

# Introduction to Machine Learning. CSCI-UA 9473, Lecture 5.

Augustin Cosse

Ecole Normale Supérieure, DMA & NYU  
Fondation Sciences Mathématiques de Paris.



2018

# What have we seen so far?

- ▶ Bayesian framework and estimators, prior, posterior, MLE, MAP
- ▶ Supervised Learning
  - ▶ Linear regression
    - ▶ Bias variance trade-off (Linear and non linear data)
    - ▶ Regularization (Ridge, Lasso, Subset Selection)
  - ▶ Linear classification
    - ▶ Discriminative vs Generative classifiers
    - ▶ Least squares
    - ▶ Logistic regression

# Today

- ▶ Geometry of separating hyperplanes (recap) and Rosenblatt perceptron
- ▶ The curse of dimensionality
- ▶ A word on non parametric classifiers
- ▶ Kernel methods
- ▶ Support vector machines (SVM)

# Parametric vs Non Parametric (I)

- ▶ Remember the difference between parametric and non parametric methods ?
- ▶ Linear regression = linear combination of a **fixed number** of (possibly non linear) basis functions

$$Y = \beta_0 + \sum_{k=1}^d \beta_k X_k$$

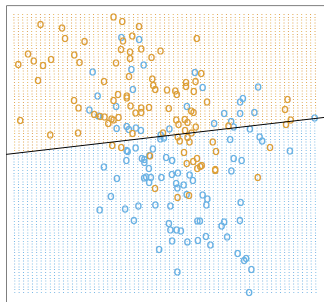
$$Y = \beta_0 + \sum_{k=1}^d \beta_k \phi_k(X)$$

- ▶ **Linearity** in the parameters leads to **interesting properties** such as **closed form solution**, **computational tractability**,...

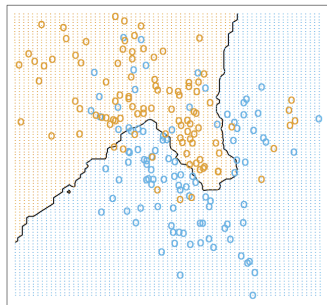
# Parametric vs Non Parametric (II)

- ▶ Today we will discuss non parametric models
- ▶ **Parametric** = fixed number of parameters, **Non parametric** = number of parameters/model complexity **grows** with the amount of **training** samples

Linear Regression of 0/1 Response



15-Nearest Neighbor Classifier



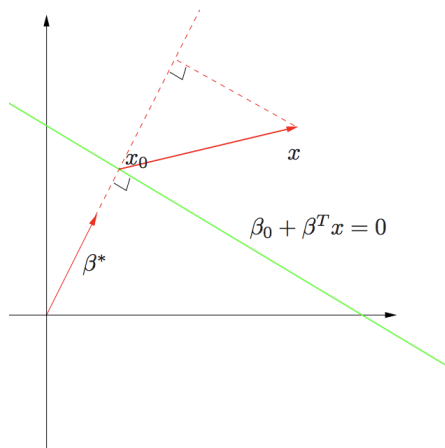
(from H.,T.,F., The Elem. of Stat. Learn.)

# Separating Hyperplanes (quick recap)

- ▶ Consider the **separating hyperplane**  $\beta_0 + \beta^T \mathbf{x}$
- ▶  $\mathbf{x}_1$  and  $\mathbf{x}_2$  **belong** to the plane if they satisfy

$$\beta_0 + \beta^T \mathbf{x}_1 = \beta_0 + \beta^T \mathbf{x}_2$$

- ▶ We thus have  $\beta^T (\mathbf{x}_1 - \mathbf{x}_2) = 0$  for all  $\mathbf{x}_1, \mathbf{x}_2$  in the plane
- ▶  $\beta$  ( $\beta^* = \beta / \|\beta\|$ ) is the vector **normal** to the hyperplane

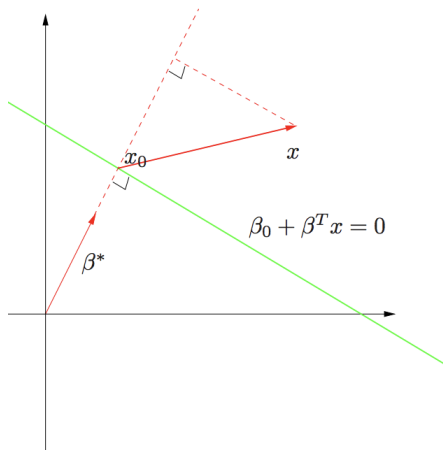


H,T,F, Elem. of Stat. Learn.

- ▶ The **signed distance** of a point  $\mathbf{x}$  to the hyperplane is defined as

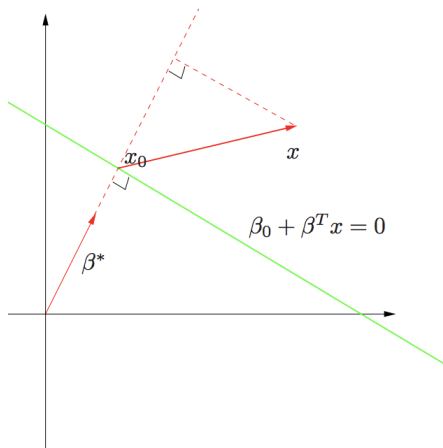
$$\begin{aligned} & (\beta^*)^T (\mathbf{x} - \mathbf{x}_0) \\ &= \frac{1}{\|\beta\|} (\beta^T \mathbf{x} + \beta_0) \end{aligned}$$

- ▶ Points that are located **above** thus lead to **positive values**  $\beta^T \mathbf{x} + \beta_0 > 0$
- ▶ Points that are located **below** lead to **negative values**  $\beta^T \mathbf{x} + \beta_0 < 0$



H,T,F, Elem. of Stat. Learn.

- ▶ A separating plane thus gives a **natural way** to associate **positive** or **negative labels** to points
- ▶ For a **two class** classification problem, we can look for the **plane** that gives **positive** labels to one class and **negative** labels to the other
- ▶ This idea leads to the **perceptron** algorithm of Rosenblatt



H,T,F, Elem. of Stat. Learn.



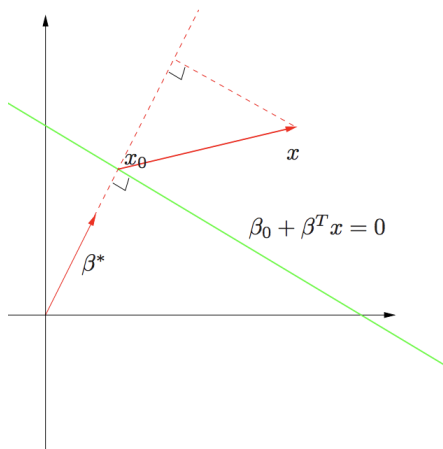
- ▶ The **perceptron** thus simply reads as

$$y(\mathbf{x}) = f(\beta_0 + \beta^T \mathbf{x})$$

Where

$$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases}$$

- ▶ During training, we associate **+1** labels to points in cluster  $\mathcal{C}_1$  and **-1** labels to points in cluster  $\mathcal{C}_2$



H,T,F, Elem. of Stat. Learn.

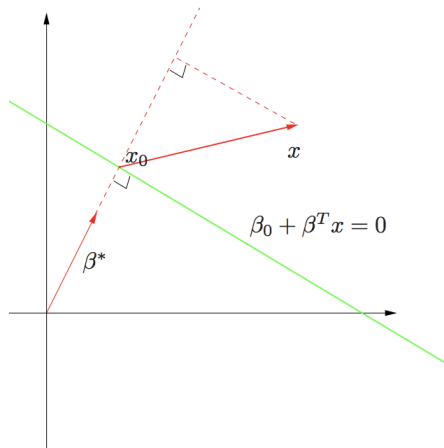
- ▶ In the perceptron, a point in  $\mathcal{C}_1$  ( $y_i = +1$ ) is thus **misclassified** if  $\beta^T \mathbf{x}_i + \beta_0 < 0$
- ▶ Generally **we would like** all points to satisfy

$$y_i(\beta^T \mathbf{x}_i + \beta_0) > 0$$

- ▶ so we minimize

$$-\sum_{i \in \mathcal{M}} y_i(\beta^T \mathbf{x}_i + \beta_0)$$

(contributions should be  $\geq 0$ )



H,T,F, Elem. of Stat. Learn.

# Perceptron

(Perceptron)

$$D(\boldsymbol{\beta}, \beta_0) = - \sum_{i \in \mathcal{M}} y_i (\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i)$$

- ▶ How do we **train** the perceptron?
- ▶ One way is to use **stochastic gradient descent** (we will come back to that idea later)
- ▶ General idea (**perceptron learning algorithm**)
  - ▶ Choose initial vector of prefactors  $\boldsymbol{\beta}$
  - ▶ Then for each misclassified points  $\mathbf{x}_n$ , do

$$\begin{bmatrix} \boldsymbol{\beta}^{k+1} \\ \beta_0^{k+1} \end{bmatrix} \leftarrow \begin{bmatrix} \boldsymbol{\beta}^k \\ \beta_0^k \end{bmatrix} - \eta \nabla D^n(\boldsymbol{\beta}, \beta_0), \quad D^n = y_n (\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_n)$$

# Perceptron: intuition

Perceptron w/  
gen. features  $\phi(X)$

$$D(\boldsymbol{\beta}, \beta_0) = - \sum_{i \in \mathcal{M}} y_i (\beta_0 + \boldsymbol{\beta}^T \boldsymbol{\phi}_i)$$

$$\begin{bmatrix} \boldsymbol{\beta}^{k+1} \\ \beta_0^{k+1} \end{bmatrix} \leftarrow \begin{bmatrix} \boldsymbol{\beta}^k \\ \beta_0^k \end{bmatrix} - \eta \nabla D^n(\boldsymbol{\beta}, \beta_0), \quad D^n = -y_n (\beta_0 + \boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n))$$

$$\text{if no intercept} \quad \boldsymbol{\beta}^{k+1} \leftarrow \boldsymbol{\beta}^k + \eta \boldsymbol{\phi}_n y_n$$

- ▶ Perceptron **convergence Theorem**: If there exists an exact solution (data is linearly separable), then the perceptron algorithm is guaranteed to find an exact solution in a finite number of steps.

# Perceptron: more intuition

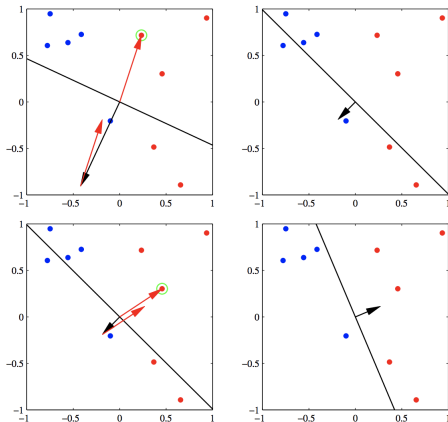


Figure 4.7 Illustration of the convergence of the perceptron learning algorithm, showing data points from two classes (red and blue) in a two-dimensional feature space  $(\phi_1, \phi_2)$ . The top left plot shows the initial parameter vector  $w$  shown as a black arrow together with the corresponding decision boundary (black line), in which the arrow points towards the decision region which classified as belonging to the red class. The data point circled in green is misclassified and so its feature vector is added to the current weight vector, giving the new decision boundary shown in the top right plot. The bottom left plot shows the next misclassified point to be considered, indicated by the green circle, and its feature vector is again added to the weight vector giving the decision boundary shown in the bottom right plot for which all data points are correctly classified.

- ▶ Let us assume **no intercept** (data has been centered)
- ▶ For each **misclassified** points  $\mathbf{x}_n$  (resp.  $\phi(\mathbf{X}_n)$ ), the algorithm **adds** the **pattern** of the misclassified point to the **weight vector**  $\beta$

$$\begin{aligned}
 & - (\beta^{k+1})^T \phi_n y_n \\
 = & - (\beta^k)^T \phi_n y_n \\
 & - (\phi_n y_n)^T \phi_n y_n \\
 < & - (\beta^k)^T \phi_n y_n
 \end{aligned}$$

Bishop, Pattern Recogn. and ML.

# Curse of dimensionality

- ▶ Two difficulties in (generalized) linear models are related to the **dimension**
  1. Fitting a **simple linear model** is **tricky** in **high dimension** because when  $d \gg n$ , the solution is usually not unique. (one possibility is to use regularization)
  2. As soon as we consider **advanced features**, the number of coefficients needed grows **non linearly** with the dimension. Ex: linear regression model made from degree-3 polynomial features

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_i + \sum_{k=1}^D \sum_{j=1}^D w_{ij} x_i x_j + \sum_{i,j,k}^D w_{i,j,k} x_i x_j x_k$$

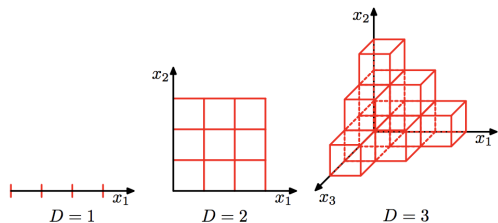
# Curse of dimensionality

- ▶ Those difficulties are part of a general phenomenon known as **curse of dimensionality**
- ▶ The expression **curse of dimensionality** which which was apparently coined by Bellman (1961) was roriginally eferring to the fact that many algorithms that work fine in low dimension become **intractable** in **high dimension**.
- ▶ In **machine Learning**, it refers to **much more** and in particular, it includes the idea that **generalization** becomes **harder** in **large dimension** because a **fixed size training set** covers a **dwindling fraction of the space** as  $d$  increases. (see Dave Donoho, *High-Dimensional Data Analysis: The Curses and Blessings of Dimensionality*, Pedro Domingos, *A few useful things to know in machine Learning*)

# Curse of dimensionality

- ▶ As an example, consider a regression model that would divide the space into subcells.
- ▶ As the dimension grows, we would need an **exponentially large** amount of training data to ensure that the **cells are not empty**

**Figure 1.21** Illustration of the curse of dimensionality, showing how the number of regions of a regular grid grows exponentially with the dimensionality  $D$  of the space. For clarity, only a subset of the cubical regions are shown for  $D = 3$ .



Bishop, Pattern Recognition and ML



# Kernels

- ▶ There are **two** main **uses** of kernels in machine learning
- ▶ Either as a way to **encode similarity** between inputs
  - ▶ **So far** we have assumed that the **data** could be **represented by** means of some **feature vector**  $X = (X_1, X_2, \dots, X_n)$
  - ▶ **Computing** explicit **features** is often **difficult** and we **only have access** to some form of **similarity/dissimilarity** between the samples (think of comparing texts for example)
  - ▶ Kernel trick = turning classification/regression models based on features to models based on similarity/kernel
- ▶ Kernels can also be used for **localization**, to build **smooth** (generative) models for regression and classification (Kernel smoothing)

# Kernels

- ▶ We define a kernel to be a function of **two arguments**  $\kappa(\mathbf{x}, \mathbf{x}')$  that is **symmetric**,  $\kappa(\mathbf{x}, \mathbf{x}') = \kappa(\mathbf{x}', \mathbf{x})$  and **non negative**  $\kappa(\mathbf{x}, \mathbf{x}') \geq 0$  (i.e the idea is to use it as a measure of similarity)
- ▶ Two important examples are the **Gaussian** kernel and the class of **RBF** kernels

$$\kappa(\mathbf{x}', \mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')\boldsymbol{\Sigma}(\mathbf{x} - \mathbf{x}')\right)$$

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

( $\sigma$  here is called the **bandwidth**)

# Kernels

- ▶ When used to encode similarity between points, Kernel are often represented through their corresponding **Gram matrix**

$$\mathbf{K} = \begin{pmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_1) & \dots & \kappa(\mathbf{x}_1, \mathbf{x}_N) \\ & \vdots & \\ \kappa(\mathbf{x}_N, \mathbf{x}_1) & \dots & \kappa(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

- ▶ If the only thing we can compute with the data is the **similarity matrix**  $\mathbf{K}$ , it is usually hard to extract the feature map from  $\mathbf{K}$ .
- ▶ However, there is an important class of kernels for which the existence of such feature map is guaranteed
- ▶ **For these kernels**, the feature vector can be computed from the Gram matrix

# Mercer Kernels

- ▶ A **Mercer Kernel** is a Kernel for which the **Gram matrix  $\mathbf{K}$**  is **positive definite** for any set of inputs  $\{\mathbf{x}_i\}$
- ▶ When the kernel is Mercer, there **always** exists a **feature** mapping  $\phi(\mathbf{x})$  such that

$$\mathbf{K}_{ij} = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

- ▶ In other words, the kernel corresponds to an inner product in some finite feature space.
- ▶ **Mercer Kernel** are especially **useful** when used in **Support Vector Machines** (SVM) because they guarantee that there exists a **unique optimal** separating **hyperplane** in the feature space.

# Kernel machines

- ▶ There are **two main approaches** to define models based on kernels
  1. Use any generalized linear model where you replace the features by kernels centered at centroids  $\mu_1, \dots, \mu_K$  (**Kernel machine**)

$$\hat{f}(\mathbf{x}) = \sum_{n=1}^K w_n \kappa(\mathbf{x}, \mu_n)$$

When using **Radial basis functions**, this model is known as **RBF network**

2. Instead of building a new model from a feature vector defined in terms of kernels, one can instead start from existing models and **replace** all **inner products**  $\langle \mathbf{x}, \mathbf{x}' \rangle$  by a call to the Kernel function  $\kappa(\mathbf{x}, \mathbf{x}')$ . This idea is known as the **Kernel trick**

# Kernel machines

- ▶ The main issue with kernel machines is how to choose the centroids
- ▶ When in low dimension, one can choose the centroids to uniformly tile the space (but this will fail in high dimension because of the curse of dimensionality)
- ▶ Another approach could be to optimize over the centroids (but the problem is highly multimodal  $\Rightarrow$  finding the optimum is hard)
- ▶ One could also start by finding clusters in the data and then assign a centroid to each cluster (That would require choosing a number of clusters)

# Kernel machines

- ▶ Finally a last approach, which is the **simplest** is to take each sample  $\mathbf{x}_i$  as a prototype. The feature vector is then given by

$$\phi(\mathbf{x}) = [\kappa(\mathbf{x}, \mathbf{x}_1), \dots, \kappa(\mathbf{x}, \mathbf{x}_n)]$$

- ▶ The advantage is that the model is now fully non parametric
- ▶ But it **requires many kernels**
- ▶ One solution to **select** a **subset** of these kernels is to use any of the **regularization** penalties that we have studied in inear regression
- ▶ This gives models known as **L1** or **L2** regularized Vector Machines

## The Kernel trick (I)

- ▶ The Kernel trick starts with existing models and tries to **replace** all **inner products** in those models with a call to the Kernel to define a corresponding version of these models that would rely only on similarity
- ▶ For this trick to work, the kernel **should** be a **Mercer kernel**
- ▶ As an example, consider **Ridge regression**. In this model the objective is given by

$$J(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|^2$$

And the optimal solution can be computed exactly (cfr assignment) as

$$\boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^T \mathbf{y} = \left( \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T + \lambda \mathbf{I}_D \right)^{-1} \mathbf{X}^T \mathbf{y}$$

Here  $\mathbf{X}$  is the  $N \times D$  design matrix encoding the points



## The Kernel trick (II)

$$\boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^T \mathbf{y} = \left( \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T + \lambda \mathbf{I}_D \right)^{-1} \mathbf{X}^T \mathbf{y}$$

- ▶ By using an inversion trick, one can write  $\boldsymbol{\beta}$  equivalently as

$$\boldsymbol{\beta} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$$

- ▶ But now  $\mathbf{X} \mathbf{X}^T$  is exactly the Gram matrix  $\mathbf{K}$ , i.e we can write

$$\boldsymbol{\beta} = \mathbf{X}^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$$

## The Kernel trick (III)

$$\boldsymbol{\beta} = \mathbf{X}^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$$

- ▶ Now let  $\boldsymbol{\beta} = \mathbf{X}^T \boldsymbol{\alpha}$  with  $\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$  and substitute this in the regression model, we get

$$\begin{aligned} \hat{f}(\mathbf{x}) &= \boldsymbol{\beta}^T \mathbf{x} = \boldsymbol{\alpha}^T \mathbf{X} \mathbf{x} = \sum_{i=1}^N \alpha_i \langle \mathbf{x}_i, \mathbf{x} \rangle \\ &= \sum_{i=1}^N \alpha_i \kappa(\mathbf{x}, \mathbf{x}_i) \end{aligned}$$

## The Kernel trick (III)

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^N \alpha_i \kappa(\mathbf{x}, \mathbf{x}_i), \quad \boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$$

- ▶ The cost of computing the dual variables  $\boldsymbol{\alpha}$  is  $\mathcal{O}(N^3)$  whereas the cost of computing the primal variables  $\boldsymbol{\beta}$  is  $\mathcal{O}(D^3)$ . The kernel method is thus essentially **useful** in **high dimension**
- ▶ However, prediction using the dual variables  $\boldsymbol{\alpha}$  takes  $\mathcal{O}(ND)$  time whether prediction using the primal variables  $\boldsymbol{\beta}$  takes  $\mathcal{O}(D)$  time. Making  $\boldsymbol{\alpha}$  **sparse** (few non zero entries) can **speed up** prediction  $\Rightarrow$  **that is precisely the point of SVMs !**

# Usage

- ▶ There are two main frameworks in which you might want to use Kernels as a way to encode similarity
  1. You have some data and you've come up with a function which you think might be a good way to encode similarity of your data
  2. You don't want to explicitly deal with feature vectors
  3.  $D > N$
- ▶ Kernels are especially useful when combined with Support Vector Machines (see later slides)

# Smoothing Kernels

- ▶ A smoothing kernel is a function  $\kappa(\mathbf{x})$  of one argument which decreases sufficiently quickly

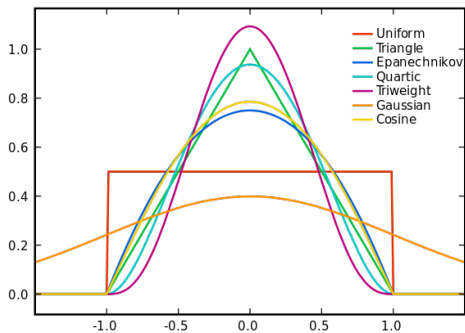
$$\int \kappa(x) dx = 1, \quad \int x\kappa(x) dx = 0, \quad \text{and} \quad \int x^2\kappa(x) dx > 0$$

- ▶ A useful example is the Gaussian kernel

$$\kappa(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

- ▶ When the inputs are vectors, we can simply take any kernel and use it with the norm  $\|\mathbf{x}\|$ , (e.g.  $\|\mathbf{x}\|^2 = \sum_{n=1}^N X_n^2$ )

# Smoothing Kernels: examples



# Smoothing Kernels

- ▶ Smoothing kernels are essentially used in local regression models to get an approximation for the conditional expectation  $\mathbb{E}\{y|\mathbf{x}\}$
- ▶ In this framework, the kernels are used to approximate  $p(\mathbf{x}, y)$

$$\mathbb{E}\{y|\mathbf{x}\} = \int yp(y|\mathbf{x}) dy = \frac{\int p(\mathbf{x}, y)y dy}{\int p(\mathbf{x}, y) dy}$$

- ▶ Then use smoothing kernels to define an approximation of the

$$p(\mathbf{x}, y) \approx \frac{1}{N} \sum_{i=1}^N \kappa_x(\mathbf{x} - \mathbf{x}_i) \kappa_y(y - y_i)$$

- ▶ The model is well defined as the Kernels integrate to 1 ( $p(\mathbf{x}, y)$  corresponds to a probability)

# Smoothing Kernels

$$p(\mathbf{x}, y) \approx \frac{1}{N} \sum_{i=1}^N \kappa_x(\mathbf{x} - \mathbf{x}_i) \kappa_y(y - y_i)$$

- ▶ Using the kernel decomposition for the joint probability distribution, we can write the conditional expectation (which gives one possible regression model) as

$$\begin{aligned} \mathbb{E}\{y|\mathbf{x}\} &= \int y p(y|\mathbf{x}) dy = \frac{\int p(\mathbf{x}, y) y dy}{\int p(\mathbf{x}, y) dy} \\ &= \frac{\frac{1}{N} \sum_{i=1}^N \kappa_x(\mathbf{x} - \mathbf{x}_i) \int y \kappa_x(y - y_i) dy}{\frac{1}{N} \sum_{i=1}^N \kappa_x(\mathbf{x} - \mathbf{x}_i) \int \kappa_y(y - y_i) dy} \\ &= \frac{\sum_{i=1}^N \kappa_x(\mathbf{x} - \mathbf{x}_i) y_i}{\sum_{i=1}^N \kappa_x(\mathbf{x} - \mathbf{x}_i)} \end{aligned}$$



## Smoothing Kernels

$$\mathbb{E} \{y|\mathbf{x}\} = \frac{\sum_{i=1}^N \kappa_x(\mathbf{x} - \mathbf{x}_i) y_i}{\sum_{i=1}^N \kappa_x(\mathbf{x} - \mathbf{x}_i)}$$

- ▶ The result above follows from the properties of smoothing kernels. In particular we use the **zero mean** property  
To get  $\int y \kappa_y(y - y_i) dy = y_i$ , let  $y' = y - y_i$  and use

$$\begin{aligned} \int x \kappa(x) dx = 0 &\Rightarrow \int (y' + y_i) \kappa_y(y) dy' \\ &= \int y' \kappa_y(y') dy' + y_i \int \kappa_y(y') dy' \\ &= 0 + y_i = y_i \end{aligned}$$

## Smoothing Kernels

- ▶ In other words, smoothing kernels give us a (non parametric) regression model of the form

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^N w_i(\mathbf{x}) y_i$$

where the weight functions are defined from the smoothing kernels as

$$w_i(\mathbf{x}) = \frac{\kappa_{\mathbf{x}}(\mathbf{x} - \mathbf{x}_i)}{\sum_{j=1}^N \kappa_{\mathbf{x}}(\mathbf{x} - \mathbf{x}_j)}$$

- ▶ The prediction is now given by a weighted combination of the outputs at the training points
- ▶ This method is known as **Nadaraya-Watson** model (a.k.a Kernel regression)

# Kernels (summary)

- ▶ Kernels as a way to replace features with similarity
  - ▶ Use Kernels machines when you don't know which features to use. We will see an example of this with SVMs.
  - ▶ Generally speaking, Kernels are interesting when  $d > n$
  - ▶ When  $d \ll n$ , (multi)linear regression is more interesting
  - ▶ When  $d \gg n$  there is no need to bring the data in higher dimensional space as finding a separating hyperplane should not be hard
- ▶ Smoothing Kernels
  - ▶ Essentially used in kernel regression or local regression models
- ▶ Kernels are also used in unsupervised learning (will discuss that later)

# Support vector machines

- ▶ Linear models have interesting computational and analytical properties but their practical applicability is limited by the curse of dimensionality
- ▶ Support vector machines are also called **sparse vector machines**
- ▶ SVM start by defining basis functions that are centered on the data and then **select** a **subset** of these during training

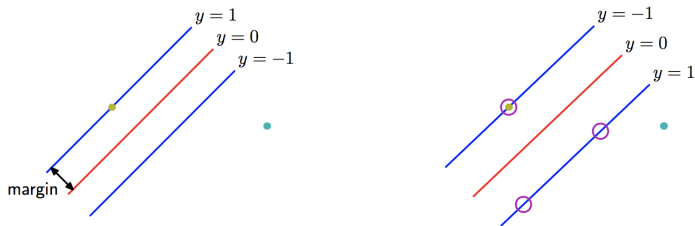
# Support vector machines

- ▶ Consider the linear regression model

$$y(\mathbf{x}) = \boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}) + \beta_0$$

- ▶ Assume we want to do classification so the labels are  $t_n = \{\pm 1\}$
- ▶ We further assume that the dataset is linearly separable in feature space so that there exists at least one separating hyperplane with  $\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) + \beta_0 > 0$  for the  $\mathbf{x}_n$  with  $t_n > 0$  and  $\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) + \beta_0 < 0$  otherwise
- ▶ When there are multiple choices we should choose the one that gives the smallest generalization error. SVM tries to achieve this through the notion of margin

# Support vector machines



**Figure 7.1** The margin is defined as the perpendicular distance between the decision boundary and the closest of the data points, as shown on the left figure. Maximizing the margin leads to a particular choice of decision boundary, as shown on the right. The location of this boundary is determined by a subset of the data points, known as support vectors, which are indicated by the circles.

(Bishop, Pattern recognition and Machine Learning)

# SVM as Maximum Margin Classifier

- ▶ Recall from the geometry of separating hyperplanes that the **distance** of a point  $\phi(\mathbf{x}_n)$  to the hyperplane  $\beta^T \mathbf{x} + \beta_0$  is defined as  $|y(\mathbf{x})|/\|\beta\|$
- ▶ When all the points are correctly classified, the sign of  $t_n$  equals the sign of  $\mathbf{y}_n = \beta^T \phi(\mathbf{x}_n) + \beta_0$  and we can thus write the distance as

$$\frac{t_n y_n}{\|\beta\|} = \frac{t_n (\beta^T \phi(\mathbf{x}_n) + \beta_0)}{\|\beta\|}$$

- ▶ The **margin** is the **perpendicular distance** of the closest point  $\phi(\mathbf{x}_n)$  to the plane

# SVM as Maximum Margin Classifier

$$\frac{t_n y_n}{\|\boldsymbol{\beta}\|} = \frac{t_n(\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) + \beta_0)}{\|\boldsymbol{\beta}\|}$$

- ▶ The **maximum margin** solution is thus given by

$$\operatorname{argmax}_{\boldsymbol{\beta}, \beta_0} \left\{ \frac{1}{\|\boldsymbol{\beta}\|} \min_n \left[ t_n(\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) + \beta_0) \right] \right\}$$

- ▶ We don't want to solve this problem because deriving a direct solution in this framework would be difficult



# SVM as Maximum Margin Classifier

$$\operatorname{argmax}_{\beta, \beta_0} \left\{ \frac{1}{\|\beta\|} \min_n \left[ t_n(\beta^T \phi(\mathbf{x}_n) + \beta_0) \right] \right\}$$

- ▶ First note that for any rescaling  $\beta \leftarrow \alpha\beta$ ,  $\beta_0 \leftarrow \beta_0\alpha$ , the objective  $t_n(\beta^T \phi(\mathbf{x}_n) + \beta_0)/\|\beta\|$  is unchanged
- ▶ We can thus focus on one of these solution (fix one particular scale for  $[\beta, \beta_0]$ ) as all the others give the same objective
- ▶ In particular we can choose to fix the scale by setting

$$t_n(\beta^T \phi(\mathbf{x}_n) + \beta_0) = 1$$

For the point that is the closest to the boundary.

# SVM as Maximum Margin Classifier

$$\operatorname{argmax}_{\beta, \beta_0} \left\{ \frac{1}{\|\beta\|} \min_n \left[ t_n (\beta^T \phi(\mathbf{x}_n) + \beta_0) \right] \right\}$$

- ▶ Now all the other points will necessarily satisfy

$$t_n (\beta^T \phi(\mathbf{x}_n) + \beta_0) \geq 1$$

- ▶ Because we fixed the distance of the closest point to the plane the original optimization problem reduces to

$$\operatorname{argmin}_{\beta} \frac{1}{2} \|\beta\|^2$$

together with the constraint  $t_n (\beta^T \phi(\mathbf{x}_n) + \beta_0) \geq 1$

# SVM as Maximum Margin Classifier

$$\begin{aligned} & \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \quad \frac{1}{2} \|\boldsymbol{\beta}\|^2 \\ & \text{subject to} \quad t_n(\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) + \beta_0) \geq 1 \end{aligned}$$

- ▶ This constrained optimization problem can be **recast** as an **unconstrained** problem by introducing multipliers  $\lambda_n \geq 0$

$$L(\boldsymbol{\beta}, \beta_0, \boldsymbol{\lambda}) = \frac{1}{2} \|\boldsymbol{\beta}\|^2 - \sum_{n=1}^N \lambda_n \left\{ t_n (\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) + \beta_0) - 1 \right\}$$

(see for example [Appendix E](#) in Bishop, Pattern Recognition and Machine Learning)

# SVM as Maximum Margin Classifier

$$L(\boldsymbol{\beta}, \beta_0, \boldsymbol{\lambda}) = \frac{1}{2} \|\boldsymbol{\beta}\|^2 - \sum_{n=1}^N \lambda_n \left\{ t_n \left( \boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n) + \beta_0 \right) - 1 \right\}$$

- ▶ To find the minimum of this function, we set the derivatives with respect to  $\boldsymbol{\beta}$  and  $\beta_0$  to zero, getting

$$\boldsymbol{\beta} = \sum_{n=1}^N \lambda_n t_n \boldsymbol{\phi}(\mathbf{x}_n)$$

$$0 = \sum_{n=1}^N t_n \lambda_n$$

- ▶ and maximize with respect to  $\lambda_n$  (large  $\lambda_n$  penalize the constraint a lot if it becomes negative)

# SVM as Maximum Margin Classifier

- ▶ Eliminating  $\beta$  and  $\beta_0$  from  $L(\beta, \beta_0, \lambda)$ , we get

$$L(\lambda) = \sum_{n=1}^N \lambda_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \lambda_n \lambda_m t_n t_m \kappa(\mathbf{x}_n, \mathbf{x}_m)$$

- ▶ with the constraints

$$\sum_{n=1}^N \lambda_n t_n = 0$$
$$\lambda_n \geq 0$$

- ▶ Maximizing  $L(\lambda)$  under the constraints above is a **quadratic programming** problem for which efficient techniques exist. Moreover when  $\kappa(\mathbf{x}_n, \mathbf{x}_m)$  is **Mercer**, there is a single solution

# SVM as Sparse Kernel Machines(I)

- ▶ Given the function  $L(\boldsymbol{\beta}, \beta_0, \boldsymbol{\lambda})$ , it is known (Karush-Kuhn-Tucker conditions) that any optimal solution must satisfy the following 3 conditions

$$\lambda_n \geq 0$$

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

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

- ▶ The last conditions have a very important consequence on SVM
- ▶ Either  $\lambda_n = 0$  or  $t_n y(\mathbf{x}_n) = 1$  (support vectors)

## SVM as Sparse Kernel Machines (II)

- ▶ Either  $\lambda_n = 0$  or  $t_n y(\mathbf{x}_n) = 1$  (support vectors)
- ▶ In particular many  $\lambda_n$  will be zero
- ▶ Using  $\boldsymbol{\beta} = \sum_{n=1}^N \lambda_n t_n \phi(\mathbf{x}_n)$ , and substituting it in  $y(\mathbf{x}) = \boldsymbol{\beta}^T \phi(\mathbf{x}) + \beta_0$ , we get the prediction model

$$y(\mathbf{x}) = \sum_{n=1}^N \lambda_n t_n \kappa(\mathbf{x}, \mathbf{x}_n) + \beta_0$$

- ▶ Which is a **combination** of the  $\lambda_n$  !

## SVM as Sparse Kernel Machines (III)

$$\text{(Karush-Kuhn-Tucker)} \quad \begin{cases} \lambda_n \geq 0 \\ t_n y(\mathbf{x}_n) - 1 \geq 0 \\ \lambda_n \{t_n y(\mathbf{x}_n) - 1\} = 0 \end{cases}$$

- ▶ Using the [Karush-Kuhn-Tucker](#) conditions, the SVM prediction model thus [reduces](#) to

$$y(\mathbf{x}) = \sum_{n \in \mathcal{S}} \lambda_n t_n \kappa(\mathbf{x}, \mathbf{x}_n) + \beta_0$$

Where  $\mathcal{S}$  are the [support vectors](#) (all remaining  $\lambda_n$ 's are 0)



## SVM as Sparse Kernel Machines (IV)

- ▶ This **sparsity** property (i.e the need to only keep a small number of support vectors) is a **key property** of SVM
- ▶ It guarantees **efficiency** of the prediction step !
- ▶ Once you know the support vectors,  $\beta$  and  $\beta_0$  can be computed using  $\beta = \sum_{n=1}^N \lambda_n t_n \phi(\mathbf{x}_n)$ , as well as the fact that **at any of the support vectors** we must have

$$t_n y_n = t_n \left( \sum_{m \in \mathcal{S}} \lambda_m t_m \kappa(\mathbf{x}_n, \mathbf{x}_m) + \beta_0 \right) = 1$$

- ▶ **Sometimes** we **average** the estimate for  $\beta_0$  over the support vectors (here  $n$ ) to get more **stability**

## Short summary

- ▶ General geometry of separating hyperplane, distance to hyperplane, perceptron
- ▶ Curse of dimensionality
- ▶ Kernels
  - ▶ As a way to encode similarity rather than features
  - ▶ As smooth interpolating functions
- ▶ SVM
  - ▶ Maximum Margin
  - ▶ Sparse Kernel machines  $\Rightarrow$  efficient prediction