$

Now let $\delta_k^\ell = \frac{\partial L}{\partial a_k^\ell}$. From (237), we can write

$$\delta_k^\ell = \sum_j \delta_j^{\ell+1} (\sigma^\ell)' w_{jm}^\ell \quad (238)$$

From (238), you see that the $\delta^{\ell+1}$ can be "backpropagated" through the network by multiplication with the weights, to obtain the δ^ℓ (we move from the output a^L to the most inner part of the network, a^1 , hence the term "backpropagation"). Once we have computed all the δ^ℓ , the derivatives are simply given by using (231) and noting that $\frac{\partial a_k^\ell}{\partial w_{kj}^\ell} = z_j^{\ell-1}$ (Eq. (232)).

The "backpropagation" algorithm can thus read as follows

- First send the training data through the network and compute all the $z_j^\ell(\dots)$ which are the outputs of every neuron.
- Then compute all the intermediate derivatives $[\sigma^\ell(a^{\ell-1})]'$ at the activations $a^{\ell-1}$ given by your training point.
- Once you have all those values, start backpropagating the δ with

$$\delta^L = \frac{\partial L}{\partial a_k^L} = \frac{\partial (y_k - t_k)^2}{\partial y_k} \frac{\partial y_k}{\partial a_k^L} \quad (239)$$

$$= 2(y_k - t_k) [\sigma^L(a_k^L)]' \quad (240)$$

and then

$$\delta_k^\ell = \sum_j \delta_j^{\ell+1} (\sigma^\ell)' w_{jm}^\ell \quad (241)$$

- Finally, once you have all the δ_k^ℓ , compute the derivatives as

$$\frac{\partial L}{\partial w_{kj}^\ell} = \frac{\partial L}{\partial a_k^\ell} z_j^\ell \quad (242)$$

13 Clustering

13.1 Hierarchical clustering

The simplest approach to clustering are the so-called hierarchical clustering algorithms which either start from a single cluster and then divide this cluster iteratively (*Divisive*) or start from the individual prototypes and then merge those prototypes with each other based on their relative similarity. We detail those algorithms below

13.1.1 Agglomerative Clustering

In *agglomerative clustering*, each of the prototypes are merged following a given criterion, until they form a single cluster. Depending on the desired number of clusters, the procedure

In each of those approaches, the combination of subclusters is motivated by the minimization of a particular notion of dissimilarity.

- In *Single Linkage* (SL), we merge at each step, we consider as dissimilarity the distance between the closest pair of points $x_i \in S_1$ and $x_j \in S_2$, and merge the two clusters that minimize this dissimilarity,

$$d_{SL}(S_1, S_2) = \min_{\substack{x_i \in S_1 \\ x_j \in S_2}} d(x_i, x_j) \quad (243)$$

- In *Complete Linkage* (CL), the dissimilarity of any pair of clusters S_1 and S_2 is defined as the distance between the furthest pair of points

$$d_{CL}(S_1, S_2) = \max_{\substack{x_i \in S_1 \\ x_j \in S_2}} d(x_i, x_j) \quad (244)$$

- Finally, an tradeoff between the two approaches can be obtained by taking the dissimilarity between any two sub-clusters to be the average of the inter-cluster distances. This approach is known as *Group Average* (GA)

$$d_{GA}(S_1, S_2) = \frac{1}{N_{S_1} N_{S_2}} \sum_{i \in S_1} \sum_{j \in S_2} d(x_i, x_j) \quad (245)$$

where we let N_{S_1} and N_{S_2} denote the number of prototypes in sub-cluster S_1 and S_2

The iterative merging of sub-clusters defines a dendrogram which can be cut (i.e the procedure can be stopped) at any level depending on the desired number of clusters. Each of the aforementioned approach also has an interpretation in terms of the diameter of each cluster.

13.1.2 Divisive Clustering

As the name indicates, *Divisive clustering* is based on recursively splitting a parent cluster C into two child subclusters C_1 and C_2 following a particular criterion. At each step, one cluster is chosen to be subdivided. This cluster is often the one with the largest diameter $\delta = \max_{i,j \in S_\infty} d(\mathbf{x}_i, \mathbf{x}_j)$ or with the largest average dissimilarity,

$$\bar{d} = \frac{1}{N_C} \sum_{i,j \in C} d(\mathbf{x}_i, \mathbf{x}_j) \quad (246)$$

The division of the parent cluster can be done for example by applying K-means to that cluster with two centroids and returning the corresponding two subclusters, An alternative, due to Macnaughton Smith et al. relies on removing from the clusters all the points that have a sufficiently high dissimilarity with respect to all the other points in the original cluster. If we let C_0 to denote the parent cluster, the methods then defines the first subcluster C_1 by removing from C_0 the point with the largest average dissimilarity

$$x^* = \operatorname{argmax}_x \frac{1}{N_{C_0} - 1} \sum_{x_j \in C_0} d(x_i, x) \quad (247)$$

Here N_{C_0} is used to denote the number of prototypes in the parent cluster C_0 . The first subcluster C_1 is thus defined as the singleton $C_1 = \{x\}$ and the algorithm proceeds by removing the next point with the highest average dissimilarity from $C_0 \setminus \{x^*\}$ and moving it to C_1 . This continues until there is no more prototype in $C_2 = C_0 \setminus C_1$ that is more similar to the points in C_1 than it is to the points in C_2 . I.e, the subdivision stops whenever there is no more y for which the difference

$$D_{C_2} - D_{C_1} = \frac{1}{|C_2| - 1} \sum_{x_j \in C_2} d(y, x_j) - \frac{1}{|C_1|} \sum_{x_j \in C_1} d(x_j, y) \quad (248)$$

is non negative.

13.2 K-means

K-means is one of the most popular and most straightforward clustering algorithm. In K-means the general idea is to find within the data, K clusters that minimize the pairwise distance between the prototypes belonging to the cluster. If we let C_1, \dots, C_K to denote the K clusters in which we want to group our prototypes $\mathbf{x}_1, \dots, \mathbf{x}_N$, and use \mathcal{C} to denote the mapping which associates

a cluster $k = 1, \dots, K$ to each prototype \mathbf{x}_ℓ , that is if the cluster of \mathbf{x}_ℓ is \mathcal{C}_k , we write $\mathcal{C}(\mathbf{x}_\ell) = k$ or $\mathcal{C}(\mathbf{x}_\ell) = \mathcal{C}_k$. K-means minimizes the following objective

$$\min_{\mathcal{C}} \sum_{k=1}^K \sum_{(i,j) \in \mathcal{C}_k} \|\mathbf{x}_i - \mathbf{x}_j\|^2 \quad (249)$$

That is, we only take into account the intra-cluster distance (the distance between points belonging to the same cluster). The main advantage of K-means lies in the use of the Euclidean distance. By using the properties of this distance, one can in fact show that the cost (249) can be reduced to finding the clustering which minimizes the sum of the distances points of a same cluster and their center of mass. To see this, first write the Euclidean distance in terms of inner products,

$$\frac{1}{2} \sum_{k=1}^K \sum_{i \in \mathcal{C}_k} \sum_{j \in \mathcal{C}_k} \|\mathbf{x}_i - \mathbf{x}_j\|^2 = \frac{1}{2} \sum_{k=1}^K \sum_{i \in \mathcal{C}_k} \sum_{j \in \mathcal{C}_k} \langle \mathbf{x}_i, \mathbf{x}_i \rangle + \langle \mathbf{x}_j, \mathbf{x}_j \rangle - 2\langle \mathbf{x}_i, \mathbf{x}_j \rangle \quad (250)$$

$$= \frac{1}{2} \sum_{k=1}^K \sum_{i \in \mathcal{C}_k} 2\langle \mathbf{x}_i, \mathbf{x}_i \rangle N_k - \frac{1}{2} \sum_k \sum_{i \in \mathcal{C}_k} 2\langle \mathbf{x}_i, \sum_{j \in \mathcal{C}_k} \mathbf{x}_j \rangle \quad (251)$$

$$= \left(\sum_{k=1}^K N_k \sum_{i \in \mathcal{C}_k} \langle \mathbf{x}_i, \mathbf{x}_i \rangle - \sum_{k=1}^K N_k \sum_{i \in \mathcal{C}_k} \langle \mathbf{x}_i, \sum_{j \in \mathcal{C}_k} \mathbf{x}_j \rangle \frac{1}{N_k} \right) \quad (252)$$

$$= \sum_{k=1}^K N_k \sum_{i \in \mathcal{C}_k} \left(\langle \mathbf{x}_i, \mathbf{x}_i \rangle - \langle \mathbf{x}_i, \sum_{j \in \mathcal{C}_k} \mathbf{x}_j \rangle \frac{1}{N_k} \right) \quad (253)$$

$$(254)$$

Then add and subtract the quantity $Q \equiv \sum_{k=1}^K \sum_{i \in \mathcal{C}_k} \frac{1}{N_k} \sum_{j \in \mathcal{C}_k} \langle \mathbf{x}_i, \mathbf{x}_j \rangle$,

$$\sum_{k=1}^K N_k \sum_{i \in \mathcal{C}_k} \left(\langle \mathbf{x}_i, \mathbf{x}_i \rangle - \langle \mathbf{x}_i, \sum_{j \in \mathcal{C}_k} \mathbf{x}_j \rangle \frac{1}{N_k} \right) + Q - Q \quad (255)$$

$$= \sum_{k=1}^K N_k \sum_{i \in \mathcal{C}_k} \left(\langle \mathbf{x}_i, \mathbf{x}_i \rangle - \langle \mathbf{x}_i, \sum_{j \in \mathcal{C}_k} \mathbf{x}_j \rangle \frac{1}{N_k} + \frac{1}{N_k^2} \langle \sum_{j \in \mathcal{C}_k} \mathbf{x}_j, \sum_{j \in \mathcal{C}_k} \mathbf{x}_j \rangle \right) \quad (256)$$

$$- \sum_{k=1}^K \sum_{i \in \mathcal{C}_k} \left(\frac{1}{N_k} \sum_{j \in \mathcal{C}_k} \langle \mathbf{x}_j, \mathbf{x}_i \rangle \right) \quad (257)$$

$$= \sum_{k=1}^K N_k \sum_{i \in \mathcal{C}_k} \left(\langle \mathbf{x}_i, \mathbf{x}_i \rangle - 2\langle \mathbf{x}_i, \frac{1}{N_k} \sum_{j \in \mathcal{C}_k} \mathbf{x}_j \rangle + \langle \frac{1}{N_k} \sum_{j \in \mathcal{C}_k} \mathbf{x}_j, \frac{1}{N_k} \sum_{j \in \mathcal{C}_k} \mathbf{x}_j \rangle \right) \quad (258)$$

$$= \sum_{k=1}^K N_k \sum_{i \in \mathcal{C}_k} \left\| \mathbf{x}_i - \frac{1}{N_k} \sum_{j \in \mathcal{C}_k} \mathbf{x}_j \right\|^2 \quad (259)$$

Formulation (259) is interesting because it gives us a heuristic to look for the optimal assignment. From (259) it indeed seems reasonable to start with a set of centroids, define the assignment by relating each prototype to its closest centroid and then update the centroid by defining them as the center of mass of each of the newly generated clusters. This idea is summarized below.

- *Assignment.* Assign each prototype to the cluster represented by the nearest centroid. I.e, for a given set of centroids, y_k , $k = 1, \dots, K$, if we let $\mathcal{C}(\mathbf{x}_\ell)$ to denote the cluster of the ℓ^{th} prototype \mathbf{x}_ℓ , then $\mathcal{C}(\mathbf{x}_k)$ is updated as follows

$$\mathcal{C}(\mathbf{x}_\ell) \leftarrow \min_k \|\mathbf{c}_k - \mathbf{x}_\ell\| \quad (260)$$

- *Update of the centroids.* Once the cluster of each point has been determined, the main advantage of K-means is that the use of the Euclidean distance enables a straightforward update of each centroid as the center of mass of the newly formed clusters:

$$\mathbf{c}_\ell = \frac{1}{N_\ell} \sum_{k \in \mathcal{C}_\ell} \mathbf{x}_k \quad (261)$$

It is worth noting that K-means might sometimes return empty clusters (some of the centroids might be moved towards a region of space without any prototypes). In this case, one can simply restart the algorithm by replacing this centroid with for example the prototype \mathbf{x}^* from the dataset that is located the furthest away from all the other (well defined) centroids. Proceeding like this will eliminate the point \mathbf{x}^* that contributes the most to the error (259).

An alternative approach could be to restart the algorithm by keeping the non zero clusters and resampling a centroid from the cluster that has the largest intra-cluster SSE. This will also reduce the error (259).

The whole process can obviously be repeated when multiple empty clusters are returned.

Several initialization approaches have been proposed for K-means. The popular ones are listed below (see [18] for more details).

- *Random partitioning.* The approach divides the dataset in K distinct clusters chosen at random.
- The *Forgy* or *Lloyd-Forgy* approach [13, 5]. The method picks K feature vectors at random to define the centroids and assigns all the remaining feature vectors to the nearest centroid. There is not recalculation of the centroids after each assignment.
- *K-means++.* The first centroid is selected at random. The next centroid is selected as the point that is the farthest from the currently selected. The selection continues until K centroids are defined.
- *MacQueen approach* [15]. Choose K instances of the database at random. Assign, following some order, the rest of the dataset to each of the centroids. After each assignment, recalculate the centroids.
- *Kaufman approach* [9]. The Kaufman initialization is a deterministic initialization which places the centroids in the areas where there is a higher density of prototypes \mathbf{x}_i . The Kaufman initialization is shown in Fig. 1.

```

1 Set  $\bar{\mathbf{x}}_0$  to be the median feature vector (i.e the one that is the most centrally
  located);
2 Set  $\mathcal{S} = \{\bar{\mathbf{x}}_0\}$  and let  $\mathcal{S}^c$  denote the set of remaining feature vectors.;
3 for all the remaining feature vectors  $\mathbf{x}_i \in \mathcal{S}^c$  do
4   | Compute  $\bar{d}_i \equiv \min_{k \in \mathcal{S}} d(\bar{\mathbf{x}}_k, \mathbf{x}_i)$ ;
5   | Set  $C_{i,j} \equiv \max(\bar{d}_i - d(\mathbf{x}_j, \mathbf{x}_i), 0)$ 
6 end
7 Choose the next centroid  $\bar{\mathbf{x}}_{k+1}$  as the one that maximizes  $\max_i C_{ij}$ ;
8 if there are  $K$  points in  $\mathcal{S}$  then
9   | Stop
10 end
11 else
12   | Set  $\mathcal{S} \leftarrow \mathcal{S} \cup \bar{\mathbf{x}}_{k+1}$  and go back to step 1.
13 end
```

Algorithm 1: Kaufman initialization

Convergence of K -means is studied in [14, 3]. Among the initialization procedures,

13.3 K-medoid

K-medoid is the restriction of K-means to centroids that are part of the original set of prototypes, \mathcal{D} . K-medoid is more flexible in the sense that it can be applied with any given distance, but the price to pay is a higher computational cost. Indeed, more general distances are not guaranteed to lead to losses, such as (259), for which the optimal centroid for a given assignment can be obtained directly as the center of mass of each cluster. In the K-medoid framework, the update of the centroids requires screening all the points within a given cluster and selecting the one that minimizes its distance to all the others, thus increasing the cost from $O(N)$ to $O(N^2)$. The algorithm is summarized by the two steps below.

- *Initialization.* Choose K points at random as the initial centroids $\{m_k\}_{k=1}^K$.
- *Assignment.* Divide the dataset according to those centroids by assigning each point as

$$\mathcal{C}(\mathbf{x}_i) = \underset{k}{\operatorname{argmin}} D(\mathbf{x}_i, \mathbf{m}_k) \quad (262)$$

- *Update of the centroids.* Update the centroids by choosing them among the feature vectors as

$$m_k \leftarrow \underset{\mathbf{x} \mid \mathcal{C}(\mathbf{x})=k}{\operatorname{argmin}} \sum_{\mathbf{x} \mid \mathcal{C}(\mathbf{x})=k} D(\mathbf{x}, \mathbf{x}_k) \quad (263)$$

K-medoid then iterates over the two steps above until the assignement stop changing.

13.4 The EM Algorithm

Both K-means and K-medoid relies on hard assignment. I.e., a point is classified as belonging or not to a given cluster. Such a hard approach leads to a lack of robustness in the sense that it does not come with a clear quantification of the uncertainty associated to the classification of each prototype. The expectation maximization (EM) approach works in a similarly way but defines the clusters based on probability distributions. The advantage of turning to probability distributions, is that each assignment now comes with a level of confidence encoded by the probabilities. In the EM algorithm, the cluster of a point is defined as the probability distribution that gives the highest probability of occurence for that point.

14 Linear latent variable models

In all the models discussed below, we are interested in finding representations of the form $x(\theta)$ for a dataset $\{\mathbf{x}_k\}_{k=1}^N$. That is we would like to understand whether most of the data can be generated from a small number of meaningful parameters.

As usual in machine learning, most of the models that we consider below have a statistical interpretation that falls within either of the Bayesian or Frequentist framework.

14.1 Factor Analysis (FA)

The simplest latent variable model assumes that the prototypes can be represented from the linear combination of a few (common) sources of variation,

$$\mathbf{x}_i = \mathbf{W}\boldsymbol{\theta}_i, \quad i = 1, \dots, N, \quad \mathbf{x}_i \in \mathbb{R}^D \quad (264)$$

\mathbf{W} is a matrix of size $D \times K$ with $K \ll N$. Assuming a sufficiently low dimensional latent representation however is not enough to make the decomposition uniquely identifiable. In particular,

for a given set of prototypes $\{\mathbf{x}_i\}_{i=1}^N$, and any rotation matrix \mathbf{R}^2 one always has

$$\mathbf{x}_i = \mathbf{W}\mathbf{R}\mathbf{R}^T\boldsymbol{\theta} = \mathbf{W}\boldsymbol{\theta} \quad (265)$$

which already gives two possible decompositions for \mathbf{x}_i . The Factor Analysis model also leads to scaling ambiguities as for any $\alpha \in \mathbb{R}$, one can always write

$$\mathbf{x}_i = (\alpha\mathbf{W})(\alpha^{-1}\boldsymbol{\theta}_i) \quad (266)$$

To ensure unique recovery of the latent factors together with the factor loading matrix \mathbf{W} , we need to enforce additional constraints that will prevent the ambiguities discussed above to occur. The most commonly used solutions are listed below and detailed in the subsequent subsections.

- Requiring \mathbf{W} to be orthonormal, i.e., $\mathbf{W}^T\mathbf{W} = \mathbf{I}$. This is the approach followed by Principal Component Analysis (PCA) which is discussed in section 14.2.
- Require \mathbf{W} to be lower triangular. I.e in this case, the first feature in \mathbf{x}_i only depends on the first factor θ_1 , the second feature is a linear combination of the first two latent factors, θ_1, θ_2 and so on. This approach usually also requires the weights from the factor loading matrix to be positive.
- Use of non-Gaussian priors for the latent factors. This, as we will see, leads to Independent Component Analysis (ICA).
- Consider sparsity promoting priors on the weights. An alternative to the lower triangular factor loading matrix is to avoid pre-specifying the location of the non zero entries in \mathbf{W} but to encourage the entries to be zero by adding an ℓ_1 regularization term.

The idea underlying Factor Analysis can be used inside a statistical framework. If we consider measurements that are generated from the latent model following a Gaussian distribution

$$\mathbf{x}_i = \boldsymbol{\mu} + \mathbf{W}\mathbf{z}_i + \varepsilon_i \quad (267)$$

Where the $\varepsilon_i \in \mathbb{R}^d$ follows a multivariate Normal distribution with covariance $\boldsymbol{\Psi}$,

$$p(\mathbf{x}_i) = \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Psi}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{z}_i)^T\boldsymbol{\Psi}^{-1}(\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{z}_i)\right\} \quad (268)$$

and follow a Frequentist approach, therefore focusing exclusively on the (log) likelihood, we have

$$-\log(p(\mathbf{x}_i|\boldsymbol{\mu}, \mathbf{W}, \mathbf{z}_i)) = -\log(\mathcal{N}(\boldsymbol{\mu} + \mathbf{W}\mathbf{z}_i, \boldsymbol{\Psi})) \quad (269)$$

$$= -\log\left(-\frac{1}{(2\pi)^{d/2}|\boldsymbol{\Psi}|^{1/2}} \exp(-(\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{z}_i)^T\boldsymbol{\Psi}^{-1}(\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{z}_i))\right) \quad (270)$$

$$= \log((2\pi)^{d/2}|\boldsymbol{\Psi}|^{1/2}) + (\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{z}_i)^T\boldsymbol{\Psi}^{-1}(\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{z}_i) \quad (271)$$

In the particular case where $\boldsymbol{\Psi} = \sigma^2\mathbf{I}$ (which corresponds to independent features), we get

$$-\log(p(\mathbf{x}_i|\boldsymbol{\mu}, \mathbf{W}, \mathbf{z}_i)) = d\log(2\pi) + \frac{d}{2}\log(\sigma) + \frac{1}{\sigma^2}\|\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{z}_i\|^2 \quad (272)$$

When $\sigma \rightarrow 0$, minimizing (272) (i.e maximizing the log likelihood), corresponds to the Factor Analysis minimization problem discussed above,

$$\min_{\mathbf{W}, \mathbf{z}_i} \|\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{z}_i\|^2 \quad (273)$$

²A rotation matrix is an orthogonal matrix with unitary determinant, i.e. $\mathbf{R}^T\mathbf{R} = \mathbf{R}\mathbf{R}^T = \mathbf{I}$

Minimizing (273) can thus be understood as neglecting the σ , or assuming no uncertainty on the prototypes. Without additional assumption on either \mathbf{W} or \mathbf{z}_i , solving this problem is hopeless. As we already mentioned, one possibility to make the problem better posed is to add assumptions on \mathbf{W} . On top of those assumptions, one can take a Bayesian perspective on the problem and consider a prior on \mathbf{z}_i . Multiplying our likelihood with such a prior gives us the possibility to average out the dependence on the \mathbf{z}_i (in the sense that $p(\mathbf{x}_i) = \int p(\mathbf{x}_i|\mathbf{z}_i)p(\mathbf{z}_i) d\mathbf{z}_i$). If we take a Gaussian prior for the \mathbf{z}_i for example, $p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i|\boldsymbol{\mu}_0, \boldsymbol{\Psi}_0)$, we can average out the dependence as

$$p(\mathbf{x}_i|\boldsymbol{\theta}) = \int \mathcal{N}(\mathbf{x}_i|\mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi})\mathcal{N}(\mathbf{z}_i|\boldsymbol{\mu}_0, \boldsymbol{\Psi}_0) d\mathbf{z}_i \quad (274)$$

Using (54), we can then write this probability distribution as

$$p(\mathbf{x}_i|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i|\mathbf{W}\boldsymbol{\mu}_0 + \boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\boldsymbol{\Psi}_0\mathbf{W}^T) \quad (275)$$

Note that is not an MAP estimation as we did not use priors on our parameters $\mathbf{W}, \boldsymbol{\mu}, \boldsymbol{\Psi}$ and $\boldsymbol{\mu}_0, \boldsymbol{\Psi}_0$. Moreover, as already explained, this procedure removes the dependence in the latent factors \mathbf{z}_0 but does not solve the ill posedness of the problem (i.e we still have too many parameters). The next sections address this last difficulty.

14.2 Principal Component Analysis (PCA)

Principal Component Analysis is the restriction of FA to orthonormal factor loading matrices \mathbf{W} . In principal component analysis, one is interested in finding a subset of vectors $\{\mathbf{v}_\ell\}_{\ell=1}^L$ that capture most of the information present in the dataset $\{\mathbf{x}_k\}_{k=1}^N$. Mathematically, we therefore want to find a model of the form $\{\boldsymbol{\mu}, \mathbf{W}, \boldsymbol{\lambda}\}$ where \mathbf{W} has orthonormal columns, i.e. $\mathbf{W}^*\mathbf{W} = \mathbf{I}$, and where any point from $\{\mathbf{x}_k\}_{k=1}^N$ can read as

$$\mathbf{x}_k \approx \boldsymbol{\mu} + \mathbf{W}\boldsymbol{\lambda} \quad (276)$$

To find such a model, a natural approach is to minimize the sum of the deviations between the prototypes and their approximations $\boldsymbol{\mu} + \mathbf{W}\boldsymbol{\lambda}$. This can be done by minimizing the ℓ_2 , as follows

$$\min_{\boldsymbol{\mu}, \boldsymbol{\lambda}, \mathbf{W}} \sum_{i=1}^N \|\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{z}_i\|^2 \quad (277)$$

Again, formulation (277) can be viewed as following from the maximization of the log-likelihood provided that the prototypes \mathbf{x}_i follow a Normal distribution with mean $\boldsymbol{\mu} + \mathbf{W}\mathbf{z}_i$.

Solving problem (277) for \mathbf{z}_i and $\boldsymbol{\mu}$ gives

$$\boldsymbol{\mu} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \quad (278)$$

$$\mathbf{z}_i = \mathbf{W}^*(\mathbf{x}_i - \boldsymbol{\mu}) \quad (279)$$

Given expressions (278) and (279), the minimization problem (277) can then be written as

$$\min_{\mathbf{W}} \sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}} - \mathbf{W}\mathbf{W}^*(\mathbf{x}_i - \bar{\mathbf{x}})\|^2 \quad (280)$$

The matrix $\mathbf{W}\mathbf{W}^*$ that appears in this expression is a projection matrix (i.e a matrix for which $\mathbf{P}^2 = \mathbf{P}$). $\mathbf{W}\mathbf{W}^*$ in particular maps every point from the dataset $\{\mathbf{x}_i\}_{i=1}^N$ onto the subspace spanned by the columns of \mathbf{W} . Indeed, note that we have

$$\mathbf{W}\mathbf{W}^*\mathbf{x}_i = \sum_{\ell=1}^L \mathbf{w}_\ell \langle \mathbf{w}_\ell, \mathbf{x}_i \rangle \quad (281)$$

Now if we develop the norm in (280) and use $\sum_{k=1}^N w_k^2 = \langle \mathbf{w}, \mathbf{w} \rangle$, we have

$$\min_{\mathbf{W}} \sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}} - \mathbf{W}\mathbf{W}^*(\mathbf{x}_i - \bar{\mathbf{x}})\|^2 \quad (282)$$

$$= \min_{\mathbf{W}} \sum_{i=1}^N \langle \mathbf{x}_i - \mathbf{W}\mathbf{W}^*\mathbf{x}_i, \mathbf{x}_i - \mathbf{W}\mathbf{W}^*\mathbf{x}_i \rangle \quad (283)$$

$$= \min_{\mathbf{W}} \sum_{i=1}^N \langle \mathbf{x}_i, \mathbf{x}_i \rangle + \sum_{i=1}^N \langle \mathbf{W}\mathbf{W}^*\mathbf{x}_i, \mathbf{W}\mathbf{W}^*\mathbf{x}_i \rangle - 2 \sum_{i=1}^N \langle \mathbf{W}\mathbf{W}^*\mathbf{x}_i, \mathbf{x}_i \rangle \quad (284)$$

$$= \min_{\mathbf{W}} \sum_{i=1}^N \|\mathbf{x}_i\|^2 + \sum_{i=1}^N \mathbf{x}_i^* \mathbf{W}\mathbf{W}^* \mathbf{x}_i - 2 \sum_{i=1}^N \mathbf{x}_i^* \mathbf{W}\mathbf{W}^* \mathbf{x}_i \quad (285)$$

Between (284) and (285), we used $\mathbf{W}^*\mathbf{W} = \mathbf{I}$. The first term in the last line can be treated as a constant as it does not depend on \mathbf{W} (i.e we can disregard it). We are thus left with

$$\min_{\mathbf{W}} - \sum_{i=1}^N \mathbf{x}_i^* \mathbf{W}\mathbf{W}^* \mathbf{x}_i \quad (286)$$

which can be written as a maximization problem

$$\max_{\mathbf{W}} \sum_{i=1}^N \mathbf{x}_i^* \mathbf{W}\mathbf{W}^* \mathbf{x}_i \quad (287)$$

If we were looking for the first component only, the problem would reduce to

$$\max_{\mathbf{w}_1} \sum_{i=1}^N \mathbf{x}_i^* \mathbf{w}_1 \mathbf{w}_1^* \mathbf{x}_i \quad (288)$$

$$= \max_{\mathbf{w}_1} \mathbf{w}_1^* \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^* \mathbf{w}_1 \quad (289)$$

$$= \max_{\mathbf{w}_1} \mathbf{w}_1^* \mathbf{X}^* \mathbf{X} \mathbf{w}_1 \quad (290)$$

In the expression above, we again used \mathbf{X} to denote the matrix encoding the prototypes $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^D$,

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} \quad (291)$$

For any given matrix $\mathbf{X}^* \mathbf{X}$, problem (290) will return the eigenvector of $\mathbf{X}^* \mathbf{X}$ corresponding to the largest eigenvalue (why?). To find the first column of \mathbf{W} , we can thus take the eigenvalue decomposition of $\mathbf{X}^* \mathbf{X}$ and retain the first eigenvector. In fact this approach continues to work for any number of components. That is to say, if we want to find the best subspace of dimension q , it suffices to take the eigenvalue decomposition of $\mathbf{X}^* \mathbf{X}$ and to retain the first L largest eigenvectors. This idea is summarized below

Theorem 11 (see for example [2]). Suppose that we are given a dataset $\{\mathbf{x}_i\}_{i=1}^N$ and that we want to find an orthonormal set of L basis vectors \mathbf{w}_j that best approximate the data, or equivalently, that minimize the average reconstruction error

$$e(\mathbf{W}, \mathbf{z}) \equiv \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{W} \mathbf{z}_i\|^2 \quad (292)$$

Then the optimal solution for \mathbf{W} is obtained by building the *empirical covariance matrix*, $\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T$, and setting $\mathbf{W} = \mathbf{V}_L$ where \mathbf{V}_L contains the L eigenvectors with largest eigenvalues of $\hat{\Sigma}$. The columns of \mathbf{V}_L are often called *principal directions*.

Equivalently one could work with the singular value decomposition (SVD) of \mathbf{X} , and retain the first L right singular vectors obtained from the matrix $\tilde{\mathbf{V}}$ in this SVD $\mathbf{X} = \tilde{\mathbf{U}} \Sigma \tilde{\mathbf{V}}^*$. We then have $\mathbf{X}^* \approx \tilde{\mathbf{V}}_L (\Sigma \tilde{\mathbf{U}}^*)_L$ where $(\Sigma \tilde{\mathbf{U}}^*)_L$ is used to denote the first L rows of $\Sigma \tilde{\mathbf{U}}^*$ and encodes the λ_k . In other words, to solve the PCA problem for a set of prototypes $\{\mathbf{x}_k\}_{k=1}^N$,

- (i) Center the dataset, $\mathbf{x}_k \leftarrow \mathbf{x}_k - \frac{1}{N} \sum_{k=1}^N \mathbf{x}_k$
- (ii) Build the matrix $\mathbf{X}^* = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ from the centered prototypes
- (iii) Compute the SVD of \mathbf{X} , $\mathbf{X} = \tilde{\mathbf{U}} \Sigma \tilde{\mathbf{V}}^*$
- (iv) Define the matrix \mathbf{V} from the first L columns of $\tilde{\mathbf{V}}$
- (v) Define the weight vectors λ_k from the k^{th} column of the matrix obtained after retaining the first L rows of $\Sigma \tilde{\mathbf{U}}^*$.

14.2.1 Probabilistic PCA

Just as for Factor Analysis, one can interpret PCA through the lens of Frequentist statistics, as a particular MLE for which the variance of the noise would be arbitrarily small and hence, the $\log(\sigma)$ negligible. When the noise ε_i cannot be neglected, or more generally if we want to follow a more rigorous statistical approach, one can again consider a prior on the latent factors, \mathbf{z}_i . We then get as in FA

$$p(\mathbf{x}_i | \boldsymbol{\theta}) = \mathcal{N}(\mathbf{W} \boldsymbol{\mu}_0 + \boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W} \boldsymbol{\Psi}_0 \mathbf{W}^T) \quad (293)$$

If we consider N observations and take $\boldsymbol{\Psi}_0 = \mathbf{I}$ and $\boldsymbol{\Psi} = \sigma^2 \mathbf{I}$ and $\boldsymbol{\mu}_0 = 0$, writing down the log-likelihood gives (you can verify this as an exercise)

$$\log(p(\mathbf{x}_i | \boldsymbol{\theta})) = -\frac{d}{2} \log(|\mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}|) - \frac{1}{2} \sum_{i=1}^N \mathbf{x}_i \left(\mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I} \right)^{-1} \mathbf{x}_i \quad (294)$$

You can check that the maxima of the likelihood are given by

$$\mathbf{W}^* = \mathbf{V} (\boldsymbol{\Lambda} - \sigma^2 \mathbf{I})^{1/2} \mathbf{R} \quad (295)$$

Where \mathbf{R} is any orthogonal matrix (see the discussion on the uniqueness of \mathbf{W} in section 14.1) and \mathbf{V} and $\boldsymbol{\Lambda}$ follow from the eigenvalue decomposition of the empirical covariance matrix $\mathbf{S} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T = \frac{1}{N} \mathbf{X}^T \mathbf{X}$. The approach in probabilistic PCA consists resolving the ambiguity on \mathbf{W} by setting $\mathbf{R} = \mathbf{I}$. We then get $\mathbf{W}^* = \mathbf{V} (\boldsymbol{\Lambda} - \sigma^2 \mathbf{I})^{1/2}$. In particular, this recovers the traditional PCA solution when $\sigma \rightarrow 0$ (note that this last approach can hence be considered as slightly more rigorous as the case $\sigma = 0$ is well defined.)

It is however also worth noting that in this slightly more rigorous framework (statistically speaking), we do not recover an orthogonal matrix when $\sigma^2 > 0$.

14.3 Independent Component Analysis (ICA)

In [Independent Component Analysis \(ICA\)](#), one is interested in finding a description of the data by means of independent factors. ICA has led to several interesting applications to source separation as well as biomedical signal processing, with a particular emphasis on EEG source separation.

Consider two signals \mathbf{s}_1 and \mathbf{s}_2 of length T (we use $s_{1,t}$ to denote the t^{th} sample from \mathbf{s}_1) that are mixed together through the matrix $\mathbf{W} \in \mathbb{R}^{2 \times 2}$. The mixings are encoded in two distinct recordings that are measured at microphones x_1 and microphone x_2 ,

$$\begin{bmatrix} x_{1,t} \\ x_{2,t} \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \begin{bmatrix} s_{1,t} \\ s_{2,t} \end{bmatrix} + \begin{bmatrix} \varepsilon_{1,t} \\ \varepsilon_{2,t} \end{bmatrix} \quad (296)$$

Following this model, at any time, the first and second microphones thus record the signals

$$x_{1,t} = w_{11}s_{1,t} + w_{12}s_{2,t} + \varepsilon_{1,t} \quad (297)$$

$$x_{2,t} = w_{21}s_{1,t} + w_{22}s_{2,t} + \varepsilon_{2,t} \quad (298)$$

Here the weights w_{ij} can be understood as encoding the distances from the microphones to the speakers. If the matrix \mathbf{W} is given, it is of course easy to recover the sources by, for example ℓ_2 norm minimization. Perhaps more surprisingly, is the fact that the recovery of the sources remain possible even when the mixing coefficients are unknown.

The original application of ICA was to the so-called cocktail party problem, where one wants to recover independent speeches from their mixtures. Among the notorious successes of ICA, one must also mention the demixing of brain activities from EEG recordings on the scalp, as well as feature extraction (e.g. the identification of shocks from frequent fluctuations in finance [\[1\]](#)).

To solve problem (296), it is possible to take a Maximum Likelihood viewpoint on the problem. Let $\mathbf{x} = \mathbf{W}\mathbf{s} = f(\mathbf{s})$. Then note that in the scalar case, if we use p_x to denote the probability density function (see definition 9) and P_x to denote the cumulative mass function (see definition 7), we can always write

$$p_x(x) = \frac{d}{dx}P_x(x) = \frac{d}{dx}P_s(f^{-1}(x)) = \frac{ds}{dx} \frac{d}{ds}P_s(s) = \frac{ds}{dx}p_s(s) \quad (299)$$

as $p_x(x)$ is a probability density (i.e it has to be non negative) we set

$$p_x(x) = p_s(s) \left| \frac{ds}{dx} \right| \quad (300)$$

When considering multiple dimensions such as in ICA, the change of variable formula (300) extends as

$$p_x(\mathbf{x}) = p_s(\mathbf{s}) \left| \det \left(\frac{\partial \mathbf{s}}{\partial \mathbf{x}} \right) \right| \quad (301)$$

Where $\frac{\partial \mathbf{s}}{\partial \mathbf{x}}$ is the Jacobian matrix (see (17)). With this in mind, we can write

$$p_x(\mathbf{W}\mathbf{s}) = p_s(\mathbf{s}) \left| \det(\mathbf{W}^{-1}) \right| \quad (302)$$

If we let $\mathbf{V} = \mathbf{W}^{-1}$, in particular we get

$$p_x(\mathbf{x}) = p_x(\mathbf{W}\mathbf{s}) = p_s(\mathbf{V}\mathbf{x}) \left| \det(\mathbf{V}) \right| \quad (303)$$

The maximum likelihood approach looks for the parameters that are maximizing the probability of observing the data (i.e. the prototypes \mathbf{x}_i , $i = 1, \dots, N$). Taking the log likelihood (note that the probability to get \mathbf{x} only depends on \mathbf{V} and \mathbf{x} from (303)) gives

$$\log(p(\mathbf{x}|\mathbf{V})) = N \log \left| \det(\mathbf{V}) \right| + \sum_{s=1}^D \sum_{t=1}^N p_s(\mathbf{v}_s^T \mathbf{x}_t) \quad (304)$$

So far, we have introduced an MLE approach to Independent Component Analysis but it is not clear why the maximum likelihood estimator should work.

In ICA we want to recover independent components s_1 and s_2 . By the Central Limit Theorem (Theorem 9), we know that a sum of (a sufficiently large number of) identically distributed, independent random variables tend to have a Gaussian distribution. In particular, from this Theorem, the random source s_i can be expected to “look” less Gaussian than any combination $a_{11}s_1 + a_{12}s_2$.

When considering the set of weighted sum of the form $\mathbf{v}^T \mathbf{x} = \mathbf{v}^T \mathbf{W} \mathbf{s}$, the one vector \mathbf{v} that should therefore correspond to a maximally non Gaussian random variable is the one that returns only one of the (non gaussian) independent sources (as opposed to a linear combination of those sources). But this optimal vector \mathbf{v} leading to the maximally independent sources (or maximally non Gaussian signal) precisely corresponds to a row (lets say the first one) of \mathbf{A}^{-1} . Indeed, for this row, which we label $\mathbf{v} = (\mathbf{A}^{-1})_{1,:}$, we have $\mathbf{v}^T \mathbf{x} = (\mathbf{A}^{-1})_{1,:} \mathbf{A} \mathbf{s} = s_1$.

When we want to recover independent sources from linear combinations of those sources, maximizing the “non Gaussian” nature of $\mathbf{w}^T \mathbf{z}$ over the vectors \mathbf{v} thus seems to be a relatively good idea. How can we then maximize this “non-Gaussianity” of $\mathbf{w}^T \mathbf{z}$? and how does this connect to the MLE which is used in the FastICA algorithm of for example scikit-learn?

To understand this, consider the following measure of (non-)gaussianity known as [negentropy](#),

$$\text{negentropy}(z) = \mathcal{H}(\mathcal{N}(\mu, \sigma^2)) - \mathcal{H}(z) \quad (305)$$

In the definition above, $\mathcal{H}(X)$ represents the entropy of a random variable X (you can view it as a measure of the uncertainty associated to the random variable),

$$\mathcal{H}(X) = - \sum_{k=1}^K p(X = k) \log_2 p(X = k) \quad (306)$$

We also used $\mu = \mathbb{E}\{z\}$ as well as $\sigma^2 = \text{var}[z]$. But what is probably more important in (305) is that the negentropy can also be understood the deviation of a random variable to the closest Gaussian distribution.

Following from the discussion above, to find the vector w such that the $w^T z$ is as far as possible from a Gaussian distribution, one could thus maximize the negentropy, i.e., for measurements $j = 1, \dots, K$,

$$\max_{\mathbf{v}} \sum_{j=1}^K \mathcal{H}(\mathcal{N}(\mu_j, \sigma_j^2)) - \mathcal{H}(z_j) \quad (307)$$

Often, when we want to recover independent component, a good approach is to center and whiten the data (see [7]). Centering means that we replace the measurements \mathbf{x} by the difference $\mathbf{x} - \mathbb{E}\mathbf{x}$. Centering the measurements also corresponds to centering the sources as we have $\mathbf{x} - \mathbb{E}\mathbf{x} = \mathbf{A}\mathbf{z} - \mathbf{A}\mathbb{E}\mathbf{z} = \mathbf{A}(\mathbf{z} - \mathbb{E}\mathbf{z})$. After estimating the centered principal components, one can simply get the original principal components by adding back the mean $\mathbf{A}^{-1}\mathbb{E}\mathbf{x}$.

Another transformation that is often applied for computational purposes, is the whitening of the data. Whitening the data means that we apply a linear transform to the measurements vector \mathbf{x} to make the components of this vector uncorrelated, i.e, whitening transforms \mathbf{x} into a vector $\tilde{\mathbf{x}}$ satisfying $\mathbb{E}\{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^*\} = \mathbf{I}$.

Whitening can be done by first computing the eigenvalue decomposition of the sample covariance $\mathbb{E}\{\mathbf{x}\mathbf{x}^T\} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ and then transforming \mathbf{x} into $\tilde{\mathbf{x}}$ as

$$\tilde{\mathbf{x}} = \mathbf{U}\mathbf{\Lambda}^{-1/2}\mathbf{U}^T \mathbf{x} \quad (308)$$

Whitening also has a consequence on the mixing matrix \mathbf{A} . Indeed, the whitening turns the measurements into

$$\tilde{\mathbf{x}} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T \mathbf{A} \mathbf{s} \quad (309)$$

which is therefore equivalent to changing replacing the mixing matrix \mathbf{A} with $\tilde{\mathbf{A}}$ defined as $\tilde{\mathbf{A}} = \mathbf{U}\mathbf{\Lambda}^{-1/2}\mathbf{U}^T$. From the whitening transformation (308), we get

$$\mathbb{E}\{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^T\} = \tilde{\mathbf{A}}\mathbb{E}\{\mathbf{s}\mathbf{s}^T\}\tilde{\mathbf{A}}^T = \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T = \mathbf{I} \quad (310)$$

The mixing matrix is now orthogonal,

$$\tilde{\mathbf{A}}^T\tilde{\mathbf{A}} = \mathbf{U}\mathbf{\Lambda}^{-1/2}\mathbf{U}^T\mathbf{A}\mathbf{A}^T\mathbf{E}\mathbf{\Lambda}^{-1/2}\mathbf{E}^T \quad (311)$$

$$= \mathbf{U}\mathbf{\Lambda}^{-1/2}\mathbf{U}^T\mathbb{E}\{\mathbf{x}\mathbf{x}^T\}\mathbf{E}\mathbf{\Lambda}^{-1/2}\mathbf{E}^T \quad (312)$$

$$= \mathbf{I}. \quad (313)$$

So in particular, whitening means that we now only need to estimate $n(n-1)/2$ parameters instead of the initial n^2 .

Once the data has been centered and whitened, the first term in (307) is constant, and maximizing the negentropy thus reduces to the maximization

$$\max_{\mathbf{V}} \sum_j -\mathcal{H}(z_j) = \sum_j \mathbb{E} \log(p(x_j)) \quad (314)$$

which is precisely the MLE. In other words, we can hope the MLE to recover the sources because it corresponds to maximizing non gaussianity of $\mathbf{w}^T\mathbf{x}$ and we know from the Central limit Theorem that the more independent sources s_i we have in a linear combination $\mathbf{w}^T\mathbf{A}\mathbf{s}$, the more Gaussian this combination will be.

So the intuition behind the MLE is that it corresponds to maximizing the deviation between the z_j and the closest Gaussian random variable.

14.3.1 Fast ICA

We now follow the notations in [16] and introduce the FastICA algorithm. Let $G(z_j) = -\log(p(z_j))$ and $g(z_j) = \frac{d}{dz_j}G(z_j)$. Maximizing the log likelihood while asking for an orthogonal matrix \mathbf{A} (note that when \mathbf{A} is orthogonal, the inverse is given by \mathbf{A}^T), lead to the following minimization

$$\min_{\mathbf{v}} \mathbb{E}\{G(\mathbf{v}^T\mathbf{x}) + \lambda(1 - \mathbf{v}^T\mathbf{v})\} \quad (315)$$

We simply replaced the sum in (303) by the expectation and neglect the first term which is constant when \mathbf{V} is orthogonal (determinant is 1).

The gradient and Hessian are given by

$$\nabla f(\mathbf{v}) = \mathbb{E}\{\mathbf{x}g(\mathbf{v}^T\mathbf{x})\} - 2\lambda\mathbf{v} \quad (316)$$

$$\nabla^2 f(\mathbf{v}) = \mathbb{E}\{\mathbf{x}\mathbf{x}^T g'(\mathbf{v}^T\mathbf{x})\} - 2\lambda\mathbf{I} \quad (317)$$

To minimize a function $f(x)$, recall that the Newton steps are defined by

$$\mathbf{x} \leftarrow \mathbf{x} - (\nabla^2 f(\mathbf{x}))^{-1}\nabla f(\mathbf{x}) \quad (318)$$

Now if we make the simplification

$$\mathbb{E}\{\mathbf{x}\mathbf{x}^T g'(\mathbf{v}^T\mathbf{x})\} \approx \mathbb{E}\{\mathbf{x}\mathbf{x}^T\}\mathbb{E}\{g'(\mathbf{v}^T\mathbf{x})\} = \mathbb{E}\{g'(\mathbf{v}^T\mathbf{x})\} \quad (319)$$

in the case of problem (315), the Newton iterations can be written as

$$\mathbf{v} \leftarrow \mathbf{v} - \frac{\mathbb{E}\{\mathbf{x}g(\mathbf{v}^T\mathbf{x}) - 2\lambda\mathbf{v}\}}{\mathbb{E}\{g'(\mathbf{v}^T\mathbf{x})\} - 2\lambda} \quad (320)$$

The expectation can be computed by Monte Carlo methods. Once the iterations have converged, one can define the optimal \mathbf{v}^* as $\mathbf{v}/\|\mathbf{v}\|$ to get the unit vector.

In the developments above we have assumed that the distribution $p(z_j)$ was known. In practice, as long as the distribution of the sources is not Gaussian it seems [16] that it is not critical for the shape of this distribution to be exactly known. Possible choices are the Laplace distribution or other smoother super-Gaussian distribution although it seems that in practice, it is common to just use $G(z) = \sqrt{z}$ or $G(z) = \log \cosh z$

14.4 From one component to multiple components

One approach to extend the algorithm given above to multiple components v_j is to run instances of that algorithm for each v_j in parallel. The issue with such an approach is that we might end up converging to the same vector, i.e. $\mathbf{v}_j = \mathbf{v}$ for all j . In this case we obviously do not recover the inverse \mathbf{A}^{-1} .

To prevent this, one possibility is to compute the component in series, and ensure that they are orthogonal to each other (i.e. decorrelated) through a Gram-Schmidt process. The idea of such a procedure is as follows. Assuming that we have computed components w_1 to w_{p-1} , the p^{th} components is then computed through the one components algorithm above combined with the following two decorrelation steps (which are applied after every Newton step (320))

$$\mathbf{w}_p \leftarrow \mathbf{w}_p - \sum_{j=1}^{p-1} \mathbf{w}_p^T \mathbf{w}_j \mathbf{w}_j \quad (321)$$

$$\mathbf{w}_p \leftarrow \frac{\mathbf{w}_p}{\|\mathbf{w}_p\|} \quad (322)$$

Finally, note that ICA does not always work for non Gaussian sources. Moreover, a main difficulty in ICA lies in the fact that we usually do not have information on the statistical nature of the sources. In particular, natural signals, such as images, often exhibit some correlation. An extreme example of such correlation is shown in Fig. 4. When signals are correlated, ICA will not be able to recover the original sources.

15 Manifold Learning

In the previous sections, we have considered some of the popular linear dimensionality reduction models (a.k.a linear latent variable models). In both ICA and PCA, the mapping between the latent space (i.e. the space of the latent variables) and the space in which the prototypes are living, can be encoded by a linear map. In this chapter, we will extend the discussion to non linear models. In particular, we will discuss methods that focus on learning the structure underlying the data even when this structure is different from the usual hyperplane.

In some datasets, although the relation between the prototypes might be lost when mapping those prototypes onto a linear subspace (for example through PCA or ICA), those prototypes might still be distributed according to a relatively simple (non linear) structure embedded in \mathbb{R}^n . We call such structure a [manifold](#).

Definition 20 ([Manifold \[20\]](#)). *A manifold \mathcal{M} is a metric space (that is a space together with a metric d on the space) such that for $\mathbf{x} \in \mathcal{M}$, there exists some neighborhood U of \mathbf{x} and some integer $n \geq 0$ such that U is homeomorphic (there exists a continuous bijective mapping between the two objects) to \mathbb{R}^n .*

A few popular manifold benchmarks are highlighted in Fig. 6.

To better understand the interest of non linear dimensionality reduction compared to the linear latent variable models, consider Fig. 5. In this figure, we compare some of the non linear

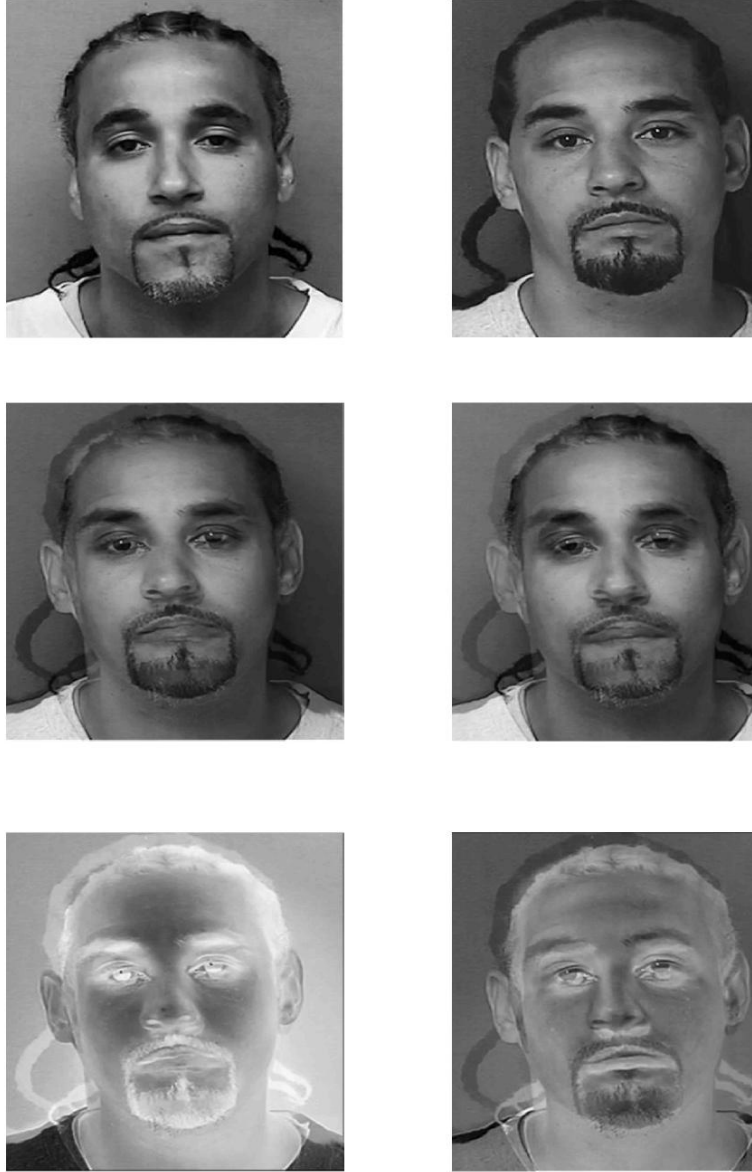


Figure 4: Example of correlated images ($\rho = 0.4$) for which ICA is unable to recover the sources. From top to bottom: Original (source) images, superposition obtained by mixing the images through the mixing coefficients $A_{11} = 1, A_{12} = 2, A_{21} = 2, A_{22} = 1.5$ and reconstruction through FastICA. Richard A. Jones (right) endured 17 years of wrongful incarceration before being finally freed from prison when Ricky L. Amos (left) was finally recognized to be the unique responsible of a robbery committed outside a Walmart in Roeland Park, Kansas. Source: CNN

dimensionality reduction methods that we will discuss below against Principal Component Analysis. The first row in the figure corresponds to applying non linear manifold learning methods, as well as PCA to the original S-shaped manifold from scikit-learn. The second row, corresponds to applying the same methods to a modified version of this manifold, i.e. if x_1, x_2, x_3 are the coordinate of a point on the original S manifold, we obtain the modified manifold as $x_1 = x_1 - 1$, $x_2 = 1.5x_2$ and $x_3 = 0.5x_3$. The modified manifold is thus “shorter and fatter”. When applying a principal component analysis to this second manifold, the PCA plane will be horizontal and the projection onto the principal plane will superimpose several of the manifold parts. The result is a complete destruction of the manifold structure.

Non linear dimensionality reduction methods are thus especially useful when considering datasets that are not distributed on a plane, but exhibit some form of (lower dimensional) structure. In such situations, the non linear methods might be able to capture such structure, while typical linear models such as PCA will not.

Among the methods presented in this chapter, *Multidimensional scaling*, as well as *Sammon’s mapping* (section 15.1) as well as *Isomap* (section 15.2) can be considered as distance preservation methods. By distance preservation, we mean that the methods reduce dimensionality of the data by using the preservation of pairwise distances as a the main criterion. If a lower dimensional representation can be built in a way that preserves pairwise distances. then the geometric structure (global shape or local neighborhood) of the underlying manifold is preserved. In the non linear case, it is usually very difficult to preserve distances (there are many different ways to define non linear manifold).

Locally linear embedding (section 15.3) and (Kohonen) *Self Organizing maps* (section 15.4) on the other hand, are considered as preserving topology of the data. By topology of a manifold, we mean the *neighborhood relationships between subregions of the manifold*. The idea of topology preserving approaches is that pairwise distances are often too constraining or gives too much information. I.e. distances depend not only on the manifold but also on the chosen embedding of the manifold. Additional arguments in favor of the use of topology in dimensionality reduction can be found in [12]. Comparative information between distances are often enough to characterize a manifold. Neighborhood relations are intrinsic to the manifold whereas distances are defined on the whole ambient space and are not particular to the manifold itself. Topology is however more difficult to understand than pairwise distances. *How can we characterize topology of the manifold?* One way is to associate the points to the vertices of a graph and encode proximity of the points by means of edge weights. *Kohonen self organizing maps* rely on a predefined lattice (build as a preliminary step to the method) whether *Locally linear embedding* and *Laplacian eigenmaps* rely on an iterative update of the lattice structure (i.e the lattice is updated while the method is running).

It is common to compare non linear latent variable models on some known manifold benchmarks. Here we will use the *S curve*, the *severed sphere* and the *Swiss roll* as those models are readily available in `scikit-learn`. Those three models are shown in Fig. 6 below.

Before discussing each of the dimension reduction approaches listed above, it is also worth noting that one of these methods, *multidimensional scaling (MDS)* cannot handle non linear embedding at least in its classical metric formulation. However, it is usually considered as the antecedent of non linear distance preserving methods [12] whence our choice to discuss this method in this chapter.

15.1 Multidimensional scaling

The Term *Multidimensional scaling* [10, 11] refers to a family of methods which construct a lower dimensional representation of the dataset from information on interpoint distances. MDS has been widely used in the humanities such as psychology, economics or sociology in which it is particularly interesting because it can be used with a notion of similarity between points. Unlike Self Organizing maps or principal curves and surfaces, MDS does not require interpoint distances but can be applied from a similarity measure only. This was then used as a quantitative approach at separating concepts.

For a given dataset $\{\mathbf{x}_i\}_{i=1}^N \subseteq \mathbb{R}^p$, we let d_{ij} to denote the distance between any two points \mathbf{x}_i and \mathbf{x}_j in the dataset,

$$d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\| \quad (323)$$

The classical multidimensional scaling approach preserves inner products instead of distances. I.e in classical multidimensional scaling, we start with similarities s_{ij} . If not given similarities explicitly, we use pairwise inner products. The *classical stress function* is thus defined as

$$E_C(\mathbf{z}_1, \dots, \mathbf{z}_N) = \sum_{i,j} (s_{i,j} - \langle \mathbf{z}_i - \bar{\mathbf{z}}, \mathbf{z}_j - \bar{\mathbf{z}} \rangle) \quad (324)$$

Where the similarity is thus often given by the centered inner products $s_{ij} \equiv \langle \mathbf{x}_i - \bar{\mathbf{x}}, \mathbf{x}_j - \bar{\mathbf{x}} \rangle$. It is naturally also possible to consider a multidimensional scaling formulation which would look for points $\mathbf{z}_1, \dots, \mathbf{z}_N \in \mathbb{R}^k$ that minimize the (metric) *stress function*

$$E_M(\mathbf{z}_1, \dots, \mathbf{z}_N) = \sum_{i \neq j} (d_{ij} - \|\mathbf{z}_i - \mathbf{z}_j\|)^2 \quad (325)$$

In this case, the minimization is known as *least squares* or *Kruskal-Shephard* scaling. Moreover, if the distances are Euclidean, they can be converted to centered inner product as follows. Let $d_{i,j}^2 = \|\mathbf{x}_i - \mathbf{x}_j\|^2$, then we have

$$d_{i,j}^2 = \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 + \|\mathbf{x}_j - \bar{\mathbf{x}}\|^2 - 2\langle \mathbf{x}_i - \bar{\mathbf{x}}, \mathbf{x}_j - \bar{\mathbf{x}} \rangle. \quad (326)$$

Now let $\mathbf{M} = \frac{1}{N} \mathbf{1}\mathbf{1}^*$ and $B_{ij} = -\frac{d_{ij}^2}{2}$. The matrix \mathbf{B} can read as

$$\tilde{\mathbf{S}} = (\mathbf{I} - \mathbf{M})\mathbf{B}(\mathbf{I} - \mathbf{M}) \quad (327)$$

Where $(\tilde{\mathbf{S}})_{i,j} = \langle \mathbf{x}_i - \bar{\mathbf{x}}, \mathbf{x}_j - \bar{\mathbf{x}} \rangle$.

A variation of this scaling which minimizes a reweighted version of (325) is known as *Sammon* mapping.

$$E_S(\mathbf{z}_1, \dots, \mathbf{z}_N) = \sum_{i \neq j} \frac{(d_{i,j} - \|\mathbf{z}_i - \mathbf{z}_j\|)^2}{d_{ij}} \quad (328)$$

The classical formulation (324) is particularly attractive not only because inner products are easy to compute but also because Classical MDS minimizes the same criterion as PCA which was discussed earlier. Juste like PCA, classical MDS relies on a simple generative model $\mathbf{z} = \mathbf{W}\mathbf{x}$ where $\mathbf{W} \in \mathbb{R}^{d \times N}$ and is such that $\mathbf{W}^T \mathbf{W} = \mathbf{I}$. For a given dataset $\{\mathbf{x}_i\}$, let $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_d]$ and $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$. The matrix encoding the inner products (a.k.a *Gram matrix*) then reads as

$$\mathbf{S} = \mathbf{Z}^T \mathbf{Z} \quad (329)$$

$$= (\mathbf{W}\mathbf{X})^T (\mathbf{W}\mathbf{X}) \quad (330)$$

$$= \mathbf{X}^T \mathbf{X} \quad (331)$$

Eqs. 329 to (331) show that the latent variables can be found in a very convenient way by simply computing the eigenvalue decomposition of the matrix \mathbf{S} encoding the inner products.

$$\mathbf{S} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T \quad (332)$$

$$= (\mathbf{\Lambda}^{1/2}\mathbf{U})^T (\mathbf{\Lambda}^{1/2}\mathbf{U}^T) \quad (333)$$

From this eigenvalue decomposition, the embedding is computed by retaining the first p eigenvalues (Note that as we are in a N -dimensional space, considering a dataset of size $m > N$, at most N eigenvalues in $\mathbf{\Lambda}$ are non zero),

$$\hat{\mathbf{Z}} = \mathbf{I}_{p \times N} \mathbf{\Lambda}^{1/2} \mathbf{U}^T \quad (334)$$

In (334), the general idea is summarized by Algorithm 2 below (the algorithm appears in the version below in [12]).

```

1 if available data consists of points  $\mathbf{x}_i$  then
2   | Stack them into a matrix  $\mathbf{Y}$  then center them, compute the pairwise inner
   | products  $\mathbf{S} = \mathbf{Y}^T \mathbf{Y}$  and go to third step;
3 else if available data consists of pairwise Euclidean distances then
4   | Square the distances and build the matrix  $\mathbf{B}$  (see discussion above);
5   | Perform double centering of  $\mathbf{B}$  to get  $\tilde{\mathbf{S}}$  following (327).;
6 Compute the eigenvalue decomposition  $\tilde{\mathbf{S}} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^*$ ;
7 The  $p$ -dimensional representation is obtained by computing the product

```

$$\hat{\mathbf{Z}} = \mathbf{I}_{p \times N} \mathbf{\Lambda}^{1/2} \mathbf{U}^* \quad (335)$$

;

Algorithm 2: Multidimensional Scaling

15.2 Isometric feature mapping (ISOMAP)

Isometric feature mappings (ISOMAP) [23] is the simplest non linear dimensionality reduction approach that uses graph distances as an approximation to geodesic distances [12]. In fact, the original ISOMAP follows from replacing the Euclidean distance in Classical multidimensional scaling by a graph distance. The complete Isometric feature mapping algorithm has three steps [23]

- During the first step, the method determines which points are neighbors on the manifold based on pairwise distances between points in the original/input space. The two most common approaches are to consider the points lying in a given Euclidean ball around each point \mathbf{x}_i , or to simply consider the neighborhood to be the K nearest neighbors. The method then builds a graph \mathcal{G} by connecting neighbouring points with edges of weights $w_{ij} = d(i, j)$.
- The second step estimates the pairwise geodesic distances $d_{\mathcal{M}}(i, j)$ of the prototypes by computing the shortest path distances in the graph \mathcal{G} .
- Finally, the last step applies classical Multidimensional scaling to the matrix of graph distances $D_{i,j} = (d_{\mathcal{G}}(i, j))$

The lower dimensional representation generated through ISOMAP is shown on the sphere, S-curves and Swiss Roll in Fig. 8

15.3 Locally linear embedding

Locally linear embedding is a topology preserving method which is originally due to Roweis and Saul [19]. In LLE the idea is to define a lower dimensional embedding but keep the local affine structure of the data. We say that LLE is based on conformal mappings [12]. A *conformal* or *biholomorphic* map is mapping that preserves local angles. The first step in LLE is to determine which angles to keep when defining the lower dimensional representation of the data. For this, LLE first finds the K -nearest neighbors to each point \mathbf{x}_i from the original dataset (other techniques are possible).

Locally linear embedding then relies on the assumption that the manifold underlying the data is locally (i.e at the level of the local K -neighborhood) linear (which seems a reasonable assumption when the dataset is large enough and not too noisy). Following this assumption, LLE then replaces each point \mathbf{x}_i by a linear combination of its K neighbors. After computing the neighborhoods, the second step in LLE then computes the linear combination of the K eighbors that best approximates every given point \mathbf{x}_i from the dataset,

$$\min_w \|\mathbf{x}_i - \sum_{k \in \mathcal{N}(\mathbf{x}_i)} w_{i,k} \mathbf{x}_k\|^2. \quad (336)$$

Here we let $\mathcal{N}(\mathbf{x}_i)$ to denote the k nearest neighbors of \mathbf{x}_i . The low dimensional representation is then given by

$$\min_{\mathbf{z}} \|\mathbf{z}_i - \sum_{k \in \mathcal{N}(\mathbf{z}_i)} w_{i,k} \mathbf{z}_k\|^2. \quad (337)$$

The low dimensional representation thus maintains the local neighborhoods. Eq. (337) is equivalent to solving the matrix equation

$$\text{Tr} [\mathbf{Z} - \mathbf{WZ}]^T (\mathbf{Z} - \mathbf{WZ}) = \text{Tr} [\mathbf{Z}^T (\mathbf{I} - \mathbf{W}) (\mathbf{I} - \mathbf{W}) \mathbf{Z}] \quad (338)$$

for which the solution is given explicitly by the eigenvectors of the sparse symmetric PSD matrix $\mathbf{M} = (\mathbf{I} - \mathbf{W})^T (\mathbf{I} - \mathbf{W})$. The best embedding of dimension p is then given by the lowest p eigenvectors. In fact the first eigenvector is usually discarded by LLE because it corresponds to the trivial $\mathbf{1}$ eigenvector encoding the free translation mode. When discarded, this eigenvector enforces a zero mean embedding.

15.4 Self Organizing maps

The idea underlying Self Organizing map goes back to the work of von der Malsburg in 1973. The idea which was then developed as a way to better understand the visual cortex then remained under the radar until the 80's, when Kohonen came with a simplification of the ideas of von der Malsburg and provided an easy implementation. Because of their relative simplicity, SOMs are used in a variety of areas including time series prediction and data visualization.

Self organizing maps are sometimes build upon the idea of Principal component analysis by extending these ideas to non linear data.

In classical quantization, a smaller set of prototypes is used to approximate (i.e to fit) a larger set of points from the original space. However, the prototypes are free to move independently of each other. The idea of SOMs is similar except that in SOM all prototypes are connected to each other through a lattice (i.e the prototypes are living on a 2D lattice) also called *Constrained topological map*. Each time a prototype is updated and moved, its neighbors on the lattice move accordingly. By doing this, we ensure that the points in the dataset corresponding to points that are close on the lattice, will be close to each other as well. Self organizing maps can thus be intuitively explained as starting from a two dimensional rectangular grid of prototypes $\mathbf{z}_j \in \mathbb{R}^N$. The prototypes can be initialized for example as lying on the two dimensional principal component plane of the data (the prototypes are then often chosen to be equispaced on that plane). The SOM procedure then slowly bends the plane/lattice to match it as much as possible with the original data.

In the case of SOM, the lattice thus plays the role of the embedding space/manifold. SOM runs through the dataset multiple times (performing multiple *epochs*). Each epoch consists of the following two steps being applied for every points \mathbf{x}_i in the considered subset \mathcal{S} [12].

- For any given point \mathbf{x}_i from the original dataset, the first step determines the index of the closest point on the lattice,

$$\gamma^* = \underset{\gamma}{\text{argmin}} d(\mathbf{x}_i, \mathbf{z}(\gamma)) \quad (339)$$

Here d is the Euclidean distance.

- The second step then updates the prototypes $\mathbf{z}(\gamma)$ as

$$\mathbf{z}(\gamma) \leftarrow \mathbf{z}(\gamma) + \alpha \mathcal{N}(\gamma, \zeta) (\mathbf{x}_i - \mathbf{z}(\gamma)) \quad (340)$$

The *learning rate* α obeys $0 \leq \alpha \leq 1$ and is often chosen to slowly decrease with the epochs (a typical choice is to decrease α from 1 to 0 over about a thousand iterations). $\mathcal{N}(\gamma, \zeta)$ is called the *neighborhood function*. This function is defined by a (threshold) parameter $\mathcal{N} = \mathcal{N}_\theta$ and was defined in the original von der Malsburg paper [24] to be the *Bubble function*,

$$\mathcal{N}(\gamma_1, \gamma_2; \theta) = \begin{cases} 0 & \text{when } d_{\mathcal{G}}(\gamma_1, \gamma_2) > \theta \\ 1 & \text{when } d_{\mathcal{G}}(\gamma_1, \gamma_2) \leq \theta \end{cases} \quad (341)$$

It is also sometimes defined from the gaussian kernel, as

$$\mathcal{N}_\theta(\gamma_1, \gamma_2) = \exp\left(-\frac{d_{\mathcal{G}}^2(\gamma_1, \gamma_2)}{2\theta^2}\right) \quad (342)$$

The distance $d_{\mathcal{G}}$ is often chosen to be any distance defined on the lower dimensional space, i.e., $d_{\mathcal{G}}(\gamma_1, \gamma_2) = d(\gamma_1, \gamma_2)$. The threshold parameter θ is also often chosen to decrease with the epochs from a starting value θ_0 to 1 (over about a thousand iterations).

At the end of the iterations, the mapping $f : \mathbf{x} \mapsto \mathbf{f}(\mathbf{x}) = \mathbf{z}(\gamma_{1,\mathbf{x}}, \dots, \gamma_{P,\mathbf{x}}) = \mathbf{z}(\gamma_{\mathbf{x}})$ from any point of the original dataset to the points on the lattice is obtained by simply taking the nearest prototype on the lattice, i.e.

$$f(\mathbf{x}_i) = \mathbf{z}(\gamma) \quad \text{with} \quad \gamma = \underset{\gamma}{\operatorname{argmin}} d(\mathbf{x}_i, \mathbf{z}(\gamma)). \quad (343)$$

The general idea is summarized in Algorithm 3 (the version appears in [12]).

```

1 Set  $k = 0$  and define the original lattice by choosing an inital  $p$ -dimensional
  subspace (for example the PCA plane) and assigning indices to the prototypes
  on the lattice;
2 Set the values of  $\alpha$  and  $\lambda$  for epoch  $k$ ;
3 for all points  $\mathbf{x}_i$  in the dataset do
4   | Compute the index  $\gamma$  of the closest point on the lattice as in (339); Then
   | update the prototypes as in (340);
5 end
6 if magnitude of the prototype update is small then
7   | Stop
8 end
9 else
10  | Set  $k \leftarrow k + 1$  and repeat steps 2 and 3
11 end

```

Algorithm 3: Self Organizing Maps

To be complete, although we only discuss the batch version in these notes, there exist both batch and online versions of the self organizing map algorithm.

16 Reinforcement learning

The main idea underlying Reinforcement Learning (RL) is that we want the agent to learn, without a teacher, but by interacting with the environment. Reinforcement learning is perhaps the approach which is the closest to how learning occurs in human beings.

Reinforcement learning relies on 5 key concepts [22]:

- The learning agent

- The environment
- The policy
- The reward
- The value

The **agent** interacts with its **environment** to achieve a goal. Therefore the agent must be able to sense the state of the environment and to take actions that affect that state. Through its repeated interactions with the environment, the agent progressively improves its performance.

The behavior of the agent is defined by the **policy**. The policy can be thought of as a mapping $\pi : E \rightarrow A$ between the description of the environment and the set of possible actions.

Reinforcement learning also associates to every action, an immediate **reward** which is a direct indication on whether an action taken by the agent was good or bad.

Although the goal of the RL agent is the maximization of the combined rewards, maximizing the immediate reward might not always be the optimal choice as it will constraint the agent to keep repeating the same action again and again, without any opportunity to “explore” new possibilities. As an illustration of this, imagine that you have just moved into a new neighborhood and you would like to find the best coffee place in that particular neighborhood. One approach would be to investigate a couple of places (lets say two) and then settle down in the place that offers the best coffee. You then continue going to that same place again and again. Drinking this coffee makes you feel good and we thus assign it a positive +1 reward. But because you never tried any other places besides the first two, you have no idea whether there exists or not a place that offers a better coffee. In particular, by further investigating, you might discover a third place, where the coffee is tasting even better than place two (i.e. corresponds to a higher, lets say +2 reward). This idea is at the core of reinforcement learning. To formalize it, in addition to the immediate reward, we consider the notion of **value**. The **value** is in some sense, a prediction of the total reward that an agent could get in the future by choosing a particular action. To go back to our coffee example, an action that has a low immediate rewards (such as taking the risk of leaving the first coffee place) might have a higher value because it is possibly followed by other actions that will have a higher reward (finding a third (much) better coffee place).

Together, the five elements above also highlight the difference between reinforcement learning on the one hand, and supervised learning or unsupervised learning on the other. In supervised learning, the knowledge is provided by an external supervisor sometimes called the “teacher” who exactly knows what the correct data should be. On that aspect, reinforcement learning is closer to practical situations, in which it is often impractical to obtain examples of desired behavior. The RL agent takes the action which he assumes to be the most beneficial, yet it has no idea what the optimal action is. In fact, the agent receives a feedback from the environment but this feedback does not give him an exact “yes or no” indication, such as the labels in supervised learning.

Reinforcement learning is also different from unsupervised learning, in the sense that it does not try to reveal structure from unlabeled data. Finding structure in the data does not provide well defined rewards such as those used to guide the agent in RL.

16.1 Multi-armed bandit

The k -bandit problem is one of the simplest reinforcement learning problem. In the k -Bandit problem, the agent has to make a choice between k possible actions, each associated to particular reward which is described by a stationary distribution depending on each particular action. The objective of the agent is then to maximize the total reward over a given period of time. The simplest possible decision model for the value is then to take the expected reward for each action, i.e.

$$v(a) = \mathbb{E} \{R|a\} \tag{344}$$

If we could compute this expectation, we would select at each step, the action which maximizes this value.

In practice, we can only estimate the average by keeping track of the reward we get each time we select a particular action. I.e. at each step, the agent (provided that it chooses action a) updates the empirical value for action a as

$$v_{\text{est}}(a) = Q_t(a) = \frac{\text{sum of rewards when } a \text{ is taken prior to } t}{\text{number of times } a \text{ was selected prior to } t} = \frac{\sum_{i=1}^{t-1} R_i \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}} \quad (345)$$

Here $\mathbb{1}_{A_i=a}$ is defined as

$$\mathbb{1}_{A_i=a} = \begin{cases} 1 & \text{if } A_i = a \\ 0 & \text{otherwise} \end{cases} \quad (346)$$

Given Q_t , the simplest action selection procedure is to take at each step the action which has the highest $Q_t(a)$. I.e.,

$$A_t = \underset{a}{\operatorname{argmax}} Q_t(a) \quad (347)$$

Such a procedure is known as [greedy action selection method](#). Following a greedy policy is guaranteed to accumulate rewards. However, as we saw in the coffee example, it might miss better opportunities by sticking to a limited subset of all possible actions. An alternative consists in following a greedy approach most of the time, yet, once in while (with probability ε), trying out a new action selected at random. This approach is known as [\$\varepsilon\$ -greedy action selection](#)

We can now introduce the [Bandit algorithm](#)

```

1 Initialize, for action  $a = 1$  to  $k$ ;
2  $Q(a) \leftarrow 0$ ;
3  $n(a) \leftarrow 0$ ;
4 Repeat;
5    $A \leftarrow \begin{cases} \underset{a}{\operatorname{argmax}} Q_t(a) & \text{with probability } 1 - \varepsilon \\ \text{a random action} & \text{with probability } \varepsilon \end{cases}$ ;
6    $n(a) \leftarrow n(a) + 1$ ;
7    $Q(a) \leftarrow Q(a) + \frac{1}{n(a)}[R(a) - Q(a)]$ ;
```

Algorithm 4: Bandit algorithm, see [22]

The last line in this algorithm is simply the update of the average. Assuming that the average reward for action a is $Q_t(a)$, and that the agent selects action a , we update this average as

$$Q_t(a) = \frac{1}{n(a)} \sum_{i=1}^{n(a)} R_i(a) \quad (348)$$

$$= \frac{1}{n(a)} \left(R_n(a) + \sum_{i=1}^{n(a)-1} R_i(a) \right) \quad (349)$$

$$= \frac{1}{n(a)} \left(R_n(a) + (n(a) - 1) \frac{1}{n(a) - 1} \sum_{i=1}^{n(a)-1} R_i(a) \right) \quad (350)$$

$$= \frac{1}{n(a)} (R_n(a) + (n(a) - 1) Q_{t-1}(a)) \quad (351)$$

$$= \frac{1}{n(a)} (R_n(a) + n(a) Q_{t-1}(a) - Q_{t-1}(a)) \quad (352)$$

$$= Q_{t-1}(a) + \frac{1}{n(a)} [R_n(a) - Q_{t-1}(a)] \quad (353)$$

16.2 Q-learning

A major breakthrough in reinforcement learning was the development of Q -learning. Q -learning is a particular instance of "off policy Temporal difference learning". In Q -learning, the agent learns the **action-value function** Q , which is a mapping between (state, action) pairs and corresponding estimated values. I.e. during training, the agent learns to associated particular states of the environment and particular actions with a corresponding value estimate. At each step, one action is selected in an ε -greedy fashion (i.e either the agent chooses the action by relying on the Q -table, or its takes the action to be random)

The action value table is then updated as follows,

$$Q[S_t, A_t] \leftarrow Q[S_t, A_t] + \alpha \left[R_{t+1} + \gamma \max_a Q[S_{t+1}, a] - Q[S_t, A_t] \right] \quad (354)$$

The table is thus updated as before, except that we now consider an additional term which takes into account the possible future increase in the value. I.e. for a given possible action choice a , we look up all the possible subsequent actions in the table. If the action a change the environment to state S'_t that means looking at all possible actions the agent can take on the environment in that particular state, so all entries of the form $Q[S'_t, a]$. Then we take our potential future gain to be the maximal value. In Q -learning, we thus try to be one step ahead compared to the traditional bandit algorithm which only considers the current value estimate.

In general, we will not learn the action-value table Q , but we will replace it by a neural network which will approximate it.

References

- [1] Andrew D Back and Andreas S Weigend. A first application of independent component analysis to extracting structure from stock returns. *International journal of neural systems*, 8(04):473–484, 1997.
- [2] Christopher M Bishop et al. Pattern recognition and machine learning (information science and statistics). 2006.
- [3] Leon Bottou and Yoshua Bengio. Convergence properties of the k-means algorithms. In *Advances in neural information processing systems*, pages 585–592, 1995.
- [4] George Casella and Roger L Berger. *Statistical inference*, volume 2. Duxbury Pacific Grove, CA, 2002.
- [5] Edward W Forgy. Cluster analysis of multivariate data: efficiency versus interpretability of classifications. *biometrics*, 21:768–769, 1965.
- [6] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- [7] Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications. *Neural networks*, 13(4-5):411–430, 2000.
- [8] Željko Ivezić, Andrew J Connolly, Jacob T VanderPlas, and Alexander Gray. *Statistics, data mining, and machine learning in astronomy: a practical Python guide for the analysis of survey data*, volume 1. Princeton University Press, 2014.
- [9] Leonard Kaufman and Peter J Rousseeuw. *Finding groups in data: an introduction to cluster analysis*, volume 344. John Wiley & Sons, 2009.
- [10] Joseph B Kruskal. Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. *Psychometrika*, 29(1):1–27, 1964.
- [11] Joseph B Kruskal. Nonmetric multidimensional scaling: a numerical method. *Psychometrika*, 29(2):115–129, 1964.

- [12] John A Lee and Michel Verleysen. *Nonlinear dimensionality reduction*. Springer Science & Business Media, 2007.
- [13] Stuart Lloyd. Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2):129–137, 1982.
- [14] J MacQueen. On convergence of k-means and partitions with minimum average variance. In *Annals of Mathematical Statistics*, volume 36, page 1084. inst mathematical statistics IMS business office-suite 7, 3401 Investment blvd, Hayward, CA 94545, 1965.
- [15] James MacQueen et al. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA, 1967.
- [16] Kevin P. Murphy. *Machine learning, a probabilistic perspective*, 2012.
- [17] Jorge Nocedal and Stephen Wright. *Numerical optimization*. Springer Science & Business Media, 2006.
- [18] José M Pena, Jose Antonio Lozano, and Pedro Larranaga. An empirical comparison of four initialization methods for the k-means algorithm. *Pattern recognition letters*, 20(10):1027–1040, 1999.
- [19] Sam T Roweis and Lawrence K Saul. Nonlinear dimensionality reduction by locally linear embedding. *science*, 290(5500):2323–2326, 2000.
- [20] Michael D Spivak. *A comprehensive introduction to differential geometry*. Publish or perish, 1970.
- [21] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [22] Richard S Sutton, Andrew G Barto, et al. *Introduction to reinforcement learning*, volume 135. MIT press Cambridge, 1998.
- [23] Joshua B Tenenbaum, Vin De Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *science*, 290(5500):2319–2323, 2000.
- [24] Chr Von der Malsburg. Self-organization of orientation sensitive cells in the striate cortex. *Kybernetik*, 14(2):85–100, 1973.

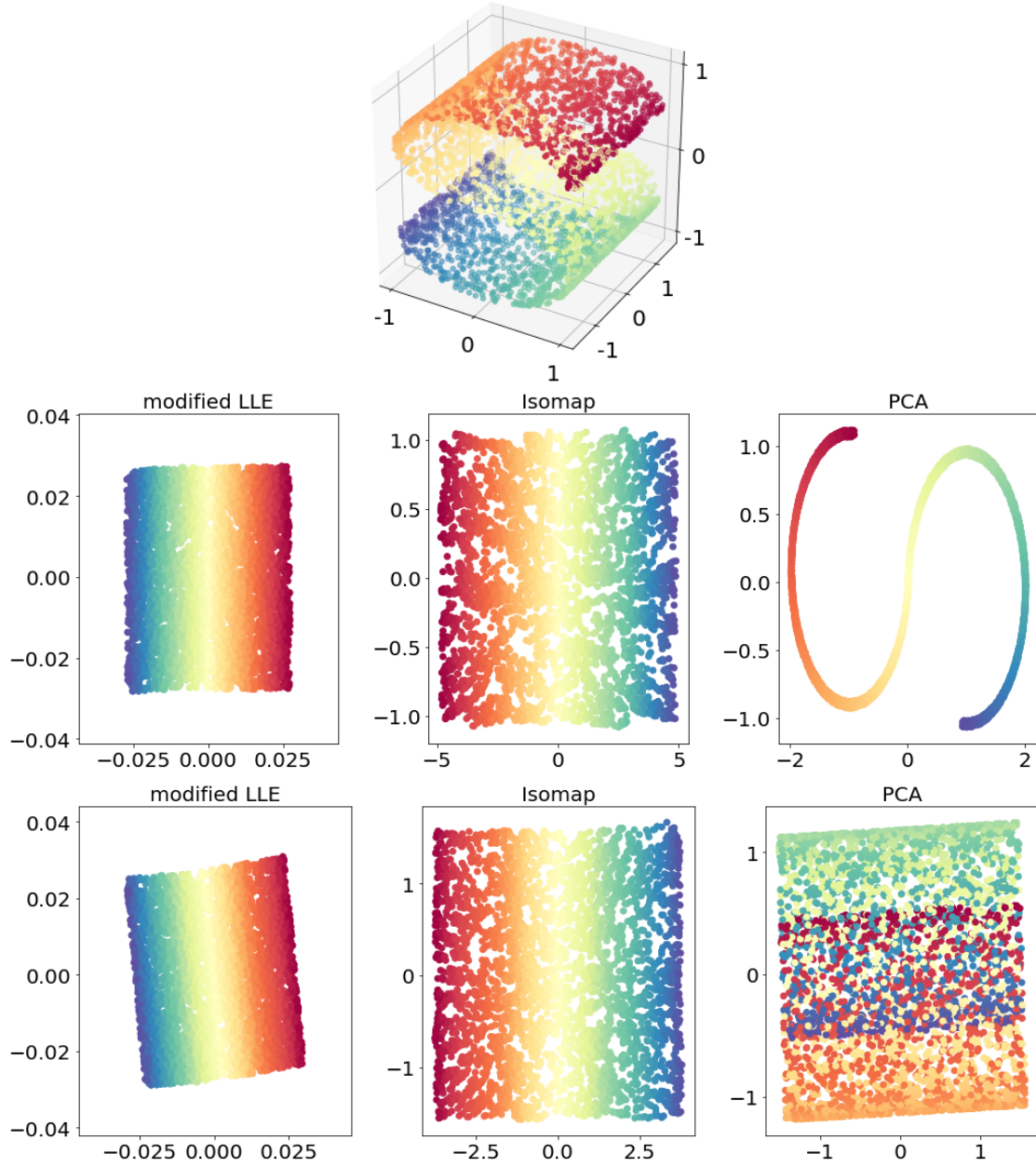


Figure 5: Various two dimensional representations of the S-shaped manifold from scikit-learn. The first row corresponds to two dimensional representations of the original manifold. In the second row, this manifold was modified following [8] to further highlight the weakness of PCA.

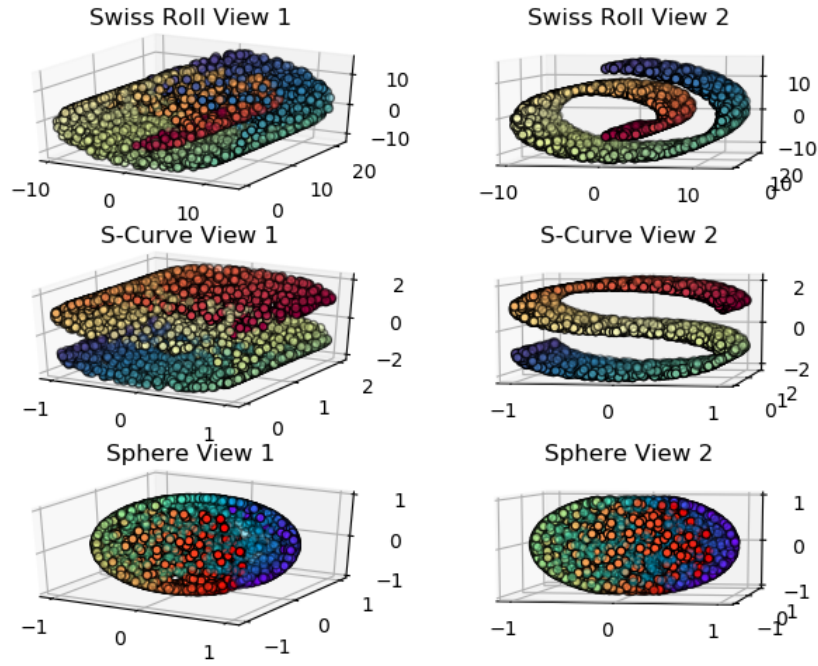


Figure 6: Usual benchmarks used for the comparison of non linear latent variable models. From the top to bottom: Swiss roll, S-manifold and sphere.

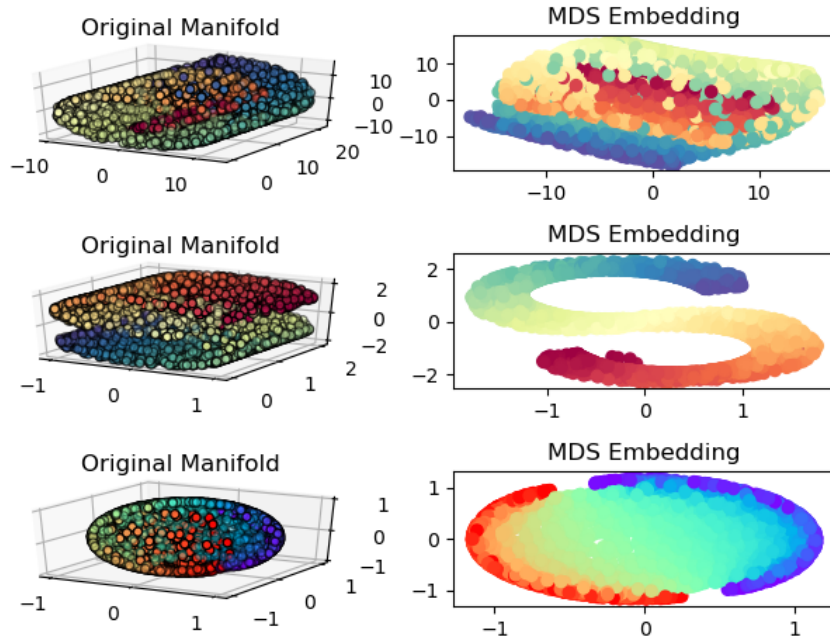


Figure 7: Two dimensional representation obtained by Multidimensional scaling on the *Swiss Roll*, the *S-curve* and the *Sphere* manifolds.

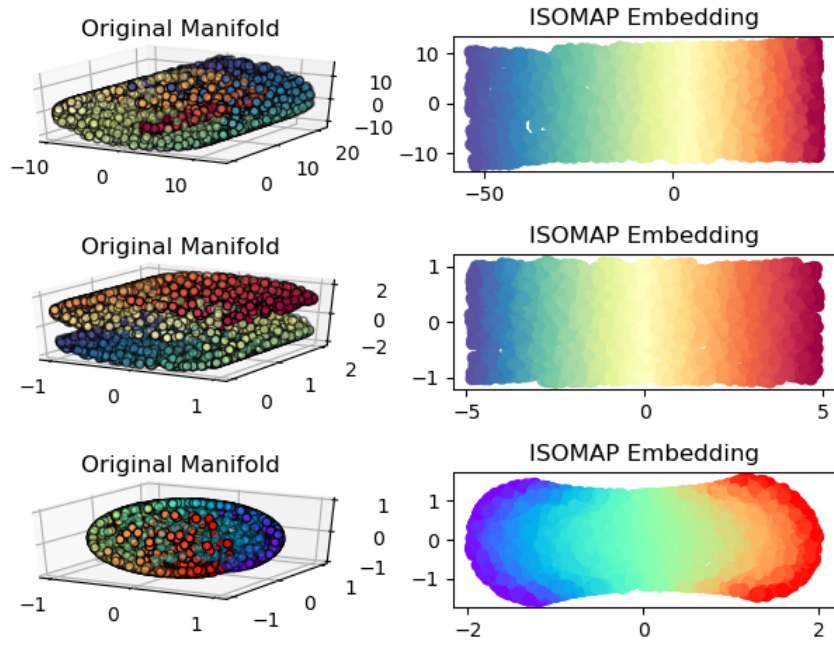


Figure 8: Two dimensional representation obtained by ISOMAP on the *Swiss Roll*, the *S-curve* and the *Sphere* manifolds.

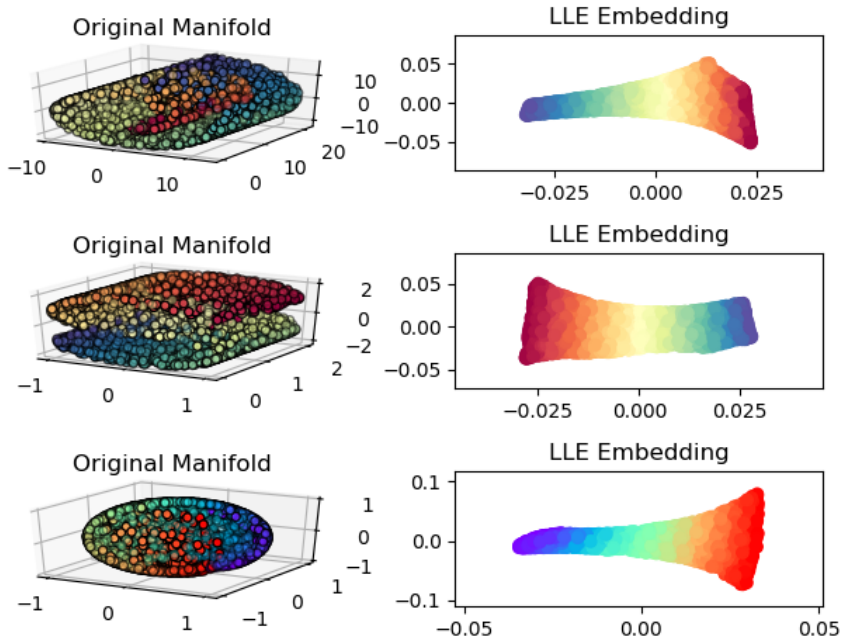


Figure 9: Two dimensional representation obtained by locally linear embedding on the *Swiss Roll*, the *S-curve* and the *Sphere* manifolds.