

Introduction to Machine Learning

Alessandro Rudi, Umut Simsekli

Notes by Antoine Groudiev

31st March 2024

Contents

1	An overview of Machine Learning	2
1.1	What is ML?	2
1.2	Topics in Machine Learning	2
1.2.1	Supervised Learning	2
1.2.2	Probabilistic approach	3
1.2.3	Unsupervised Learning	3
2	Linear Least Squares Regression	3
2.1	Introduction	3
2.2	General setup and notations for supervised learning	4
2.3	Ordinary Least Squares Estimator (OLS)	5
2.4	Statistical analysis	5
2.5	Ridge regression	5
3	Logic regression and convex analysis	5
3.1	6
4	Convex analysis and convex optimization	7
4.1	Constrained optimization problems	7
4.2	Optimization algorithms for unconstrained convex optimization	7
4.2.1	Gradient Descent	7
4.2.2	Stochastic Gradient Descent	7
5	Kernels	7
5.1	Introduction to kernels	7
5.2	Representer theorem	8
5.2.1	Theorem statement	8
5.2.2	Finite dimensional representation of the learning problem	8
5.3	Properties of kernels	9
6	Elements of Statistical Machine Learning	10
7	Model-Based Machine Learning	10
8	Maximum Likelihood	10
9	Unsupervised Learning	10

10 MCMC Sampling	10
11 Neural Networks	10

Introduction

This document is Antoine Groudiev's class notes while following the class *Introduction to Machine Learning* (Apprentissage Statistique) at the Computer Science Department of ENS Ulm. It is freely inspired by the class notes of Pierre Gaillard Alessandro Rudi, and Umut Şimşekli.

1 An overview of Machine Learning

1.1 What is ML?

Considering a problem, such as image classification: given an input image of a dog or a cat, the program is asked to determine whether the image is a dog or a cat. Conventional programming would hardcode the solution to this problem. But this process takes time and is not easily generalisable. Instead, an ML model is trained on a dataset to produce a program to solve the problem.

Many successful applications of Machine Learning are:

- Face recognition
- Spam filtering
- Speech recognition
- Self-driving systems; pedestrian detection

1.2 Topics in Machine Learning

1.2.1 Supervised Learning

Example (Classification). Features $x \in \mathbb{R}^d$, labels $y \in \{1, \dots, k\}$

Definition (Regression). Features $x \in \mathbb{R}^d$, labels $y \in \mathbb{R}$. To tackle such problem, we look for a parametrized function $f_\theta(x_i) \simeq y_i$ for some f_θ in a function space

$$\mathcal{F} = \{f_\theta : \theta \in \Theta\}$$

Our goal is therefore to find the best function in \mathcal{F} such that f "fits" the training data. For example, we can say that f "fits" the training data when

$$\frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2$$

is "small". Such a function is not interesting in general, like for classification.

Definition (Loss function). Assumes that the features are in \mathcal{X} and the labels are in \mathcal{Y} . We introduce the more general *loss function* notion:

$$l : \mathcal{Y}^2 \rightarrow \mathbb{R}_+$$

For a regression task, we can use $l(\hat{y}, y) = (\hat{y} - y)^2$. For a classification task, $l(\hat{y}, y) = \mathbb{1}_{\hat{y} \neq y}$.

Therefore, for a regression problem, we might choose:

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i)$$

In the parametric case, when $\mathcal{F} = \{f_\theta : \theta \in \Theta\}$, we might minimize with respect to θ :

$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n l(f_\theta(x_i), y_i)$$

1.2.2 Probabilistic approach

Let $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ be the feature space. Let D be a distribution on \mathcal{Z} ; we make the assumption that the training data is iid from D :

$$(x_i, y_i) \sim D$$

and the same thing hold for the test data:

$$(\tilde{x}_i, \tilde{y}_i) \sim D$$

According to the Strong Law of Large Numbers, the test loss converges almost surely:

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n l(f_\theta(\tilde{x}_i), \tilde{y}_i) = \mathbb{E}_{(x,y) \sim D}[l(f_\theta(x), y)] =: R(\theta) = R(f_\theta)$$

where $R(\theta)$ is the *population risk*.

Definition (Risk minimization).

1.2.3 Unsupervised Learning

Example (Clustering).

Example (Dimensionality reduction). We are given features $x \in \mathbb{R}^d$ and labels $y \in \{0, 1\}$ which form a "training" dataset $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$. We assume that $d \gg 1$; our goal is to find $d' \ll d$ such that $(x_1, y_1, \dots,)$

2 Linear Least Squares Regression

2.1 Introduction

In this chapter, we will study the simple but still widely used problem of *Linear Least Square Regression*. We are given a set of points, which we assume to be sampled from some distribution: there exists some function which generated these points, and we want to retrieve or at least approximate this unknown function. To do so, we will naturally look for the function which best fits the points; nevertheless, assuming that it is unlikely that the function is overly-complicated, we will only approximate it using *linear function*. Finally, to choose which linear function "fits best" the data, we will introduce the mean square error, which we will minimize to find our linear approximation function.

Formally, our objective is to find a function f such that it explains well the distribution $(y_i)_{1 \leq i \leq n}$ as a function of $(x_i)_{1 \leq i \leq n}$, that is $y_i \sim f(x_i)$. To do this, we can choose a *function space* \mathcal{F} and solve the empirical risk minimization problem:

$$\hat{f}_n \in \operatorname{argmin}_{f \in \mathcal{F}} \hat{R}_n(f) := \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

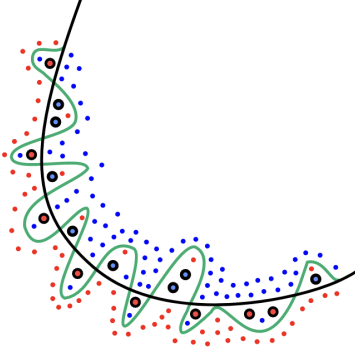


Figure 1: Example of overfitting

Care must be taken when selecting the function space to avoid overfitting¹: for example, if we were to choose $\mathcal{F} := \mathbb{R}_d[X]$ for some $d \in \mathbb{N}$, it is in our best interest to keep d small to avoid getting a function f which fits perfectly the points but diverges between them. Although the empirical mean square error \hat{R}_n decreases when the function space \mathcal{F} becomes larger (i.e. larger polynomial degrees), the \hat{f}_n estimator loses its predictive power. \hat{f}_n will not necessarily perform well on new data. In what follows, we will consider the linear function space, containing functions of the form $f : x \mapsto ax + b$, which is the simplest.

2.2 General setup and notations for supervised learning

Definition (Training data set). The *training data set*, often denoted $D_n := \{(x_i, y_i) \mid i \in \llbracket 1, n \rrbracket\}$, is the set of some observations $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$. We will often make the assumption that the observations (x_i, y_i) are realizations of i.i.d. random variables from a distribution ν .

The distribution ν is unknown to the statistician; it's a matter of learning it from the data.

Definition (Learning rule). A *learning rule* \mathcal{A} is a function that associates to training data D_n a prediction function \hat{f}_n :

$$\mathcal{A} : \bigcup_{n \in \mathbb{N}} (\mathcal{X} \times \mathcal{Y})^n \longrightarrow \mathcal{Y}^{\mathcal{X}}$$

$$D_n \longmapsto \hat{f}_n$$

The estimated function \hat{f}_n is constructed to predict a new output y from a new input x , where (x, y) is a pair of *test data*, i.e. not necessarily observed in the training data. The function \hat{f}_n is an estimator because it depends on the data D_n and not on unobserved parameter, such as the distribution ν . If D_n is random, it is also a random function.

Definition (Squared Loss Risk). Given an estimator \hat{f}_n , we define its risk:

$$\mathcal{R}(\hat{f}_n) := \mathbb{E} \left[(Y - \hat{f}_n(X))^2 \mid D_n \right] \quad \text{where } (X, Y) \sim \nu \quad (2.2.1)$$

In practice, the statistician cannot compute the risk, since one cannot access the distribution ν . Therefore, a common method in supervised machine learning is to replace the risk (defined using the distribution ν through the expectation \mathbb{E}) by the empirical risk (defined using the training data set).

Definition (Squared Loss Empirical risk). Given an estimator \hat{f}_n and a data set $D_n = \{(x_i, y_i) \mid i \in \llbracket 1, n \rrbracket\}$, we define its *empirical risk*:

$$\hat{\mathcal{R}}_n(f) := \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 \quad (2.2.2)$$

¹We say that the function *overfits* the data when it corresponds too closely to the specific set of data (i.e. on the training data), such that it fails to fit additional data (i.e. test data).

However, one must be careful about overfitting, the case where $\hat{\mathcal{R}}_n(f)$ is much lower than $\mathcal{R}(f)$, as discussed previously. In this chapter, we will study the performance of the least square estimator in the case of the linear model.

Definition (Linear model). When $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \mathbb{R}$, the simplest interesting function space is the set of affine functions. To ease the notation, we assume that the first components of the inputs is 1 so that it is sufficient to consider linear functions. Therefore, the function space is:

$$\mathcal{F} := \left\{ x \mapsto \theta^\top x \mid \theta \in \mathbb{R}^d \right\} \quad (2.2.3)$$

i.e. linear functions parametrized by $\theta \in \mathbb{R}^d$.

Remark. Choosing a linear model, the empirical risk minimization corresponds to the problem of minimizing the following quantity over $\theta \in \mathbb{R}^d$:

$$\hat{\mathcal{R}}_\theta := \frac{1}{n} \sum_{i=1}^n (y_i - \theta^\top x_i)^2$$

This expression can be rewritten using matrix notation. We let $Y = (Y_1, \dots, Y_n)^\top \in \mathbb{R}^n$ be the vector of outputs and $X \in \mathbb{R}^{n \times d}$ the matrix of inputs, whose rows are x_i^\top . X is called the design matrix. The empirical risk is therefore given by:

$$\hat{\mathcal{R}}_n(\theta) = \frac{1}{n} \|Y - X\theta\|_2^2$$

2.3 Ordinary Least Squares Estimator (OLS)

2.4 Statistical analysis

2.5 Ridge regression

3 Logic regression and convex analysis

Recap of important notions and notations

We are given an input space X and an output space Y . We want to learn the relationship between input and output, modelised by a probability distribution $\rho \in \mathbb{P}(X \times Y)$. Thus, we try to find the best function $f_\star : X \rightarrow Y$, given a loss function $l : Y \times Y \rightarrow \mathbb{R}$. Therefore, f_\star is often defined by:

$$f_\star = \operatorname{argmin}_{f: X \rightarrow Y} \mathbb{E}_{X,Y} [l(f(X), Y)]$$

where

$$\mathbb{E}_{X,Y} [g(X, Y)] = \int_{\mathbb{R}^2} g(x, y) \cdot d\rho(x, y)$$

In practice, you only know some samples $D_N = [(x_1, y_1), \dots, (x_N, y_N)]$ with $(x_i, y_i) \sim \rho$, making it impossible to choose such an f_\star . Therefore, we try to find a good model \hat{f}_{D_N} , such that

$$\lim_{N \rightarrow +\infty} \mathcal{E}(\hat{f}_{D_N}) - \mathcal{E}(f) = 0$$

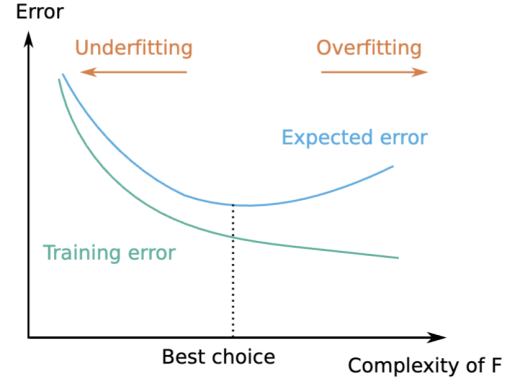


Figure 2: Overfitting and underfitting according to the complexity of \mathcal{F}

Such a result will often be given by a *learning rate function* $c(N)$, with

$$\mathbb{E}_{D_N}[\mathcal{E}(\hat{f}_{D_N}) - \mathcal{E}(f)] \leq c(N) = o(1)$$

The function \hat{f}_{D_N} can be chosen such that it minimizes the empirical error:

$$\hat{f}_{D_N} = \operatorname{argmin}_{f \in \mathcal{H}} \hat{\mathcal{E}}(f) = \operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N l(f(x_i), y_i)$$

3.1

We consider the case where $X = \mathbb{R}^d$ and $Y = \mathbb{R}$. We define the loss l to be the least squares, $l(y, y') = (y - y')^2$, and we choose our functions to be of the form of $f_\star = \theta_\star^T X$. In this case, ERM is OLS:

$$\hat{\theta}_N = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N (\theta^T x_i - y_i)^2$$

We can also define $\hat{\theta}_{N,\lambda}$ to be:

$$\hat{\theta}_{N,\lambda} = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N (\theta^T x_i - y_i)^2 + \lambda \|\theta\|^2$$

This allows to regulate the "complexity" of the function to avoid overfitting. This is called Tikhonov regularization. In this case, we have

$$\mathbb{E}_{\hat{Y}}[\hat{\theta}_N - \mathcal{E}(\theta_\star)] = \frac{\sigma^2 d}{N}$$

and therefore the optimal function is

$$\hat{f}_{N,\lambda} = \operatorname{argmin}_{f \in \mathcal{H}} \hat{\mathcal{E}}(f) + \lambda R(f)$$

We define $X \in \mathbb{R}^{N \times d} := (x_1^T, \dots, x_N^T)$, and $\hat{Y} = (\hat{y}_1, \dots, \hat{y}_n)$. We this notation, we have

$$\hat{\theta}_{N,\lambda} = \frac{1}{N} \|X\theta - \hat{Y}\|^2 + \lambda \|\theta\|^2$$

Thus, we have

$$\begin{aligned} \nabla \mathcal{L}(\theta) &:= \frac{2}{N} X^T X \theta - 2 \frac{X^T \hat{Y}}{N} + 2\lambda \theta = 0 \\ \left(\frac{X^T X}{N} + \lambda \right) \theta &= X^T \hat{Y} \end{aligned}$$

therefore,

$$\hat{\theta}_{N,\lambda} = \left(\frac{X^T X}{N} + \lambda I \right)^{-1} \frac{X^T \hat{Y}}{N} = (X^T X + \lambda N I)^{-1} X^T \hat{Y}$$

We introduce the singular value decomposition of X :

$$X = U \Sigma V^T$$

where $U^T U = U U^T = I_N$, $V^T V = V V^T = I_d$, and Σ is diagonal with $\forall i, \Sigma_{i,i} \geq 0$. In this case,

$$\begin{aligned} X^T X + \lambda N I_d &= V \Sigma U^T U \Sigma V^T + \lambda N I_d \\ &= V \underbrace{(\Sigma^2 + \lambda N I)}_{\text{invertible}} V^T \end{aligned}$$

4 Convex analysis and convex optimization

4.1 Constrained optimization problems

4.2 Optimization algorithms for unconstrained convex optimization

4.2.1 Gradient Descent

4.2.2 Stochastic Gradient Descent

5 Kernels

5.1 Introduction to kernels

In this course, we often focused on prediction methods which are *linear*, that is, the input data are vectors and the prediction function is linear (e.g. $f(x) = w^\top x$ for $w \in \mathbb{R}^d$). In this situation with given data (x_i, y_i) , the vector w can be obtained by minimizing

$$\hat{L}(w) = \frac{1}{n} \sum_{i=1}^n l(y_i, w^\top x_i) + \lambda \Omega(w)$$

Classical examples are logistic regression or least-squares regression. These methods look at first sight of limited practical significance, because input data may not be vectors, and relevant prediction functions may not be linear.

The goal of kernel methods is therefore to go beyond these limitations while keeping the good aspects. The underlying principle is to replace x by a function $\varphi(x) \in \mathbb{R}^d$, *explicitly* or *implicitly*, and consider linear predictions in $\Phi(x)$, i.e. $f(x) = w^\top \varphi(x)$. We call $\varphi(x)$ the *feature* associated to x .

Example (Linear regression). In the case of linear regression, $\varphi(x) = x$ for $x \in \mathbb{R}^d$. As expected, this gives us linear models:

$$f(x) = w^\top x = \sum_{j=1}^d w_j x_j$$

Example (Polynomial regression of degree r). With $x \in \mathbb{R}$, we have $\varphi(x) \in \mathbb{R}^{r+1}$ defined by:

$$\varphi(x) = (1, x, x^2, \dots, x^r)$$

Therefore, the prediction functions will be general polynomials of degree at most r :

$$f(x) = w^\top \varphi(x) = \sum_{j=0}^r (\varphi(x))_j = \sum_{j=1}^r w_j x^j$$

Example (Polynomial multivariate regression of degree r). We consider $x \in \mathbb{R}^d$ and

$$\varphi(x) = (x_1^{\alpha_1}, \dots, x_d^{\alpha_d})$$

with $\sum_{i=1}^d \alpha_i = r$. In this situation, $p = \binom{d+r-1}{r}$ might be too big for an explicit representation to be feasible;

Example (Generic set of functions). Let $\phi_1, \phi_r : \mathbb{R}^d \rightarrow \mathbb{R}$ be a set of functions of interest (e.g. a subset of the Fourier basis); we define $\varphi(x) = (\phi_1(x), \dots, \phi_r(x))$ to have:

$$f(x) = w^\top \varphi(x) = \sum_{j=1}^r w_j \phi_j(x)$$

5.2 Representer theorem

5.2.1 Theorem statement

Given a dataset x_1, \dots, x_n , we are able to compute the observed feature maps $\varphi(x_1), \dots, \varphi(x_n)$. We can ask ourselves if there exists an easier representation for w in terms of these observed feature maps, i.e. we want to know if it is possible to characterize the minimum \hat{w} of:

$$\hat{L}(w) = \frac{1}{n} \sum_{i=1}^n l(y_i, w^\top \varphi(x_i)) + \lambda w^\top w$$

in the form of $\hat{w} = \sum_{i=1}^n \alpha_i \varphi(x_i)$, with $\alpha_i \in \mathbb{R}$. The following theorem guarantees such characterization under basic properties of \hat{L} .

Theorem (Representer theorem). Let $\varphi : \mathcal{X} \rightarrow \mathbb{R}^d$. Let $(x_1, \dots, x_n) \in \mathcal{X}^n$ and assume that $\Psi : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$ is strictly increasing with respect to the last variable. Then, the minimum of

$$\hat{L}(w) := \Psi(w^\top \varphi(x_1), \dots, w^\top \varphi(x_n), w^\top w)$$

is obtained for

$$w = \sum_{i=1}^n \alpha_i \varphi(x_i)$$

for some $\alpha \in \mathbb{R}^n$.

Proof.

□

Corollary. For $\lambda > 0$,

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n l(y_i, w^\top \varphi(x_i)) + \frac{\lambda}{2} w^\top w$$

is obtained for

$$w = \sum_{i=1}^n \alpha_i \varphi(x_i).$$

Note that there is no assumption on l , and in particular, no convexity assumption. This result is extendable to Hilbert spaces (RKHS), as we will see in the next section.

5.2.2 Finite dimensional representation of the learning problem

Using the representer theorem, we know that the minimum of \hat{L} is of the form $w = \sum_{i=1}^n \alpha_i \varphi(x_i)$; we can therefore directly optimize this characterization, i.e. we can then write:

$$\min_{w \in \mathbb{R}^r} \frac{1}{n} \sum_{i=1}^n l(y_i, w^\top \varphi(x_i)) + \frac{\lambda}{2} w^\top w = \min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n l(y_i, (K\alpha)_i) + \frac{\lambda}{2} \alpha^\top K \alpha$$

where K is an $n \times n$ matrix with values

$$K_{i,j} = \varphi(x_i)^\top \varphi(x_j).$$

Indeed,

$$\varphi(x_i)^\top w = \sum_{j=1}^n \alpha_j \varphi(x_i)^\top \varphi(x_j) = (K\alpha)_i$$

moreover,

$$\|w\|^2 = w^\top w = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \varphi(x_i)^\top \varphi(x_j) = \alpha^\top K \alpha$$

We finally have a closed form representation for the function evaluation. Defining the *kernel function* $k(x, x') := \varphi(x)^\top \varphi(x')$, we have:

$$f(x) = w^\top \varphi(x) = \sum_{i=1}^n \alpha_i \varphi(x_i)^\top \varphi(x) = \sum_{i=1}^n \alpha_i k(x_i, x).$$

Remark (Kernel trick). *The whole learning problem can be written in terms of the kernel k ; indeed, f depends only on k , \hat{L} depends on K with $K_{i,j} = k(x_i, x_j)$. Therefore, we have the so-called kernel trick, i.e. we do not need to compute explicitly the features φ to be able to represent and solve the learning problem, we just need to be able to compute their inner product.*

Example (Power of the kernel trick with infinite dimensional feature maps). Consider $\mathcal{X} = [-1, 1]$ and the feature map

$$\varphi(x) = (1, x, x^2, \dots).$$

The resulting model space would have the form

$$f(x) = \sum_{j=0}^{+\infty} w_j x^j$$

with $\sum_{j=0}^{+\infty} w_j^2 < +\infty$. This model space is the set of analytic function on \mathcal{X} , which is a very rich space. In particular, it is dense in the space of continuous functions. However, it is not possible to compute $\varphi(x)$ explicitly since it is infinite dimensional. The kernel trick provides an elegant way to compute the solution of the learning problem in closed form; indeed, the inner product can be computed in closed form in $O(1)$:

$$k(x, x') = \varphi(x)^\top \varphi(x') = \sum_{j=0}^{+\infty} x^j x'^j = \frac{1}{1 - xx'}$$

Therefore, the kernel trick allow to replace \mathbb{R}^d by \mathbb{R}^n , which is interesting when d is very large. Furthermore, it allows to separate the representation problem (design a kernel on a set \mathcal{X}), algorithms, and analysis (which only use the kernel matrix K).

5.3 Properties of kernels

Since the learning problem is completely defined in terms of the kernel function, the explicit knowledge of the feature map is not required anymore. In particular, given a function $k : X \times X \rightarrow \mathbb{R}$, to use it in a learning problem, we need to be sure that it is a *positive definite kernel*, i.e. that there exists a feature map φ such that

$$\forall x, x' \in X, \quad k(x, x') = \varphi(x)^\top \varphi(x')$$

Kernel functions admits many characterizations, which we will now present.

Property 1 (Characterization in terms of positive-definiteness). k is a positive definite kernel if and only if the kernel matrix k is positive semi-definite (i.e. all its eigenvalues are non-negative).

Theorem (Aronszajn). k is a positive definite kernel if and only if there exists a Hilbert space \mathcal{F} , and $\varphi : X \rightarrow \mathcal{F}$ such that

$$\forall x, y \in X, \quad k(x, y) = \langle \varphi(x), \varphi(y) \rangle$$

If such objects exist, \mathcal{F} is called the *feature space* and φ the *feature map*.

Property 2. The sum and product of kernels are kernels.

Example (Linear kernel). The linear kernel corresponds to $\varphi = x \mapsto x$:

$$k(x, y) = x^\top y$$

Example (Polynomial kernel). The kernel $k(x, y) = (x^\top y)^r$ can be expanded as:

$$k(x, y) = \left(\sum_{i=1}^d x_i y_i \right)^r = \sum_{\alpha_1 + \dots + \alpha_p = r} \binom{r}{\alpha_1, \dots, \alpha_p} \underbrace{(x_1 y_1)^{\alpha_1} \dots (x_p y_p)^{\alpha_p}}_{(x_1^{\alpha_1} \dots x_p^{\alpha_p})(y_1^{\alpha_1} \dots y_p^{\alpha_p})}$$

Example (Translation-invariant kernels on a bounded interval).

Example (Translation-invariant kernels on \mathbb{R}^d).

6 Elements of Statistical Machine Learning

7 Model-Based Machine Learning

8 Maximum Likelihood

9 Unsupervised Learning

10 MCMC Sampling

11 Neural Networks