

Introduction to Machine Learning & Deep Learning - Part 2

Sorbonne Université- Master informatique DAC et Master mathématiques et applications M2A-

P. Gallinari, patrick.gallinari@lip6.fr, <http://www-connex.lip6.fr/~gallinar/>

Nicolas Baskiotis, Edouard Oyallon, Laure Soulier, nom.prenom@lip6.fr

Année 2021-2022

Kernels etc.

Kernels

Support Vector Machines

Gaussian Processes

Kernel Methods – a brief introduction

Introducing kernels

- ▶ The concept of kernels is important in machine learning
- ▶ It allows to derive general families of ML methods
 - ▶ Applicable to generic ML problems: supervised, unsupervised, ranking, ..
 - ▶ That can be used on different types of data (vectors, strings, graphs, ...)
- ▶ It provides a general framework for the formal analysis of complex algorithms
 - ▶ e.g. NN in the infinite limit (infinite number of hidden cells) can be modeled and then analyzed as kernel methods
- ▶ Kernels, and over all support vector machines have been one of the main ML paradigm in 1995-2005.
 - ▶ The concept allows to make use and to formalize several important ideas concerning e.g. optimization (convex optimization), generalization
 - ▶ Most kernel methods are not well adapted to high dimensional spaces and large datasets, they failed in this sense but remain an important concept

Intuition (1) – kernels as similarity measures

- ▶ Kernel exploit similarity measures between data representations
 - ▶ Expressed as dot products in a feature space
- ▶ **Feature space** - Let X be a set (e.g. the set of objects to be classified), we will represent these objects in a **feature space** \mathcal{H} , which is a **vector space equipped with a dot product**.
 - ▶ For that we will use a map Φ :
$$\Phi: X \rightarrow \mathcal{H}$$
$$x \mapsto \Phi(x)$$
- ▶ **Similarity measure** - we define a similarity measure via the dot product in \mathcal{H} :
 - ▶ $K(x, x') = \langle \Phi(x), \Phi(x') \rangle$
 - ▶ In the following, $K(\cdot, \cdot)$ will be called a **kernel**
 - ▶ Note: X can be any set, and not only a subset of \mathbb{R}^n
 - i.e. it may be endowed with a dot product itself or not, e.g. think of X as a set of books or proteins
 - Even when $X \subset \mathbb{R}^n$, i.e. a dot product space, the mapping Φ will allow us to define more complex (non linear) representations of $x \in X$

Intuition (2) – machine learning algorithms and dot products

- ▶ Several machine learning algorithms can be expressed using dot products in a feature space
 - ▶ We introduce two simple examples
 - ▶ Perceptron
 - ▶ Linear regression
 - ▶ This idea can be generalized to many families of supervised and unsupervised methods

Intuition (2) – machine learning algorithms and dot products

Example 1: Perceptron dual formulation for binary classification

Training set $D = \{(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^N, y^N)\}$, $\mathbf{x}^i \in \mathbb{R}^n$, $y^i \in \{-1, 1\}$, hyp: the classes are linearly separables

Perceptron – primal formulation Initialize $\mathbf{w}(0) = \mathbf{0}$ Repeat (t) choose example, $(\mathbf{x}(t), y(t))$ if $y(t)\mathbf{w}(t) \cdot \mathbf{x}(t) \leq 0$ then $\mathbf{w}(t+1) = \mathbf{w}(t) + y(t)\mathbf{x}(t)$ until convergence	Decision function- primal $F(\mathbf{x}) = \text{sgn}\left(\sum_{j=0}^n w_j x_j\right),$ $\mathbf{w} = \sum_{i=1}^N \alpha_i y^i \mathbf{x}^i$ α_i : number of times for which the algorithm made a classification error on \mathbf{x}^i
Perceptron – dual formulation Initialize $\alpha = \mathbf{0}, \alpha \in \mathbb{R}^N$ Repeat (t) choose an example, $(\mathbf{x}(t), y(t))$ let $k: \mathbf{x}(t) = \mathbf{x}^k$ if $y(t) \sum_{i=1}^N \alpha_i y^i \mathbf{x}^i \cdot \mathbf{x}(t) \leq 0$ then $\alpha_k = \alpha_k + 1$ until convergence	Decision function - dual $F(\mathbf{x}) = \text{sgn}\left(\sum_{i=1}^N \alpha_i y^i \mathbf{x}^i \cdot \mathbf{x}\right)$ Gram matrix \mathbf{K} : matrix $N \times N$ with term $i, j: \mathbf{K}_{ij} = \mathbf{x}^i \cdot \mathbf{x}^j$ similarity matrix between the training data

Intuition (2) – machine learning algorithms and dot products

Example 1: Perceptron dual formulation for binary classification

- ▶ In the dual formulation of the Perceptron
 - ▶ The decision function writes as $F(x) = \text{sgn}(\sum_{i=1}^N \alpha_i y^i K(x^i, x))$
 - ▶ With the kernel $K(x^i, x) = \langle x^i, x \rangle$, i.e. the kernel is computed directly in the input domain
 - ▶ What if we make use of another similarity function $K(x^i, x)$ instead of the canonical dot product?
 - ▶ The α_i s can be considered as a dual representation of the hyperplane normal vector

Intuition (2) – machine learning algorithms and dot products

Example 2: dual formulation for regression

- ▶ Training examples

- ▶ $D = \{(x^1, y^1), \dots, (x^N, y^N)\}$, we denote $X = \{x^1, \dots, x^N\}$

- ▶ Let us consider a linear model for regression

- ▶ $f(x) = w \cdot x$
 - ▶ Let $x^\perp \in X^\perp$, with X^\perp the orthogonal set of X
 - ▶ $(w + x^\perp) \cdot x^i = w \cdot x^i, \forall x^i \in X$
 - ▶ Adding to w a component outside the space generated by X , has no effect on the linear regression prediction for all the **data in the training set**
 - ▶ If the training criterion only depends on the regression performed on the training data, as is usually the case, it is not needed to consider components of w outside the space generated by X
 - ▶ w can thus be written under the form
 - ▶ $w = \sum_{i=1}^N \alpha_i x^i$
 - ▶ The parameters $\alpha_i, i = 1 \dots N$ are called dual parameters
 - ▶ The regression function can then be directly written under a dual form using dot product:
 - ▶ $f(x) = \sum_{i=1}^N \alpha_i \langle x^i, x \rangle$

Intuition (2) – machine learning algorithms and dot products

Example 2: dual formulation for regression

- ▶ What if we make use of another similarity function $K(x^i, x)$ instead of the canonical dot product?
 - ▶ More generally, let us consider a regression defined through the mapping $\phi(x)$:
 - ▶ $f(x) = \mathbf{w} \cdot \phi(x)$
 - ▶ The solution will be in the space spanned by $\{\phi(x^1), \dots, \phi(x^N)\}$
 - ▶ $\mathbf{w} = \sum_{i=1}^N \alpha_i \phi(x^i)$
 - ▶ $f(x) = \sum_{i=1}^N \alpha_i \langle \phi(x^i), \Phi(\mathbf{x}) \rangle = \sum_{i=1}^N \alpha_i K(x^i, x)$
 - ▶ $K(x^i, x^j) = \langle \phi(x^i), \phi(x^j) \rangle = K_{ij}$
 - ▶ $K = [K_{ij}]$ is the Gram matrix

Intuition – Summary

- ▶ Linear ML methods have a dual representation and can be formulated using dot products in a vector space

- ▶ Examples: adaline, regression, ridge regression, etc
- ▶ The information on the training data is provided by the Gram matrix K :

$$K = (K_{ij})_{i,j=1\dots N} = (K(x^i, x^j))_{i,j=1\dots N}$$

- ▶ With

$$K(x, x') = \langle \Phi(x), \Phi(x') \rangle$$

$$\Phi: X \rightarrow \mathcal{H}$$

$$x \mapsto \Phi(x)$$

- ▶ Such a function $K(.,.)$ defined by a dot product in a feature space will be called a **kernel**
- ▶ For supervised problems, the decision/ regression function $F(x)$ writes as a **linear combination of scalar products**:

$$F(x) = \sum_{i=1}^N \alpha_i K(x^i, x)$$

Kernels

- ▶ After this informal introduction, we will introduce some formal arguments for characterizing kernels that admit a dot product representation in a feature space
- ▶ We first introduce some examples motivating the usefulness of kernels
- ▶ We then address the following question:
 - ▶ What kind of function $K(x, x')$ admits a representation as a dot product in a feature space $K(x, x') = \langle \Phi(x), \Phi(x') \rangle$

Definitions

▶ Gram matrix

- ▶ Given a function $K: X \times X \rightarrow \mathbb{R}$, and a dataset $X = \{x^1, \dots, x^N\}$, the $N \times N$ matrix with element $K_{ij} = K(x^i, x^j)$ is called the Gram matrix of K with respect to X

▶ Positive semi-definite matrix

- ▶ A symmetric matrix \mathbf{K} is positive semi-definite if its eigenvalues are all non negative – or equivalently if $x^T \mathbf{K} x \geq 0 \ \forall x \in X$

Positive definite kernels

- ▶ A **positive definite kernel** on set X , is a function $K: X \times X \rightarrow \mathbb{R}$

- ▶ that is symmetric:

$$K(x, x') = K(x', x)$$

- ▶ Which satisfies, $\forall N \in \mathbb{N}, \forall (x^1, \dots, x^N) \in X^N$ and $\forall (a_1, \dots, a_N) \in \mathbb{R}^N$:

$$\sum_{i=1}^N \sum_{j=1}^N a_i a_j K(x^i, x^j) \geq 0$$

- ▶ Note:

- ▶ this is the general definition of a positive definite function
 - ▶ Positive definiteness allows an easy characterization of kernels

- ▶ Alternative definition with the similarity matrix of a p.d. kernel

- ▶ A kernel K is p.d. if and only if, $\forall N \in \mathbb{N}, \forall (x^1, \dots, x^N) \in X^N$, the similarity matrix $K_{ij} = K(x^i, x^j)$ is **positive semi-definite**
 - ▶ Note: this should be true $\forall N \in \mathbb{N}$

Examples of p.d. kernels

► Linear kernel

- Let $X = \mathbb{R}^n$, the function $K: X^2 \rightarrow \mathbb{R}$:

$$(x, x') \rightarrow K(x, x') = \langle x, x' \rangle_{\mathbb{R}^n}$$

is a p.d. kernel

Proof

- $\langle x, x' \rangle_{\mathbb{R}^n} = \langle x', x \rangle_{\mathbb{R}^n}$
- $\sum_{i=1}^N \sum_{j=1}^N a_i a_j \langle x^i, x^j \rangle_{\mathbb{R}^n} = \left\| \sum_{i=1}^N a_i x^i \right\|_{\mathbb{R}^n}^2 \geq 0$

More general kernels

- ▶ More generally: kernels as dot product in an inner product space

- ▶ Lemma

- ▶ Let X be any set, $\Phi: X \rightarrow \mathbb{R}^n$, the function $K: X^2 \rightarrow \mathbb{R}$:
 $(x, x') \rightarrow K(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathbb{R}^n}$

is a p.d. kernel

Proof: same as above

$$\begin{aligned} \langle \Phi(x), \Phi(x') \rangle_{\mathbb{R}^n} &= \langle \Phi(x'), \Phi(x) \rangle_{\mathbb{R}^n} \\ \sum_{i=1}^N \sum_{j=1}^N a_i a_j \langle \Phi(x^i), \Phi(x^j) \rangle_{\mathbb{R}^n} &= \left\| \sum_{i=1}^N a_i \Phi(x^i) \right\|_{\mathbb{R}^n}^2 \geq 0 \end{aligned}$$

More general kernels

Example: Polynomial Kernel

- ▶ Consider a 2 dimensional input space $X \subset \mathbb{R}^2$ and
- ▶ $\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3, \Phi(x) = \Phi(x_1, x_2) = (x_1^2, x_2^2, \sqrt{2}x_1 x_2)$
 - $K(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathbb{R}^3}$
 - $K(x, x') = x_1^2 x_1'^2 + 2x_1 x_2 x_1' x_2' + x_2^2 x_2'^2$
 - $K(x, x') = \langle x, x' \rangle_{\mathbb{R}^2}^2$
- ▶ Note:
 - ▶ $K(x, x') = \langle \Phi(x), \Phi(x') \rangle$ can be computed directly as $\langle x, x' \rangle_{\mathbb{R}^2}^2$ without explicitly evaluating their coordinate in the feature space
 - ▶ Cheaper to compute in the original space than in the feature space
 - ▶ The same kernel is obtained with $\Phi(x_1, x_2) = (x_1^2, x_2^2, x_1 x_2, x_2 x_1)$ and a dot product in \mathbb{R}^4
 - ▶ Shows that the feature space is not uniquely determined by the kernel function

Example: Polynomial Kernel

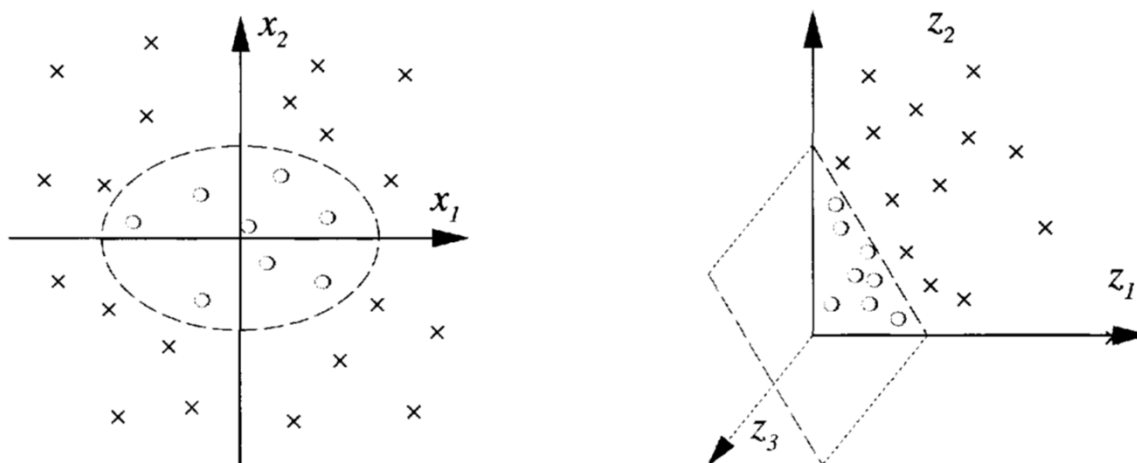


Figure 2.1 Toy example of a binary classification problem mapped into feature space. We assume that the true decision boundary is an ellipse in input space (left panel). The task of the learning process is to estimate this boundary based on empirical data consisting of training points in both classes (crosses and circles, respectively). When mapped into feature space via the nonlinear map $\Phi_2(x) = (z_1, z_2, z_3) = ([x]_1^2, [x]_2^2, \sqrt{2} [x]_1[x]_2)$ (right panel), the ellipse becomes a hyperplane (in the present simple case, it is parallel to the z_3 axis, hence all points are plotted in the (z_1, z_2) plane). This is due to the fact that ellipses can be written as linear equations in the entries of (z_1, z_2, z_3) . Therefore, in feature space, the problem reduces to that of estimating a hyperplane from the mapped data points. Note that via the polynomial kernel (see (2.12) and (2.13)), the dot product in the three-dimensional space can be computed without computing Φ_2 . Later in the book, we shall describe algorithms for constructing hyperplanes which are based on dot products (Chapter 7).

Scholkopf et al.
2002

Characterization of kernels

- ▶ Up to now kernels have been characterized by explicitly defining a mapping in a feature space and then computing an inner product in this space
- ▶ We will introduce an alternative characterization of a kernel
 - ▶ It is one of the main theoretical tools to characterize kernels
 - ▶ Without explicitly defining the feature space (i.e. Φ)

Characterization of kernels

Definitions and properties

► Inner product

- Let \mathcal{H} a vector space over \mathbb{R} , a function $\langle ., . \rangle_{\mathcal{H}}$ is said to be an inner product on \mathcal{H} if
 - $\langle \alpha_1 f_1 + \alpha_2 f_2, g \rangle_{\mathcal{H}} = \alpha_1 \langle f_1, g \rangle_{\mathcal{H}} + \alpha_2 \langle f_2, g \rangle_{\mathcal{H}}$ linear (bilinear)
 - $\langle f, g \rangle_{\mathcal{H}} = \langle g, f \rangle_{\mathcal{H}}$ symmetric
 - $\langle f, f \rangle_{\mathcal{H}} \geq 0$ and $\langle f, f \rangle_{\mathcal{H}} = 0$ iff $f = 0$
 - We can then define a norm on \mathcal{H} as $\|f\|_{\mathcal{H}} = \sqrt{\langle f, f \rangle_{\mathcal{H}}}$
- \mathcal{H} endowed with an inner product is an inner product space

► Hilbert space

- Is an inner product space \mathcal{H} with the additional properties that it is separable and complete i.e. any Cauchy sequence in \mathcal{H} converges in \mathcal{H}
 - A Cauchy sequence (f_n) is a sequence whose elements become progressively arbitray close to each other

$$\lim_{m>n, n \rightarrow \infty} \|f_n - f_m\|_{\mathcal{H}} = 0$$

- \mathcal{H} is separable if for any $\epsilon > 0$ there exists a finite set of elements of \mathcal{H} , $\{f_1, \dots, f_N\}$ such that for all $f \in \mathcal{H}$,

$$\min_i \|f_i - f\|_{\mathcal{H}} < \epsilon$$

Characterization of kernels

Definitions and properties

- ▶ Cauchy-Schwartz inequality for dot products

- ▶ In an inner product space

- ▶ $\langle x, x' \rangle^2 \leq \|x\|^2 \|x'\|^2$

- ▶ Cauchy-Schwartz inequality for kernels

- ▶ If K is a p.d. kernel and $x_1, x_2 \in X$, then:

- $|K(x^1, x^2)|^2 \leq K(x^1, x^1) \cdot K(x^2, x^2)$

Characterization of kernels

- ▶ Theorem

- ▶ $K: X \times X \rightarrow \mathbb{R}$ is a p.d. kernel on X if and only if there exists a Hilbert space \mathcal{H} and a mapping $\Phi: X \rightarrow \mathcal{H}$ such that:

$$\forall x, x' \in X, K(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}}$$

- ▶ Central result that establish a link between kernels defined as dot products in a feature vector space and positive definite functions
- ▶ In order to demonstrate this result, we explicitly construct the feature -Hilbert- space

Characterization of kernels

- ▶ Assumption: K is a p.d. kernel
- ▶ Objective: construct an appropriate Hilbert space and a mapping Φ
- ▶ Defining the mapping Φ

- ▶ Let us define $\Phi: X \rightarrow \mathbb{R}^X$, where $\mathbb{R}^X := \{f: X \rightarrow \mathbb{R}\}$ is the space of functions mapping X into \mathbb{R} as:

$$\Phi: X \rightarrow \mathbb{R}^X$$

$$x \mapsto K(\cdot, x)$$

$\Phi(x)$ denotes a function that assigns a value $K(x', x)$ to $x' \in X$, i.e. $\Phi(x)(\cdot) = K(\cdot, x)$

To each point x in the X space, one associates a function $\Phi(x) = K(\cdot, x)$

This function will be a point in a vector space

See Fig. next slide

Characterization of kernels

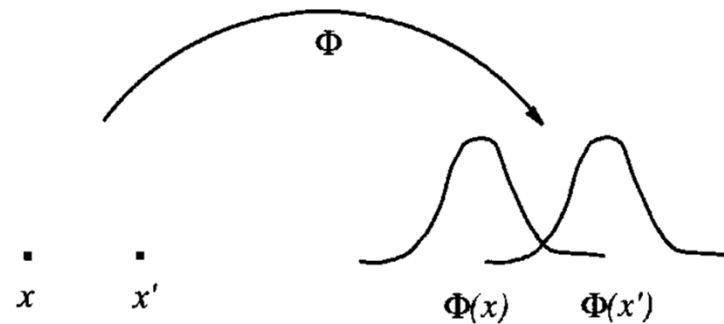


Figure 2.2 One instantiation of the feature map associated with a kernel is the map (2.21), which represents each pattern (in the picture, x or x') by a kernel-shaped *function* sitting on the pattern. In this sense, each pattern is represented by its similarity to *all* other patterns. In the picture, the kernel is assumed to be bell-shaped, e.g., a Gaussian $k(x, x') = \exp(-\|x - x'\|^2 / (2 \sigma^2))$. In the text, we describe the construction of a dot product $\langle \cdot, \cdot \rangle$ on the function space such that $k(x, x') = \langle \Phi(x), \Phi(x') \rangle$.

Fig, Scholkopf
et al. 2002

Characterization of kernels

- ▶ Construction of the feature space
- ▶ Let us consider the space of functions
- ▶ $\mathcal{H} = \left\{ \sum_{i=1}^m \alpha_i K(\cdot, x^i) : m \in \mathbb{N}, x_i \in X, \alpha_i \in \mathbb{R}, i = 1 \dots m \right\}$
 - ▶ Note: here $m \in \mathbb{N}, x^i \in X, \alpha_i \in \mathbb{R}$ are arbitrary,
 - ▶ \mathcal{H} is closed under multiplication by a scalar and addition of functions and is then a vector space
 - ▶ We define the dot product onto \mathcal{H} :
 - ▶ Let $f(\cdot) = \sum_{i=1}^l \alpha_i K(\cdot, x^i)$ $g(\cdot) = \sum_{j=1}^m \beta_j K(\cdot, x'^j)$
 - ▶ $\langle f, g \rangle = \sum_{i=1}^l \sum_{j=1}^m \alpha_i \beta_j K(x^i, x'^j) = \sum_{i=1}^l \alpha_i g(x^i) = \sum_{j=1}^m \beta_j f(x'^j)$
 - ▶ From these equalities, $\langle \cdot, \cdot \rangle$ is symmetric, bilinear
 - ▶ Since K is p.d. for any $f(\cdot) = \sum_{i=1}^l \alpha_i K(\cdot, x^i)$, one has:
$$\langle f, f \rangle = \sum_{i,j=1}^l \alpha_i \alpha_j K(x^i, x^j) \geq 0$$
 - ▶ Note: this means that $\langle \cdot, \cdot \rangle$ is itself a p.d. kernel on the space of functions

Characterization of kernels

- ▶ Reproducing property of the kernel
 - ▶ $\langle f, K(\cdot, x) \rangle = \sum_{i=1}^l \alpha_i K(x, x^i) = f(x)$
 - ▶ Particular case: $\langle K(\cdot, x), K(\cdot, x') \rangle = K(x, x')$ or $\langle \Phi(x), \Phi(x') \rangle = K(x, x')$
 - ▶ Using the reproducing property and Cauchy Schwartz:
 - ▶ $|f(x)|^2 = |\langle f, K(\cdot, x) \rangle|^2 \leq K(x, x) \cdot \langle f, f \rangle$
 - ▶ Then $\langle f, f \rangle = 0$ implies $f = 0$
 - ▶ This establishes that $\langle \cdot, \cdot \rangle$ is a dot product
 - ▶ It remains to show that space \mathcal{H} is also complete and separable
 - ▶ See e.g. (Shawe Taylor et al. 2004)
- ▶ Summary
 - ▶ Given a p.d. kernel K , one has built \mathcal{H} an associated Hilbert space in which the reproducing property holds, and a mapping Φ
 - ▶ \mathcal{H} is called the Reproducing Kernel Hilbert Space (RKHS) of K
 - ▶ We give the formal definition of a RKHS later

Characterization of kernels

- ▶ Conversely
- ▶ Given a mapping Φ from X to a dot product space, we can get a p.d. kernel via $k(x, x') = \langle \Phi(x), \Phi(x') \rangle$
- ▶ Proof
 - ▶ $\forall \alpha_i \in \mathbb{R}, x^i \in X, i = 1 \dots m$, we have
 - ▶ $\sum_{i,j} \alpha_i \alpha_j K(x^i, x^j) = \langle \sum_i \alpha_i \Phi(x^i), \sum_j \alpha_j \Phi(x^j) \rangle = \|\sum_i \alpha_i \Phi(x^i)\|^2 \geq 0$

Characterization of kernels

Summary

- ▶ This characterization allows us
 - ▶ to give an equivalent definition of p.d. kernels as functions with the property that there exists a map Φ into a dot product space such that $k(x, x') = \langle \Phi(x), \Phi(x') \rangle$ holds
 - ▶ To construct kernels from feature maps
 - ▶ $k(x, x') = \langle \Phi(x), \Phi(x') \rangle$
- ▶ is at the base of the kernel trick

Kernel Trick

- ▶ Given an algorithm which is formulated in terms of a p.d. kernel, K , one can construct an alternative algorithm by replacing K by another p.d. kernel K'
- ▶ Intuition
 - ▶ The original algorithm is a dot product based algorithm on vectors $\Phi(x^1), \dots, \Phi(x^m)$, when K is replaced by K' , the algorithm is the same but operates on $\Phi'(x^1), \dots, \Phi'(x^m)$
 - ▶ The best known application of the trick is when K is the dot product in the input domain. It can be replaced by another kernel, e.g. non linear. Most of the linear data analysis algorithms (PCA, ridge regression, etc) can then be automatically « kernalized ».
 - ▶ Any algorithm that process finite dimensional vectors that is expressed in terms of pairwise inner products, can be applied to infinite-dimensional vectors in the feature space of a p.d. kernel, by replacing each inner product by a kernel evaluation

Reproducing Kernel Hilbert Spaces - RKHS

- ▶ Let X be a non empty set and \mathcal{H} a Hilbert space of functions with inner product $\langle ., . \rangle$. Then \mathcal{H} is called a RKHS if there exists a function $K: X \times X \rightarrow \mathbb{R}$ with the following properties:
 - ▶ K has the reproducing property
 - ▶ $\langle f, k(x, .) \rangle = f(x) \forall f \in \mathcal{H}$
 - ▶ In particular
 - ▶ $\langle k(x, .), k(x', .) \rangle = k(x, x')$
 - ▶ $\forall x \in X, K(x, .) \in \mathcal{H}$
- ▶ K is called a reproducing kernel
- ▶ Property
 - ▶ The RKHS determines uniquely K and reciprocally
 - ▶ A function $K: X \times X \rightarrow \mathbb{R}$ is positive definite iff it is a reproducing kernel!

RKHS example – The linear kernel

- ▶ Let $X = \mathbb{R}^n$ and consider the linear kernel

- ▶ $K(x, x') = \langle x, x' \rangle_{\mathbb{R}^n}$

- ▶ The RKHS of the linear kernel is the set of linear functions

$$\mathcal{H} = \{f_w(x) = \langle w, x \rangle_{\mathbb{R}^n} : w \in \mathbb{R}^n\}$$

- ▶ Inner product is defined as

$$\forall v, w \in \mathbb{R}^n, \langle f_v, f_w \rangle_{\mathcal{H}} = \langle v, w \rangle_{\mathbb{R}^n}$$

- ▶ The corresponding norm is

$$\forall w \in \mathbb{R}^n, \|f\|_{\mathcal{H}} = \|w\|_{\mathbb{R}^n}$$

Infinite dimensional feature space

► Lemma

- Let $D = \{x^1, \dots, x^N\}$ distinct points in X , and $\sigma \neq 0$. The matrix K given by $K_{ij} := \exp(-\frac{\|x^i - x^j\|^2}{2\sigma^2})$ has full rank.
- This means that the points $\Phi(x^1), \dots, \Phi(x^N)$ are linearly independent (since $K = \Phi^T \Phi$ with Φ the matrix with column vectors the $\Phi(x^i)$).
- Then they span an N dimensional subspace of \mathcal{H} .
- Since this is true for all N , i.e. no restriction on the number of training examples, the feature space is then of **infinite dimension**

How to build new kernels

- ▶ Kernels can be built from combinations of known ones
- ▶ Let K_1, K_2 be kernels defined on a metric space X^2 , K_3 defined on the Hilbert space \mathcal{H} , the following combinations are kernels:
 - ▶ $K(x, z) = K_1(x, z) + K_2(x, z)$
 - ▶ $K(x, z) = K_1(x, z) \cdot K_2(x, z)$
 - ▶ $K(x, z) = aK_1(x, z)$
 - ▶ $K(x, z) = K_3(\phi(x), \phi(z))$
 - ▶

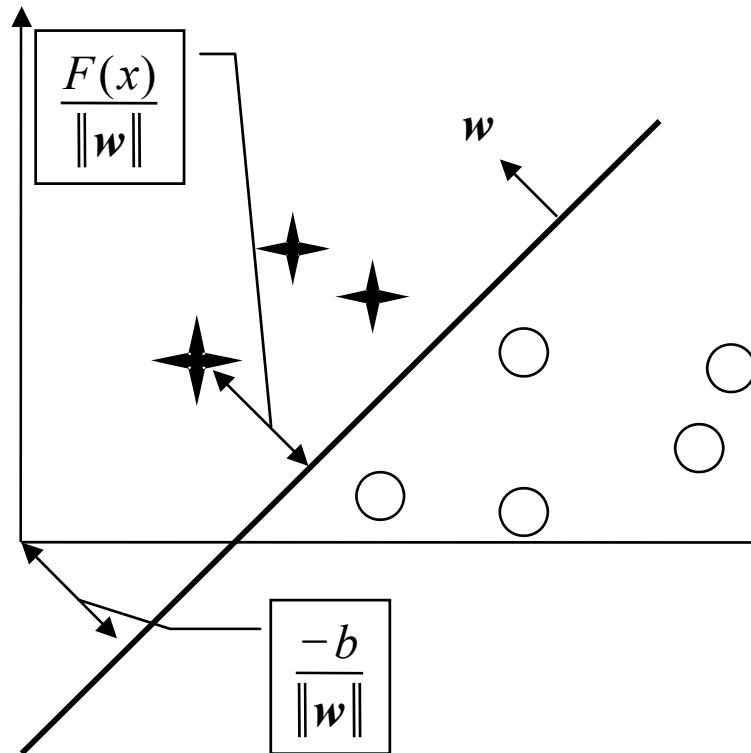
Support vector machines

Support vector machines - SVMs

- ▶ SVMs have been developed initially for classification and then extended to regression and density estimation
- ▶ The course is a brief introduction to the algorithmic concepts behind SVMs
- ▶ We will consider 2 classes classification
 - ▶ And develop the concept in the case of linear machines
 - ▶ The development of SVMs in the mid 90s mainly concerned non linear embeddings and dealing with structured data (strings, graphs)
 - ▶ However the main algorithmic concept could be understood in the simplified linear setting
 - ▶ As for any kernel method, non linear SVMs are not well adapted to large dimensional spaces and large datasets
 - ▶ They remain an important concept

Margin

► Margin definition



► Given dataset

$$D = \{(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^N, y^N)\}$$

with $y^i \in \{-1, 1\}$ and hyperplane H

$$F(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b = 0$$

► The geometric margin for \mathbf{x}^i is:

$$\triangleright M(\mathbf{x}^i) = y^i \left(\frac{\mathbf{w} \cdot \mathbf{x}^i}{\|\mathbf{w}\|} + \frac{b}{\|\mathbf{w}\|} \right) = y^i \frac{F(\mathbf{x}^i)}{\|\mathbf{w}\|}$$

► Margin of \mathbf{w} w.r.t. dataset D

$$\text{Min}_{\mathbf{x}^i \in D} M(\mathbf{x}^i)$$

► Maximal margin hyperplane

$$\text{Max}_{\mathbf{w}} (\text{Min}_{\mathbf{x}^i \in D} M(\mathbf{x}^i))$$

Geometric margin vs functional margin

- ▶ **Geometric margin**

- ▶ $y^i \frac{F(x^i)}{\|w\|}$

- ▶ **Functional margin**

- ▶ $y^i F(x^i)$

- ▶ Replacing w with kw $k \in \mathbb{R}$ does not change the decision function or the geometric margin, but changes the functional margin
- ▶ For SVMs, one will set the functional margin to 1 (arbitrary value) and one will optimize the geometric margin

Linear separation with optimal hyperplane (1974)

- ▶ Hyp : the training set D is linearly separable
 - ▶ $D = \{(x^1, y^1), \dots, (x^N, y^N)\}$ with $y^i \in \{-1, 1\}$
- ▶ We consider linear decision functions : $F(x) = w \cdot x + b$
- ▶ Learning problem: what is the « optimal hyperplane » H^* :
 - ▶ That separates D , i.e. $y^i F(x^i) \geq 1, \forall i$
 - ▶ With a maximal geometric margin $M = \min_{(x^i, y^i) \in D} \frac{y^i F(x^i)}{\|w\|} = \frac{1}{\|w\|}$
- ▶ This is the **primal formulation** of the linear SVM:
 - ▶ Minimize $\|w\|^2$
 - ▶ Under the constraint $y^i F(x^i) \geq 1 \quad \forall i = 1, \dots, N$

-
- ▶ For linear kernels, one usually solves the primal problem
 - ▶ There exist several fast gradient algorithms
 - ▶ For non linear kernels, one usually solves a dual formulation of the problem
 - ▶ In the following, one introduces some basic notions on constraint optimization in order to describe this dual formulation

Interlude

Optimisation
under equality and inequality constraints

Optimisation under equality and inequality constraints

- ▶ Let
 - ▶ $f, g_{i,i=1,\dots,k}, h_{j,j=1,\dots,m}$ be functions defined on \mathbb{R}^n taking their value in \mathbb{R}
- ▶ Let us consider the following optimization problem (pb. (0)) :

$$\text{Min } f(\mathbf{w}), \mathbf{w} \in \Omega \subset \mathbb{R}^n$$

Under constraints

$$g_i(\mathbf{w}) \leq 0, i = 1, \dots, k$$

$$\text{denoted } \mathbf{g}(\mathbf{w}) \leq 0$$

$$h_j(\mathbf{w}) = 0, j = 1, \dots, m$$

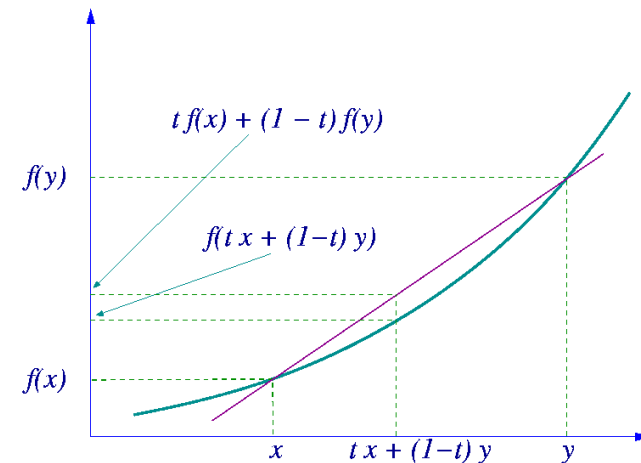
$$\text{denoted } \mathbf{h}(\mathbf{w}) = 0$$

▶ Definitions

- ▶ **Objective function** $f(\mathbf{w})$
- ▶ **Admissible region** $R = \{\mathbf{w} \in \Omega: \mathbf{g}(\mathbf{w}) \leq 0, \mathbf{h}(\mathbf{w}) = 0\}$, region of Ω where f is defined and the constraints verified
- ▶ \mathbf{w}^* is a **global minimum** if there does not exist another point \mathbf{w} such that $f(\mathbf{w}) < f(\mathbf{w}^*)$, it is a **local optimum** if $\exists \epsilon > 0: f(\mathbf{w}) \geq f(\mathbf{w}^*)$, on the subspace $\|\mathbf{w} - \mathbf{w}^*\| < \epsilon$
- ▶ A constraint $g_i(\mathbf{w}) \leq 0$ is said **active** if the solution \mathbf{w}^* verifies $\mathbf{g}(\mathbf{w}^*) = 0$ and **inactive** otherwise
- ▶ The optimal value of the objective function ($f(\mathbf{w})$ solution of pb. (0)) is called the **value of the primal optimization problem**.

Optimization – convex functions

- ▶ $f(w)$ is **convex** for $w \in \mathbb{R}^n$ if
 - ▶ $\forall t \in [0,1], \forall w, v \in \mathbb{R}^n, \forall t \in [0,1], f(tw + (1-t)v) \leq tf(w) + (1-t)f(v)$



- ▶ A set $\Omega \subset \mathbb{R}^n$ is **convex** if $\forall w, v \in \mathbb{R}^n, \forall t \in [0,1], tw + (1-t)v \in \Omega$
- ▶ If a function is convex, any local minimum is a global minimum
- ▶ An optimization problem where Ω is convex, the objective function is convex and the constraints are convex is said convex

Unconstrained optimisation

- ▶ Theorem Fermat

- ▶ A necessary and sufficient condition w^* to be a minimum of $f(w)$, $f \in C^1$ is $\frac{\partial f(w^*)}{\partial w} = 0$
- ▶ If f is convex, the condition is sufficient

Optimization with equality constraints

Lagrangian

- ▶ Optimization with equality constraints (pb (I)):

$$\text{Min } f(\mathbf{w}), \mathbf{w} \in \Omega \subset \mathbb{R}^n$$

Under constraints

$$h_j(\mathbf{w}) = 0, \quad j = 1, \dots, m \quad \text{denoted } \mathbf{h}(\mathbf{w}) = 0$$

- ▶ We define the **Lagrangian** $L(\mathbf{w}, \boldsymbol{\beta})$ associated to this problem as:

$$L(\mathbf{w}, \boldsymbol{\beta}) = f(\mathbf{w}) + \sum_{j=1}^m \beta_j h_j(\mathbf{w})$$

- ▶ the β_j are the Lagrange coefficients
- ▶ Note
 - ▶ If \mathbf{w}^* is a solution to pb. (I), it may happen that $\frac{\partial f(\mathbf{w}^*)}{\partial \mathbf{w}} \neq 0$

Optimization with equality constraints

Th. Lagrange

► Theorem Lagrange

- A necessary condition for \mathbf{w}^* , to be a solution of (pb. (I)), with $f, h_i \in C^1$ is

- $\frac{\partial L(\mathbf{w}^*, \boldsymbol{\beta}^*)}{\partial \mathbf{w}} = 0$

- $\frac{\partial L(\mathbf{w}^*, \boldsymbol{\beta}^*)}{\partial \boldsymbol{\beta}} = 0$

- If $L(\mathbf{w}, \boldsymbol{\beta}^*)$ is a convex function of \mathbf{w} , the condition is sufficient

► Note

- The first equation gives a new system of equations
 - The second one gives constraints

Optimization under equality and inequality constraints – Generalized Lagrangian

- ▶ The generalized Lagrangian for pb. (0) is:

$$L(\mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = f(\mathbf{w}) + \sum_{i=1}^k \alpha_i g_i(\mathbf{w}) + \sum_{j=1}^m \beta_j h_j(\mathbf{w})$$

- ▶ Remember the primal formulation of pb. (0):

$\text{Min } f(\mathbf{w}), \mathbf{w} \in \Omega \subset \mathbb{R}^n$

under constraints

$g_i(\mathbf{w}) \leq 0, i = 1, \dots, k$

$h_j(\mathbf{w}) = 0, j = 1, \dots, m$

denoted $\mathbf{g}(\mathbf{w}) \leq 0$

denoted $\mathbf{h}(\mathbf{w}) = 0$

Optimization under equality and inequality constraints – Lagrange Duality

► Reformulation of the primal problem

- Let $\theta_P(\mathbf{w}) = \max_{\alpha, \beta; \alpha \geq 0} L(\mathbf{w}, \alpha, \beta) = f(\mathbf{w}) + \max_{\alpha, \beta; \alpha \geq 0} (\alpha \cdot g(\mathbf{w}) + \beta \cdot h(\mathbf{w}))$,
 - P in θ_P holds for Primal
 - Note: $\theta_P(\mathbf{w})$ is a function of \mathbf{w} , not of α, β
- $\theta_P(\mathbf{w}) = \begin{cases} f(\mathbf{w}) & \text{if } \mathbf{w} \text{ is admissible} \\ +\infty & \text{otherwise, i. e. } \mathbf{w} \text{ violates any of the constraints} \end{cases}$
 - This is because $\max_{\alpha, \beta; \alpha \geq 0} (\alpha g(\mathbf{w}) + \beta h(\mathbf{w})) = 0$ for a point \mathbf{w} admissible i.e. that satisfies the primal constraints
 - $\theta_P(\mathbf{w})$ takes the same value as the objective function $f(\mathbf{w})$ for any \mathbf{w} admissible
- The original primal problem (pb 0) can be reformulated as:
 - $\min_{\mathbf{w} \in \Omega} \theta_P(\mathbf{w})$ or $\min_{\mathbf{w} \in \Omega} \max_{\alpha, \beta; \alpha \geq 0} L(\mathbf{w}, \alpha, \beta)$
 - $\min_{\mathbf{w} \in \Omega} \theta_P(\mathbf{w})$ is called the **value of the primal**

Optimization – Lagrange duality

► Dual formulation of the optimization problem

- Let $\theta_D(\alpha, \beta) = \min_{w \in \Omega} L(w, \alpha, \beta)$
 - D in θ_D holds for dual
- The dual optimization problem for (pb 0) is:
 - $\max_{\alpha, \beta, \alpha_i \geq 0} \theta_D(\alpha, \beta) = \min_{w \in \Omega} L(w, \alpha, \beta)$
under constraint $\alpha \geq 0$
- $\max \theta_D(\alpha, \beta)$ is called the **value of the dual**
- Note : $\min_{w \in \Omega} L(w, \alpha, \beta)$ is a function of α, β only

► Property (weak duality)

- **the value of the dual is bounded above by the value of the primal**
- $\max_{\alpha, \beta; \alpha \geq 0} \min_{w \in \Omega} L(w, \alpha, \beta) \leq \min_{w \in \Omega} \max_{\alpha, \beta; \alpha \geq 0} L(w, \alpha, \beta)$

For some cases there is equality, this is called strong duality

Optimisation : weak duality

- ▶ $\max_{\alpha, \beta; \alpha \geq 0} \min_{v \in \Omega} L(v, \alpha, \beta) \leq \min_{w \in \Omega} \max_{\alpha, \beta; \alpha \geq 0} L(w, \alpha, \beta):$
 - ▶ $\theta_D(\alpha, \beta) = \min_{v \in \Omega} L(v, \alpha, \beta):$
 - ▶ $\theta_D(\alpha, \beta) \leq L(w, \alpha, \beta)$ for any w admissible
 - ▶ $L(w, \alpha, \beta) = f(w) + \alpha \cdot g(w) + \beta \cdot h(w)$ with $\alpha \geq 0, g(w) \leq 0, h(w) = 0$
 - ▶ $L(w, \alpha, \beta) \leq f(w)$
 - ▶ $\theta_D(\alpha, \beta) \leq f(w)$
 - ▶ $\max_{\alpha, \beta; \alpha \geq 0} \min_{v \in \Omega} L(v, \alpha, \beta) \leq f(w)$ for any w admissible
 - ▶ $\theta_P(w) = \max_{\alpha, \beta; \alpha \geq 0} L(w, \alpha, \beta)$
 - ▶ $\theta_P(w) = f(w) + \max_{\alpha, \beta; \alpha \geq 0} (\alpha \cdot g(w) + \beta \cdot h(w))$
 - ▶ $\theta_P(w) = \begin{cases} f(w) & \text{if } w \text{ is admissible} \\ +\infty & \text{otherwise} \end{cases}$ since $\max_{\alpha, \beta; \alpha \geq 0} (\alpha g(w) + \beta h(w)) = 0$ for an admissible point
 - ▶ $\theta_P(w) = f(w)$ for w admissible
- ▶ Hence the inequality of weak duality

► Theorem: strong duality

► For an optimization problem

$\text{Min } f(\mathbf{w}), \mathbf{w} \in \Omega \subset \mathbb{R}^n$ convex and $f \in C^1$ convex

under constraints

$$g_i(\mathbf{w}) \leq 0, i = 1, \dots, k \quad \text{denoted } \mathbf{g}(\mathbf{w}) \leq 0$$

$$h_j(\mathbf{w}) = 0, j = 1, \dots, m \quad \text{denoted } \mathbf{h}(\mathbf{w}) = 0$$

where the g_i and the h_j are affines ($h_j(\mathbf{w}) = A_j \mathbf{w} + b_j$) (hence convex)

- The values of the primal and of the dual are equal
- The conditions for the existence of an optimum are given by the theorem of Kuhn and Tucker

Optimization

Theorem Kuhn and Tucker

- ▶ Let us consider (pb. (0)) with Ω convex and $f \in C_1$ convex, g_i, h_j affines ($h = A \cdot w + b$)
- ▶ A necessary and sufficient condition for w^* to be an optimum is that there exist α^* and β^* :

$$\left\{ \begin{array}{l} \frac{\partial L(w^*, \alpha^*, \beta^*)}{\partial w} = 0 \\ \frac{\partial L(w^*, \alpha^*, \beta^*)}{\partial \beta} = 0 \\ \alpha_i^* g_i(w^*) = 0, i = 1..k \\ g_i(w^*) \leq 0, i = 1..k \\ \alpha_i^* \geq 0, i = 1..k \end{array} \right.$$

- ▶ Under the assumptions of **convexity**, the dual formulation is an alternative to the primal formulation which may be simpler to handle (e.g. non linear SVMs)

Optimization

► Note

- The third condition ($\alpha_i^* g_i(\mathbf{w}^*) = 0, i = 1..k$), called complementary condition of Karush-Kuhn-Tucker implies that for an active constraint $\alpha_i^* \geq 0$ when for an inactive constraint $\alpha_i^* = 0$
 - Either a constraint is **active** ($\alpha_i^* \geq 0$ and $g_i(\mathbf{w}^*) = 0$), \mathbf{w}^* is a frontier point of the admissible region
 - Or it is **inactive** ($\alpha_i^* = 0$) and \mathbf{w}^* is in the admissible region
- If the solution point \mathbf{w}^* is in the admissible region (inactive constraint) then the conditions of optimality are given by the Fermat theorem and $\alpha_i^* = 0$. If it is on the frontier (active constraint), the conditions of optimality are given by the theorem of Lagrange with $\alpha_i^* > 0$.

► End of the interlude

SVM – primal and dual formulations

For a linear kernel

- ▶ SVM

- ▶ Ω, f and the constraints are convex, L is quadratic
- ▶ One considers the case, $D = \{(x^i, y^i)\}_{i=1\dots N}$ linearly separable, $y^i \in \{-1, 1\}$

- ▶ Primal pb

$$\begin{aligned} & \text{Min}(w \cdot w) && \text{(i.e. max the margin)} \\ & \text{under constraints} \\ & y^i(w \cdot x^i + b) \geq 1, i = 1..N \end{aligned}$$

- ▶ Lagrangian primal

$$L(w, b, \alpha) = \frac{1}{2} w \cdot w - \sum_{i=1}^N \alpha_i (y^i (w \cdot x^i + b) - 1)$$

- ▶ Lagrangian dual

$$L(w, b, \alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y^i y^j \alpha_i \alpha_j (x^i \cdot x^j)$$

- ▶ With $\alpha_i \geq 0$ in both cases

SVM – primal and dual formulations

For a linear kernel

► Derivation of the dual formulation

- Let us start from the primal formulation of the Lagrangian

- $L(w, b, \alpha) = \frac{1}{2} w \cdot w - \sum_{i=1}^N \alpha_i (y^i (w \cdot x^i + b) - 1)$

- Let us first minimize $L(w, b, \alpha)$ w.r.t. w, b in order to get θ_D

- $\frac{\partial L(w, b, \alpha)}{\partial w} = w - \sum_{i=1}^N \alpha_i y^i x^i = 0$

- $w = \sum_{i=1}^N \alpha_i y^i x^i$

- $\frac{\partial L(w, b, \alpha)}{\partial b} = \sum_{i=1}^N \alpha_i y^i = 0$

- Let us plug back $w = \sum_{i=1}^N \alpha_i y^i x^i$ in the primal form of $L(w, b, \alpha)$

- We get

- $L(w, b, \alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N y^i y^j \alpha_i \alpha_j (x^i \cdot x^j) - b \sum_{i=1}^N \alpha_i y^i$

- $L(w, b, \alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N y^i y^j \alpha_i \alpha_j (x^i \cdot x^j)$

- This expression of the dual form only depends on the variables α, β

SVM – primal and dual formulations

For a linear kernel

- ▶ Dual problem

- ▶ $\text{Max}_{\alpha} L(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N d^i d^j \alpha_i \alpha_j (x^i \cdot x^j)$

under constraints:

$$\begin{cases} \sum_{i=1}^N d^i \alpha_i = 0 \\ \alpha_i \geq 0, i = 1..N \end{cases}$$

- ▶ Constraint $\alpha_i \geq 0, i = 1..N$ has always been present
- ▶ $\sum_{i=1}^N d^i \alpha_i = 0$ comes from the derivation of the dual form

- ▶ This is a quadratic optimization problem under constraints

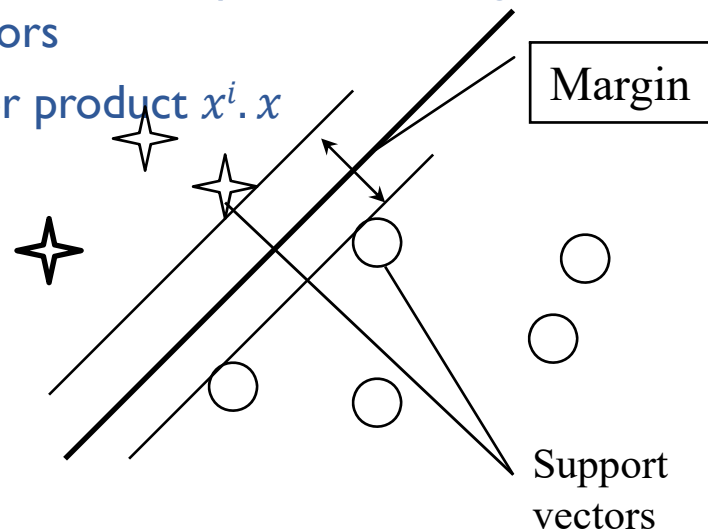
SVM – solution

For a linear kernel

- ▶ Solution : \mathbf{w}^* only depends on the support vectors, i.e. the vectors on the margin that verify: $d^i F^*(x^i) = 1$
- ▶ The decision function takes the form:

$$F(x, \alpha^*, \beta^*) = \sum_{i \text{ support vector}} y^i \alpha_i^* (x^i \cdot x) + b^*$$

- ▶ Note: Whatever the dimension of the space is, the degree of freedom is equal to the number of support vectors
- ▶ F^* only depends on the inner product $x^i \cdot x$



SVM – primal and dual formulations

For a linear kernel

- ▶ Let α, b be solution of the dual problem
- ▶ Then the decision function writes
 - ▶ $F(x, \alpha^*, b^*) = \sum_{i=1}^N \alpha_i^* y^i x^i + b^*$
 - ▶ Remember the KKT condition $\alpha_i^* g_i(\mathbf{w}^*) = 0, i = 1..k$:
 - ▶ $\alpha_i^* (y^i (w^* \cdot x^i + b^*) - 1) = 0, i = 1..k$ (in the primal problem definition)
 - ▶ This means that only the points x^i for which $y^i (w^* \cdot x^i + b^*) = 1$ have a coefficient $\alpha_i^* \neq 0$
 - ▶ i.e. $F(x, \alpha^*, b^*) = \sum_{i \text{ support vector}} \alpha_i^* y^i x^i + b^*$

Support vector machine

Non linear kernels

- ▶ The dot product in the input space X is replaced by a dot product in the feature space \mathcal{H}
- ▶ This dot product can be computed by a kernel function
- ▶ X can be any set – not only a inner product space
- ▶ The above formulation can be transposed by replacing $\langle x^i, x \rangle$ with $K(x^i, x)$

$$w = \sum_{x^i \text{ support vector}} y^i \alpha_i \Phi(x^i) \quad \Phi: R^n \rightarrow \mathcal{H} \quad F(x) = \sum_{x^i \text{ support vector}} y^i \alpha_i \Phi(x^i) \cdot \Phi(x) + b$$

$$\Phi(x) \cdot \Phi(x') = K(x, x')$$
$$F(x) = \sum_{x^i \text{ support vector}} y^i \alpha_i K(x, x^i) + b$$

Support vector machine

Non linear kernels

- ▶ In practice, this formulation does not lead to stable solutions, and other formulations that weaken the constraints (soft margins) are used.

Support vector machine

Optimization algorithms – rebirth of stochastic algorithms

- ▶ Standard algorithm for SVMs
 - e.g. Sequential Minimal Optimization (SMO)
- ▶ Versus stochastic algorithms – (Bottou 2007)
 - Task : Document classification - RCV1 documents belonging to the class CCAT (2 classes classification task)
 - Programs [SVMLight](#) and [SVMPerf](#) are well known SVM solvers written by [Thorsten Joachims](#). SVMLight is suitable for SVMs with arbitrary kernels. Similar results could be achieved using [Chih-Jen Lin](#)'s [LibSVM](#) software. SVMPerf is a specialized solver for linear SVMs. It is considered to be one of the most efficient optimizer for this particular problem.

Algorithm (hinge loss)	Training Time	Primal cost	Test Error
SVMLight	23642 secs	0.2275	6.02%
SVMPerf	66 secs	0.2278	6.03%
Stochastic Gradient (svmsgd)	1.4 secs	0.2275	6.02%
Stochastic Gradient (svmsgd2)	1.4 secs	0.2275	6.01%

Gaussian process regression

Motivations

- ▶ Most ML algorithm for regression predict a mean value
- ▶ Gaussian processes are Bayesian methods that allow us to predict, not only a mean value, but a distribution over the output values
 - ▶ In regression, for each input value x , the predicted distribution is Gaussian and is then fully characterized by its mean and variance

Gaussian distributions refresher

- ▶ **Multivariate Gaussian distribution** $x \sim \mathcal{N}(\mu, \Sigma), x \in R^n$
 - ▶ $p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{1/2}} \exp(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu))$
- ▶ **Summation (a)**
 - ▶ Let x and y two random variables with the same dimensionality, $p(x) = \mathcal{N}(\mu_x, \Sigma_x)$ and $p(y) = \mathcal{N}(\mu_y, \Sigma_y)$
 - ▶ Then their sum is also Gaussian: $p(x + y) = \mathcal{N}(\mu_x + \mu_y, \Sigma_x + \Sigma_y)$
- ▶ **Marginalization (b)**
 - ▶ Let $x, p(x) = \mathcal{N}(\mu, \Sigma)$, consider a partition of x into two sets of variables $x = \begin{pmatrix} x_a \\ x_b \end{pmatrix}$.
 - ▶ Let us denote $\mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}, \Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$
 - ▶ Then the marginals are also Gaussians, e.g.: $p(x_a) = \int_{x_b} p(x_a, x_b; \mu, \Sigma) dx_b = \mathcal{N}(\mu_a, \Sigma_{aa})$,
 - ▶ Σ being symmetric, $\Sigma_{ab} = \Sigma_{ba}$
- ▶ **Conditioning (c)**
 - ▶ The conditionals are also Gaussians
 - ▶ $p(x_a | x_b) = \mathcal{N}(\mu_{a|b}, \Sigma_{a|b})$ with $\mu_{a|b} = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (x_b - \mu_b)$ and $\Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}$
- ▶ **Marginalization bis (d)**
 - ▶ Let x and y two random vectors such that $p(x) = \mathcal{N}(\mu, \Sigma_x)$ and $p(y|x) = \mathcal{N}(Ax + b, \Sigma_y)$
 - ▶ The marginal of y is $p(y) = \int p(y|x)p(x)dx = \mathcal{N}(A\mu + b, \Sigma_y + A\Sigma_x A^T)$

Introducing the Gaussian processes

From Bayesian linear regression to Gaussian processes

- ▶ Consider the linear parameter model:
 - ▶ $y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$
 - ▶ where $\mathbf{w} \in R^M$, $\phi(\mathbf{x}) \in R^M$ are M fixed basis functions
 - ▶ For example, ϕ could be a linear function $\phi(\mathbf{x}) = (\mathbf{x}, 1)$ or ϕ could be a vector of gaussian kernels $\phi_i(\mathbf{x}) = \exp(-\frac{(\mathbf{x}-\mu_i)^2}{2s^2})$
- ▶ We consider a Bayesian setting
 - ▶ With \mathbf{w} following a prior distribution given by an isotropic Gaussian
 - ▶ $p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \alpha^{-1}I)$
 - α^{-1} is the precision parameter = the inverse variance
 - ▶ For any value of \mathbf{w} , $y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$ defines a specific function of \mathbf{x}
 - ▶ $p(\mathbf{w})$ thus defines a distribution over functions $y(\mathbf{x})$

Introducing the Gaussian processes

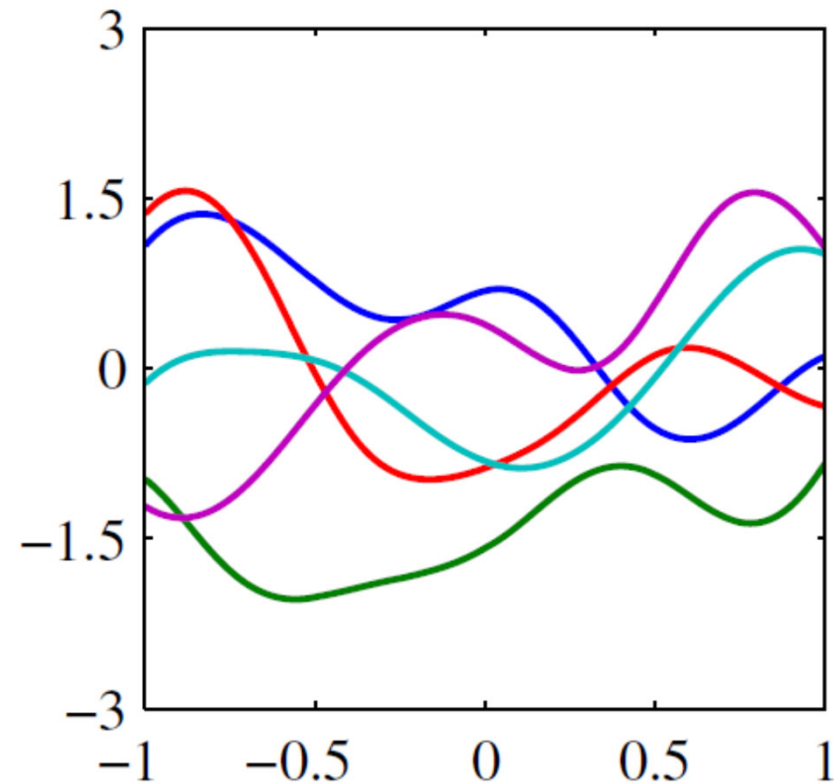
From Bayesian linear regression to Gaussian processes

- ▶ How to characterize the distribution over functions $y(x)$?
 - ▶ In practice, we will want to evaluate $y(x)$ at specific values x
 - ▶ e.g. at the training points or for a test point
 - ▶ Let us consider a finite data sample x^1, \dots, x^N
 - ▶ Let us denote $\mathbf{y} = (y^1, \dots, y^N)^T$, with $y^i = y(x^i)$
 - ▶ We want to characterize the distribution of \mathbf{y}
 - ▶ $\mathbf{y} = \Phi \mathbf{w}$, with $\Phi = [\phi(x^1), \dots, \phi(x^N)]^T$ called the design matrix $\Phi_{ij} = \phi_j(x^i)$
 - ▶ \mathbf{w} is $M \times 1$, Φ is $N \times M$, \mathbf{y} is $N \times 1$
 - ▶ \mathbf{y} being a linear combination of Gaussian variables (the elements of \mathbf{w}) is itself Gaussian and fully characterized by its mean and variance
 - ▶ $E[\mathbf{y}] = \Phi E[\mathbf{w}] = \mathbf{0}$
 - ▶ $Cov[\mathbf{y}] = E[\mathbf{y}\mathbf{y}^T] = \Phi E[\mathbf{w}\mathbf{w}^T] \Phi^T = \frac{1}{\alpha} \Phi \Phi^T = \mathbf{K}$
 - ▶ \mathbf{K} is a **Gram matrix** with elements $K_{nm} = k(x^n, x^m) = \frac{1}{\alpha} \phi(x^n)^T \phi(x^m)$
 - $k(x, x')$ is the **kernel function**
- $\mathbf{y} = \mathcal{N}(\mathbf{0}, \mathbf{K})$
- ▶ This is a first example of Gaussian process, defined by a linear model
 - ▶ Usually, the kernel function is not defined through basis functions, but directly by specifying a Kernel function, e.g. a Gaussian kernel

Introducing the Gaussian processes

From Bayesian linear regression to Gaussian processes

- ▶ Samples of functions drawn from Gaussian processes for a « Gaussian Kernel »
 - ▶ $k(x, x') = \exp(-\frac{\|x-x'\|^2}{2\sigma^2})$
 - ▶ We specify a set of input points $x = (x^1, \dots, x^N)$ in $[-1, 1]$ and an $N \times N$ covariance matrix K .
 - ▶ We draw a vector (y^1, \dots, y^N) from the Gaussian defined by $\mathbf{y} = \mathcal{N}(\mathbf{0}, K)$
- ▶ Bishop C. PRML



Gaussian processes

► Definition

- A stochastic process is a collection of random variables $\{f(x); x \in \mathcal{X}\}$ indexed by elements of set \mathcal{X} (in the following one will consider $\mathcal{X} = R$).
 - This is a probability distribution over the functions $f(x)$
- A Gaussian process is a stochastic process such that the set of values of $f(x)$ evaluated at any number of points x^1, \dots, x^N is jointly Gaussian, i.e.:

$$\begin{bmatrix} f(x^1) \\ \vdots \\ f(x^N) \end{bmatrix} \sim N \left(\begin{bmatrix} m(x^1) \\ \vdots \\ m(x^N) \end{bmatrix}, \begin{bmatrix} k(x^1, x^1) & \dots & k(x^1, x^N) \\ \vdots & \ddots & \vdots \\ k(x^N, x^1) & \dots & k(x^N, x^N) \end{bmatrix} \right)$$

► Properties

- A Gaussian process is entirely specified by its
 - Mean **function** $m(x) = E[f(x)]$
 - Covariance **function** $k(x, x') = E[(f(x) - m(x))(f(x') - m(x'))]$
- One denotes $f \sim GP(m, k)$ meaning that f is distributed as a GP with mean m and covariance k functions

Gaussian processes

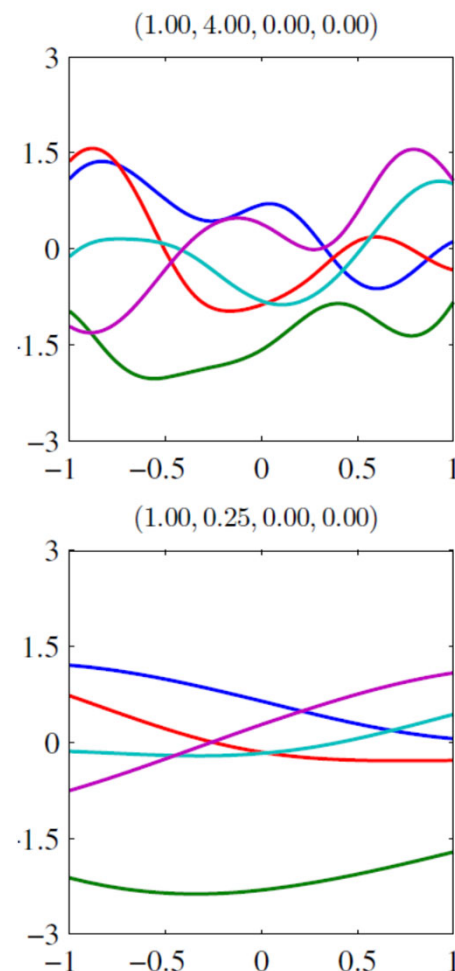
► Intuition

- Gaussian distributions model finite collections of real valued variables
- Gaussian processes extend multivariate gaussians to infinite collections of real-valued variables
 - GP are distributions over random functions
 - Let H be a class of functions $f: X \rightarrow Y$. A random function $f(\cdot)$ from H is a function which is randomly drawn from H
 - Intuitively, one can think of $f(\cdot)$ as an infinite vector drawn from an infinite multivariate Gaussian. Each dimension of the Gaussian corresponds to an element x from the index and the corresponding component of the random vector is the value $f(x)$
- What could be the functions $m(\cdot)$ and $k(\cdot, \cdot)$?
 - Any real valued function $m(\cdot)$ is acceptable
 - K should be a valid covariance matrix corresponding to a Gaussian distribution
 - This is the case if K is positive semi-definite (remember conditions for valid kernels)
 - Any valid kernel can be used as a covariance function

Gaussian processes

► Example

- Zero mean Gaussian process $GP(0, k(\cdot, \cdot))$ defined for functions $h: X \rightarrow R$
- $k(x, x') = \exp(-\frac{\theta_1}{2} ||x - x'||^2)$
- The function values are distributed around 0
- $f(x)$ and $f(x')$ will have a high covariance $k(x, x')$ if x and x' are nearby and a low covariance otherwise
 - i.e. they are locally smooth



Bishop PRML, Top $\theta_1 = 4$, bottom $\theta_1 = 0.25$

Gaussian processes for regression

- ▶ We consider a Gaussian process regression model (1 dimensional for simplification)
 - ▶ $y = f(x) + \epsilon$, with $x \in R^n$ and $y \in R$
 - ▶ $\epsilon \sim \mathcal{N}(0, \sigma^2)$ independently chosen for each observation accounts for the noise at each observation
 - ▶ Let us consider a set of training examples $S = \{(x^1, y^1), \dots, (x^N, y^N)\}$ from an unknown distribution
 - ▶ Let us denote $Y = (y^1, \dots, y^N)^T$ and $F = (f^1, \dots, f^N)^T$ with $f^i = f(x^i)$
 - ▶ From the definition of a Gaussian process, one assume a prior distribution over functions $f(\cdot)$. We assume a zero mean Gaussian process prior:
 - ▶ $p(F) = \mathcal{N}(0, K)$ with K a Gram matrix defined by a kernel function $K_{ij} = k(x_i, x_j)$
- ▶ We will
 - ▶ Characterize the joint distribution of $Y = (y^1, \dots, y^N)^T$
 - ▶ In order to define the predictive distribution for test points $p(y_{N+1}|Y)$

Gaussian processes for regression

Characterizing the joint distribution of $Y = (y^1, \dots, y^N)^T$

▶ The joint distribution of $Y = (y^1, \dots, y^N)^T$ is

▶ $p(Y) = \int p(Y|F)p(F)df = \mathcal{N}(O, C)$

▶ With the covariance matrix C defined as $C(x^i, x^j) = k(x^i, x^j) + \frac{1}{\sigma^2} \delta_{ij}$

▶ δ_{ij} is the Kronecker symbol

▶ Demonstration

▶ We will first show $p(Y|F) = \mathcal{N}(F, \frac{1}{\sigma^2} I_N)$

▶ $p(Y|F) = p(y^1, \dots, y^N | F)$

▶ $p(Y|F) = \prod_{i=1}^N p(y^i | f^i)$

▶ $p(Y|F) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{1}{2\sigma^2} (y^i - f^i)^2)$

▶ $p(Y|F) = \left(\frac{1}{2\pi\sigma^2}\right)^{N/2} \exp(-\frac{1}{2\sigma^2} \|Y - F\|^2)$

▶ $p(Y|F) = \mathcal{N}(F, \frac{1}{\sigma^2} I_N)$

Gaussian processes for regression

Characterizing the joint distribution of $Y = (y^1, \dots, y^N)^T$

► Demonstration of $p(Y) = \int p(Y|F)p(F)df = \mathcal{N}(O, C)$

► $p(F) = \mathcal{N}(0, K)$

► $p(Y|F) = \mathcal{N}(F, \frac{1}{\sigma^2} I_N)$

► $p(Y) = \int p(Y|F)p(F)df$

► By property (d) in Gaussian refresher we get:

► $p(Y) = \mathcal{N}\left(O, \frac{1}{\sigma^2} I_N + K\right) = \mathcal{N}(O, C)$

► With $C_{ij} = k(x^i, x^j) + \frac{1}{\sigma^2} \delta_{ij}$

Gaussian processes for regression

Predictive distribution

- ▶ For the regression, our goal is to predict the value y for a new observation x
 - ▶ Let us consider a training set $D = \{(x^i, y^i); i = 1 \dots N\}$, and denote $Y^N = (y^1, \dots, y^N)^T$, let y^{N+1} the value one wants to predict for observation x^{N+1} , $Y^{N+1} = (Y^N, y^{N+1})^T$
- ▶ Let us first explicit the joint distribution over Y^{N+1}
 - ▶ $p(Y^{N+1}) = \mathcal{N}(0, C_{N+1})$ with $C_{N+1} = \begin{pmatrix} C_N & k \\ k^T & c \end{pmatrix}$
 - ▶ C_N the covariance matrix of Y_N
 - ▶ $k \in R^N$ $k_i = k(x^i, x^{N+1}); i = 1 \dots N$
 - ▶ $c = k(x^{N+1}, x^{N+1}) + \sigma^2 \in R$
 - ▶ **Proof**
 - ▶ This is a direct application of the result shown before $p(Y) = \mathcal{N}(0, C)$

Gaussian processes for regression

Predictive distribution

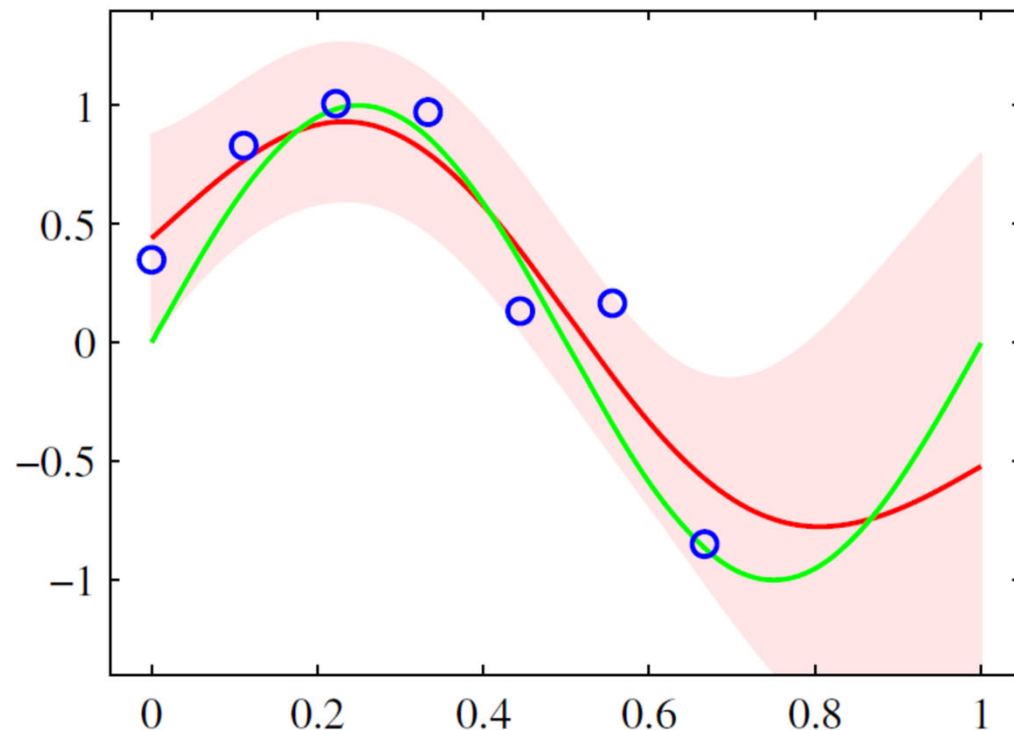
- ▶ Prediction is achieved via the conditional distribution $p(y_{N+1}|Y)$
 - ▶ By definition of a Gaussian process, $p(y_{N+1}|Y, X)$ is a Gaussian.
 - ▶ Its mean and covariance are given by:
 - ▶ $m(x_{N+1}) = k^T C_N^{-1} Y$
 - ▶ $\sigma^2(x_{N+1}) = c - k^T C_N^{-1} k$
 - ▶ **Proof**
 - ▶ This is a direct application of property (c) (conditioning)
- ▶ **Property**
 - ▶ $m(x_{N+1})$ writes as $m(x_{N+1}) = \sum_{i=1}^N a_i k(x_i, x_{N+1})$
 - ▶ With a_i the i^{th} component of $C_N^{-1} Y$

Gaussian processes for regression

Predictive distribution

- ▶ This means that for any new datum x^{N+1} , one can compute
 - ▶ A mean prediction $m(x_{N+1})$
 - ▶ An uncertainty associated to this prediction $\sigma^2(x_{N+1})$

Illustration of Gaussian process regression applied to the sinusoidal data set in Figure A.6 in which the three right-most data points have been omitted. The green curve shows the sinusoidal function from which the data points, shown in blue, are obtained by sampling and addition of Gaussian noise. The red line shows the mean of the Gaussian process predictive distribution, and the shaded region corresponds to plus and minus two standard deviations. Notice how the uncertainty increases in the region to the right of the data points.



Learning hyperparameters

- ▶ The kernel functions can be chosen a priori
- ▶ Alternatively, they may be defined as parametric functions (e.g. squared exponential kernel as in the example) and the parameters may be learned e.e. by maximum likelihood
 - ▶ Log likelihood for a Gaussian process regression model
 - ▶ $\log p(Y|\theta) = -\frac{1}{2}\log|C_N| - \frac{1}{2}Y^T C_N^{-1}Y - \frac{N}{2}\log(2\pi)$
 - ▶ Training can be performed using gradient descent on the parameters θ