

# Lecture 9

## Kernel Methods

Luigi Freda

ALCOR Lab  
DIAG  
University of Rome "La Sapienza"

December 29, 2016

## 1 Intro

- Why Kernels?
- Kernel Characterization

## 2 Kernels Functions

- RBF Kernels
- Kernels for Comparing Documents
- Mercer Kernels

## 3 Kernel-based Models

- Kernel Machines
- Kernel Trick
- Kernelized 1NN Classification
- Kernelized Ridge Regression

## 4 Support Vector Machine (SVM)

- Loss Functions
- SVM for Regression
- SVM for Classification
- Large Margin Principle

# Outline

## 1 Intro

- Why Kernels?
- Kernel Characterization

## 2 Kernels Functions

- RBF Kernels
- Kernels for Comparing Documents
- Mercer Kernels

## 3 Kernel-based Models

- Kernel Machines
- Kernel Trick
- Kernelized 1NN Classification
- Kernelized Ridge Regression

## 4 Support Vector Machine (SVM)

- Loss Functions
- SVM for Regression
- SVM for Classification
- Large Margin Principle

# Kernels

Why?

- so far, we have been assuming that each **object** that we wish to classify or cluster or process in anyway can be **represented** as a **fixed-size feature vector**  $\mathbf{x}_i \in \mathbb{R}^D$
- this may require to **preprocess** raw data in order to obtain fixed-size feature vectors
- for certain kinds of objects, it is not clear how to best represent them as fixed-sized feature vectors
- for example, how do we represent
  - ① a text document or protein sequence, which can be of **variable length**?
  - ② a molecular structure, which has **complex 3d geometry**?
  - ③ an evolutionary tree, which has **variable size** and **shape**?

# Kernels

## Why?

- common approach: assume that we have some way of measuring the **similarity between objects**, that **doesn't require preprocessing** them into feature vector format
- for example, when comparing strings, we can compute the edit **distance** between them.
- let  $\kappa(\mathbf{x}, \mathbf{x}') \geq 0$  be some **measure of similarity** between objects  $\mathbf{x}, \mathbf{x}' \in \chi$ , where  $\chi$  is some **abstract space**; we will call  $\kappa$  a **kernel function**
- we will now see together some algorithms that can be written purely in terms of kernel function computations
- we can use such algorithms when we don't have access to (or choose not to look at) the **"inside" of the objects  $\mathbf{x}_i$**  that we are processing

# Outline

## 1 Intro

- Why Kernels?
- Kernel Characterization

## 2 Kernels Functions

- RBF Kernels
- Kernels for Comparing Documents
- Mercer Kernels

## 3 Kernel-based Models

- Kernel Machines
- Kernel Trick
- Kernelized 1NN Classification
- Kernelized Ridge Regression

## 4 Support Vector Machine (SVM)

- Loss Functions
- SVM for Regression
- SVM for Classification
- Large Margin Principle

# Kernels

## General Characterization

- we define a **kernel function** to be a **real-valued function** of **two arguments**,  $\kappa(\mathbf{x}, \mathbf{x}') \in \mathbb{R}$ , for  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$
- typically the function is **symmetric**, i.e.

$$\kappa(\mathbf{x}, \mathbf{x}') = \kappa(\mathbf{x}', \mathbf{x})$$

and **non-negative**, i.e.

$$\kappa(\mathbf{x}, \mathbf{x}') \geq 0$$

- in general  $\kappa(\mathbf{x}, \mathbf{x}')$  can be interpreted as a **measure of similarity** (but this may also **not be required**)

# Outline

- 1 Intro
  - Why Kernels?
  - Kernel Characterization
- 2 Kernels Functions
  - RBF Kernels
  - Kernels for Comparing Documents
  - Mercer Kernels
- 3 Kernel-based Models
  - Kernel Machines
  - Kernel Trick
  - Kernelized 1NN Classification
  - Kernelized Ridge Regression
- 4 Support Vector Machine (SVM)
  - Loss Functions
  - SVM for Regression
  - SVM for Classification
  - Large Margin Principle



- a **radial basis function** or **RBF kernel**  $\kappa(\mathbf{x}, \mathbf{x}') \in \mathbb{R}$  is only a function of  $\|\mathbf{x} - \mathbf{x}'\|$

$$\kappa(\mathbf{x}, \mathbf{x}') = \varphi(\|\mathbf{x} - \mathbf{x}'\|)$$

- a typical example is the **Squared Exponential kernel** (SE kernel) or **Gaussian Kernel**

$$\kappa(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^T \Sigma^{-1}(\mathbf{x} - \mathbf{x}')\right) = \exp\left(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|_{\Sigma}^2\right)$$

- if  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_D)$  we obtain the **ARD kernel** (Automatic Relevance Determination)

$$\kappa(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2} \sum_{j=1}^D \frac{(x_j - x'_j)^2}{\sigma_j^2}\right)$$

$\sigma_j$  can be interpreted as defining the **characteristic length scale** of dimension  $j$

- if  $\Sigma = \sigma \mathbf{I}$  we obtain the **isotropic kernel**

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

where  $\sigma$  is called the **bandwidth**

# Outline

- 1 Intro
  - Why Kernels?
  - Kernel Characterization
- 2 **Kernels Functions**
  - RBF Kernels
  - **Kernels for Comparing Documents**
  - Mercer Kernels
- 3 **Kernel-based Models**
  - Kernel Machines
  - Kernel Trick
  - Kernelized 1NN Classification
  - Kernelized Ridge Regression
- 4 **Support Vector Machine (SVM)**
  - Loss Functions
  - SVM for Regression
  - SVM for Classification
  - Large Margin Principle

- when performing document classification or retrieval, it is useful to have a way of **comparing two documents**  $\mathbf{x}_i$  and  $\mathbf{x}_{i'}$
- if we use a **bag of words representation**, where  $x_{ij}$  is the number of times words  $j$  occurs in document  $i$ , we can use the **cosine similarity**

$$\kappa(\mathbf{x}_i, \mathbf{x}_{i'}) = \frac{\mathbf{x}_i^T \mathbf{x}_{i'}}{\|\mathbf{x}_i\|_2 \|\mathbf{x}_{i'}\|_2}$$

this quantity measures the cosine of the angle between  $\mathbf{x}_i$  and  $\mathbf{x}_{i'}$  when interpreted as vectors

- since  $\mathbf{x}_i$  is a count vector ( $x_{ij} \geq 0$ ), the cosine similarity  $\kappa(\mathbf{x}_i, \mathbf{x}_{i'}) \in [0, 1]$
- $\kappa(\mathbf{x}_i, \mathbf{x}_{i'}) = 0$  means the vectors are orthogonal and therefore have no words in common

# Outline

- 1 Intro
  - Why Kernels?
  - Kernel Characterization
- 2 **Kernels Functions**
  - RBF Kernels
  - Kernels for Comparing Documents
  - **Mercer Kernels**
- 3 **Kernel-based Models**
  - Kernel Machines
  - Kernel Trick
  - Kernelized 1NN Classification
  - Kernelized Ridge Regression
- 4 **Support Vector Machine (SVM)**
  - Loss Functions
  - SVM for Regression
  - SVM for Classification
  - Large Margin Principle

# Mercer (Positive Definite) Kernels

- some methods require that the kernel function satisfies the **requirement** that the **Gram matrix**

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

is **positive definite** for any set of inputs  $\{\mathbf{x}_i\}_{i=1}^N$

- Mercer kernels** or **positive definite kernels** satisfy the requirement  $\mathbf{K} > 0$
- it can be shown that the Gaussian kernel and the cosine similarity kernel are Mercer kernels

# Mercer (Positive Definite) Kernels

- the importance of Mercer kernels is the following result, known as **Mercer's theorem**
- if the Gram matrix is positive definite, i.e.  $\mathbf{K} > 0$  for any set of inputs  $\{\mathbf{x}_i\}_{i=1}^N$ , we can compute an eigenvector decomposition

$$\mathbf{K} = \begin{bmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_1) & \dots & \kappa(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ \kappa(\mathbf{x}_N, \mathbf{x}_1) & \dots & \kappa(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix} = \mathbf{U}^T \mathbf{\Lambda} \mathbf{U}$$

where  $\mathbf{\Lambda}$  is a diagonal matrix of eigenvalues  $\lambda_i > 0$

- now consider an element of  $\mathbf{K}$

$$k_{ij} = (\mathbf{\Lambda}^{1/2} \mathbf{U}_{:i})^T (\mathbf{\Lambda}^{1/2} \mathbf{U}_{:j}) = (\mathbf{\Lambda}^{1/2} \mathbf{u}_i)^T (\mathbf{\Lambda}^{1/2} \mathbf{u}_j)$$

- let us define  $\phi(\mathbf{x}_i) \triangleq \mathbf{\Lambda}^{1/2} \mathbf{U}_{:i}$ , then we can write

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

- hence the entries  $k_{ij}$  can be computed by performing an inner product of some **new feature vectors**  $\phi(\mathbf{x})$  that are implicitly defined by the eigenvectors in  $\mathbf{U}$

# Mercer (Positive Definite) Kernels

- in general, if the **kernel is Mercer** then there exists a function  $\phi$  mapping  $\mathbf{x} \in \chi$  to  $\phi(\mathbf{x}) \in \mathbb{R}^D$  such that

$$\kappa(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

where  $\phi$  depends on the **eigen functions** of  $\kappa$  (where  $D$  is potentially infinite)

- conclusion:** since we are able to define a similarity distance  $\kappa(\mathbf{x}, \mathbf{x}')$  in terms of an inner product, i.e.  $\kappa(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$ , the result is that we are **implicitly transforming** each **raw data sample**  $\mathbf{x} \in \chi$  into a new **feature vector**  $\phi(\mathbf{x})$  without any need to explicitly represent it

# Mercer (Positive Definite) Kernels

- for example, the **polynomial kernel**  $\kappa(\mathbf{x}, \mathbf{x}') = (\gamma \mathbf{x}^T \mathbf{x}' + r)^M$ , where  $r > 0$  is a Mercer kernel
- in this case one can show that  $\phi(\mathbf{x})$  contains all the terms up to degree  $M$
- for example with  $M = 2$  and  $\gamma = r = 1$ , we have

$$(\mathbf{x}^T \mathbf{x}' + 1)^2 = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

with  $\phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2]^T \in \mathbb{R}^6$

- the **Gaussian kernel** is also a Mercer kernel
- the feature map  $\phi$  of a **Gaussian kernel** lives in an **infinite dimensional space**: in such a case, it is clearly infeasible to explicitly represent the feature vectors
- recall that

$$\exp(x) = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

which means that we are dealing with a "polynomial degree"  $M \rightarrow \infty$



- an example of non-Mercer kernel is the **sigmoid kernel**

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

- in general, **verifying that a kernel is a Mercer kernel** is difficult, and requires techniques from functional analysis
- however, one can show that it is possible to **build up new Mercer kernels** from simpler ones using a set of **standard rules**
- for example, if  $\kappa_1$  and  $\kappa_2$  are both Mercer, so is

$$\kappa(\mathbf{x}, \mathbf{x}') = \kappa_1(\mathbf{x}, \mathbf{x}') + \kappa_2(\mathbf{x}, \mathbf{x}')$$

- deriving the feature vector  $\phi$  implied by a kernel is in general quite difficult, and only possible if the kernel is Mercer
- however, deriving a kernel from a feature vector  $\phi$  is easy

$$\kappa(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$$

- if  $\phi(\mathbf{x}) = \mathbf{x}$ , we get the **linear kernel**, defined by

$$\kappa(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

- this is useful if the original features are individually informative and the decision boundary is likely to be representable as a linear combination of the original features
- of course, when data is not linearly separable, non-linear kernels are required

- the **Matern kernel** is commonly used in Gaussian process regression and has the following form

$$\kappa(r) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left( \frac{\sqrt{2\nu}r}{l} \right)^\nu K_\nu \left( \frac{\sqrt{2\nu}r}{l} \right)$$

where  $r \triangleq \|\mathbf{x} - \mathbf{x}'\|$  with  $\nu > 0$ ,  $l > 0$  and  $K_\nu$  is a modified Bessel function

# String Kernel

- the real power of kernels arises when the inputs are structured objects
- we now describe one way of comparing two strings  $\mathbf{x}$  and  $\mathbf{x}'$  of lengths  $D$ ,  $D'$  using a string kernel
- the two strings are defined over the 20 letter alphabet  $\mathcal{A} = \{A, R, N, D, C, E, Q, G, H, I, L, K, M, F, P, S, T, W, Y, V\}$
- let  $\mathbf{x}$  be the following sequence of length 110

*IPTSALVKETLALLSTHRTLIIANETLRIPVPVHKNHQLCTEEIFQGIGTLESQ  
TVQGGTVERLFKNLSLIKKYIDGQKKKCGEERRRVNQFLDYLQEFLGVMNTEWI*

and let  $\mathbf{x}'$  be the following sequence of length 153

*PHRRDLCSRSIWLARKIRSDLTALTESYVKHQGLWSELTEAERLQENLQAYRTFHV  
LLARLLEDQQVHFTPTGDFHQAIHTLLLQVAAFAYQIEELMILLEYKIPRNEADG  
MLFEKKLWGLKVLQELSQWTVRSIHDLRFISSHQTGIP*

- these strings have the substring *LQE* in common
- we can define the **similarity of two strings** to be the number of substrings they have in common

# String Kernel

- more formally and more generally, let us say that  $s$  is a substring of  $\mathbf{x}$  if we can write  $\mathbf{x} = usv$  for some (possibly empty) strings  $u$ ,  $s$  and  $v$
- now let  $\phi_s(x)$  denote the number of times that substring  $s$  appears in string  $\mathbf{x}$
- we define the kernel between two strings  $\mathbf{x}$  and  $\mathbf{x}'$  as

$$\kappa(\mathbf{x}, \mathbf{x}') = \sum_{s \in \mathcal{A}^*} w_s \phi_s(\mathbf{x}) \phi_s(\mathbf{x}')$$

where  $w_s \geq 0$  and  $\mathcal{A}^*$  is the set of all strings (of any length) from the alphabet  $\mathcal{A}$

- this is a Mercer kernel

# Outline

- 1 Intro
  - Why Kernels?
  - Kernel Characterization
- 2 Kernels Functions
  - RBF Kernels
  - Kernels for Comparing Documents
  - Mercer Kernels
- 3 Kernel-based Models
  - **Kernel Machines**
  - Kernel Trick
  - Kernelized 1NN Classification
  - Kernelized Ridge Regression
- 4 Support Vector Machine (SVM)
  - Loss Functions
  - SVM for Regression
  - SVM for Classification
  - Large Margin Principle

- we define a **kernel machine** to be a Generalized Linear Model (GLM) where the input feature vector has the form

$$\phi(\mathbf{x}) = [\kappa(\mathbf{x}, \mu_1), \dots, \kappa(\mathbf{x}, \mu_K)]$$

where  $\mu_k \in \chi$  are the set of  $K$  **centroids**

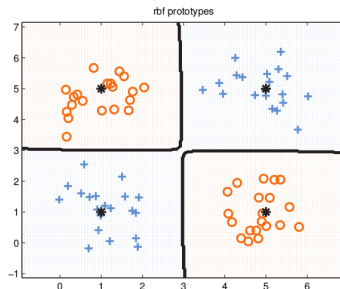
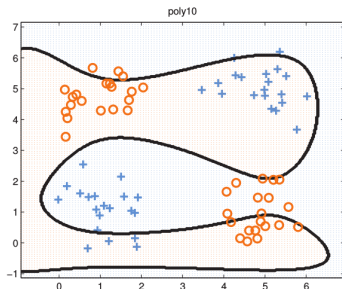
- open question: how to choose the centroids  $\mu_k$ ?
- the above vector  $\phi(\mathbf{x})$  is called **kernelized feature vector**
- kernel machines do not require that the kernel are Mercer
- if  $\kappa$  is an RBF kernel, the corresponding kernel machine is called a **RBF network**

# Kernel Machines

- we can use the **kernelized feature vector for logistic regression** by defining

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

this provides a simple way to define a **non-linear decision boundary**



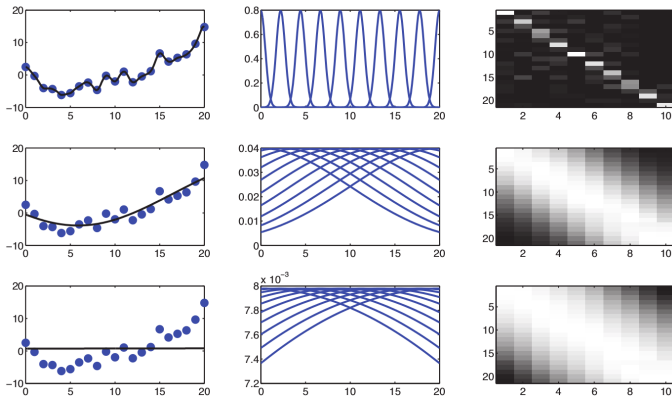
- left*: fitting a linear logistic regression classifier using degree 10 polynomial expansion.
- right*: same model, but using an RBF kernel with centroids specified by the 4 black crosses



- we can also use the **kernelized feature vector inside a linear regression model** by defining

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}), \sigma^2)$$

# Kernel Machines



- left column: fitted function where  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$  and  $y_i, x_i \in \mathbb{R}$
- middle column: RBF basis functions  $\kappa(x, \mu_i)$  evaluated on a grid ( $K = 10$  functions uniformly spaced)
- right column: design matrix
- top to bottom we show different bandwidths:  $\sigma = 0.1, \sigma = 0.5, \sigma = 50$

how to choose the centroids  $\mu_k$ ?

$$\phi(\mathbf{x}) = [\kappa(\mathbf{x}, \mu_1), \dots, \kappa(\mathbf{x}, \mu_K)]$$

- if  $D$  is small, a simple solution is to **uniformly tile**/grid the space occupied by the data (recall the curse of dimensionality)
- another approach is to **find clusters** in the data (but how to pick the number of clusters)
- a simpler approach is to make **each sample  $\mathbf{x}_i$  a prototype**

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

and use a sparse-promoting prior for  $\mathbf{w}$  to efficiently select subset of training exemplars  $\mathbf{x}_i$

- another very popular approach is the **Support Vector Machine (SVM)** which modify the likelihood term instead of using a sparsity-promoting prior

# Outline

- 1 Intro
  - Why Kernels?
  - Kernel Characterization
- 2 Kernels Functions
  - RBF Kernels
  - Kernels for Comparing Documents
  - Mercer Kernels
- 3 Kernel-based Models
  - Kernel Machines
  - **Kernel Trick**
  - Kernelized 1NN Classification
  - Kernelized Ridge Regression
- 4 Support Vector Machine (SVM)
  - Loss Functions
  - SVM for Regression
  - SVM for Classification
  - Large Margin Principle

- rather than defining our feature vector in terms of kernels

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

we can instead **work with the original feature vectors  $\mathbf{x}$** , but modify the algorithm so that it **replaces all inner products** of the form  $\langle \mathbf{x}, \mathbf{x}' \rangle = \mathbf{x}^T \mathbf{x}'$  with a call to the kernel function,  $\kappa(\mathbf{x}, \mathbf{x}')$

- this is called the **kernel trick**
- it turns out that many algorithms can be kernelized in this way
- the use of Mercer kernel is required for this trick to work

# Outline

- 1 Intro
  - Why Kernels?
  - Kernel Characterization
- 2 Kernels Functions
  - RBF Kernels
  - Kernels for Comparing Documents
  - Mercer Kernels
- 3 Kernel-based Models
  - Kernel Machines
  - Kernel Trick
  - Kernelized 1NN Classification
  - Kernelized Ridge Regression
- 4 Support Vector Machine (SVM)
  - Loss Functions
  - SVM for Regression
  - SVM for Classification
  - Large Margin Principle

# Kernelized 1NN Classification

- a 1NN classifier computes the Euclidean distance of a test vector to all the training points, find the closest one, and look up its label
- this can be kernelized by observing that

$$\|\mathbf{x} - \mathbf{x}'\|_2^2 = \langle \mathbf{x} - \mathbf{x}', \mathbf{x} - \mathbf{x}' \rangle = \langle \mathbf{x}, \mathbf{x} \rangle + \langle \mathbf{x}', \mathbf{x}' \rangle - 2\langle \mathbf{x}, \mathbf{x}' \rangle$$

- in this way we can redefine the distance by using the chosen kernel

# Outline

- 1 Intro
  - Why Kernels?
  - Kernel Characterization
- 2 Kernels Functions
  - RBF Kernels
  - Kernels for Comparing Documents
  - Mercer Kernels
- 3 Kernel-based Models
  - Kernel Machines
  - Kernel Trick
  - Kernelized 1NN Classification
  - Kernelized Ridge Regression
- 4 Support Vector Machine (SVM)
  - Loss Functions
  - SVM for Regression
  - SVM for Classification
  - Large Margin Principle



# Kernelized Ridge Regression

## The Primal Problem

- $\mathbf{x} \in \mathbb{R}^D$  and  $\mathbf{X} \in \mathbb{R}^{N \times D}$
- minimization problem

$$J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - (w_0 + \mathbf{w}^T \mathbf{x}_i))^2 + \lambda \|\mathbf{w}\|_2^2 = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

- the solution is

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^T \mathbf{y}$$

- but here we do not have yet inner products to replace with kernel functions

# Kernelized Ridge Regression

## The Dual Problem

- by using the matrix inversion lemma

$$(\mathbf{E} - \mathbf{F}\mathbf{H}^{-1}\mathbf{G})^{-1}\mathbf{F}\mathbf{H}^{-1} = \mathbf{E}^{-1}\mathbf{F}(\mathbf{H} - \mathbf{G}\mathbf{E}^{-1}\mathbf{F})^{-1}$$

and setting  $\mathbf{E} = \mathbf{I}_D$ ,  $\mathbf{H} = \mathbf{I}_N$ ,  $\mathbf{F} = -\mathbf{X}^T$  and  $\mathbf{G} = \mathbf{X}$  we can pass from

$$\mathbf{w} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_D)^{-1}\mathbf{X}^T\mathbf{y}$$

to

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

# Kernelized Ridge Regression

## The Dual Problem

- given

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

- we can replace  $\mathbf{X}\mathbf{X}^T$  with  $\mathbf{K}$  since  $k_{ij} = \mathbf{x}_i^T \mathbf{x}_j$
- we can define the **dual variables**

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

- hence, one has

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha} = \sum_{i=1}^N \alpha_i \mathbf{x}_i$$

the solution vector  $\mathbf{w}$  is just a linear sum of the  $N$  training vectors

- if we plug this in at test time to compute the predictive mean

$$\hat{f}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \sum_{i=1}^N \alpha_i \mathbf{x}_i^T \mathbf{x} = \sum_{i=1}^N \alpha_i \kappa(\mathbf{x}, \mathbf{x}_i)$$

- this kind of technique can be used to kernelize many other linear models such as logistic regression

# Outline

- 1 Intro
  - Why Kernels?
  - Kernel Characterization
- 2 Kernels Functions
  - RBF Kernels
  - Kernels for Comparing Documents
  - Mercer Kernels
- 3 Kernel-based Models
  - Kernel Machines
  - Kernel Trick
  - Kernelized 1NN Classification
  - Kernelized Ridge Regression
- 4 Support Vector Machine (SVM)
  - Loss Functions
  - SVM for Regression
  - SVM for Classification
  - Large Margin Principle

- consider the regularized **empirical risk function**

$$J(\mathbf{w}, \lambda) = \sum_{i=1}^N L(y_i, \hat{y}_i) + \lambda \|\mathbf{w}\|^2$$

where  $\hat{y}_i = \mathbf{w}^T \mathbf{x}_i + w_0$

- if  $L(y_i, \hat{y}_i) = (y_i - \hat{y}_i)^2$  we have a **quadratic loss** and the problem becomes a **ridge regression**
- if  $L(y_i, \hat{y}_i) = -\log p(y_i | \mathbf{x}_i, \mathbf{w}_i) = -\log(\text{sigm}(y_i \eta_i)) = \log(1 + e^{-y_i \eta_i})$  with  $y_i \in \{-1, 1\}$  and  $\eta_i = \mathbf{w}^T \mathbf{x}_i + w_0$ , we have a **log-loss** and the problem becomes a **logistic regression**
- we have seen how to kernelize a model but we want also a sparse solution for efficiency reasons
- if we replace the quadratic/log-loss with some other loss function, we can ensure that the **solution is sparse**, so that predictions only depend on a subset of the training data, known as **support vectors**
- this combination of the kernel trick plus a modified loss function is known as a **support vector machine** or SVM

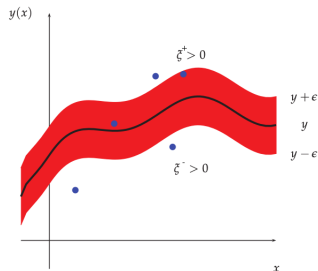
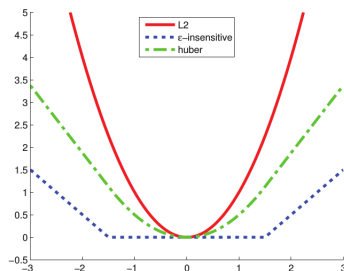
# Outline

- 1 Intro
  - Why Kernels?
  - Kernel Characterization
- 2 Kernels Functions
  - RBF Kernels
  - Kernels for Comparing Documents
  - Mercer Kernels
- 3 Kernel-based Models
  - Kernel Machines
  - Kernel Trick
  - Kernelized 1NN Classification
  - Kernelized Ridge Regression
- 4 Support Vector Machine (SVM)
  - Loss Functions
  - SVM for Regression
  - SVM for Classification
  - Large Margin Principle

- let's consider a regression problem  $J(\mathbf{w}, \lambda) = \sum_{i=1}^N L(y_i, \hat{y}_i) + \lambda \|\mathbf{w}\|^2$  with  $\hat{y}_i = \mathbf{w}^T \mathbf{x} + w_0$
- if we use the **epsilon insensitive loss function**

$$L_{\epsilon}(y, \hat{y}) \triangleq \begin{cases} 0 & \text{if } \|y - \hat{y}\| < \epsilon \\ \|y - \hat{y}\| - \epsilon & \text{otherwise} \end{cases}$$

it means that any point lying inside an  $\epsilon$ -tube around the prediction is not penalized

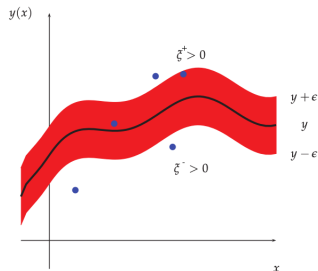
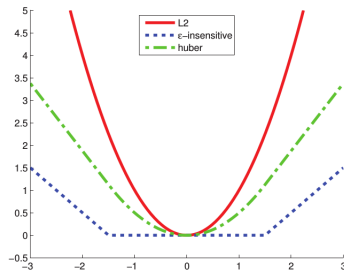


- the objective function is usually written as

$$J = C \sum_{i=1}^N L_{\epsilon}(y_i, \hat{y}_i) + \frac{1}{2} \|\mathbf{w}\|^2$$

with  $\hat{y}_i = f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + w_0$  and  $C = \frac{1}{\lambda}$  is regularization constant

- this objective function is **convex** and **unconstrained** but **not differentiable**





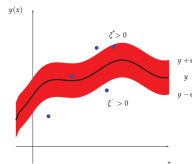
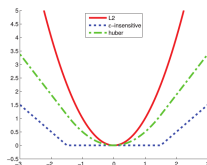
- one popular approach is to formulate the problem as a **constrained optimization problem**
- in particular, we can introduce the **slack variables**  $\xi_i^+ \geq 0$  and  $\xi_i^- \geq 0$  to represent the degree to which each point  $\mathbf{x}_i$  lies outside the tube

$$y_i \leq f(\mathbf{x}_i) + \epsilon + \xi_i^+$$

$$y_i \geq f(\mathbf{x}_i) - \epsilon - \xi_i^-$$

- the problem can be restated as a standard quadratic program in  $2N + D + 1$  variables

$$J = C \sum_{i=1}^N (\xi_i^+ + \xi_i^-) + \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t. } \xi_i^+ \geq 0 \text{ and } \xi_i^- \geq 0$$



- the problem

$$J = C \sum_{i=1}^N (\xi_i^+ + \xi_i^-) + \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t. } \xi_i^+ \geq 0 \text{ and } \xi_i^- \geq 0$$

- it is possible to show that the optimal solution has the form

$$\hat{\mathbf{w}} = \sum_i \alpha_i \mathbf{x}_i \quad \text{with } \alpha_i \geq 0$$

- it turns out that in the solution the vector  $\alpha$  is **sparse**, because we don't care about errors which are smaller than  $\epsilon$
- the  $\mathbf{x}_i$  for which  $\alpha_i > 0$  are called the **support vectors**: these are points for which the errors lie on or outside the tube
- all other vectors can be neglected when computing  $\hat{\mathbf{w}}$

- at test time we evaluate the  $y$  function as

$$\hat{y}(\mathbf{x}) = \hat{w}_0 + \hat{\mathbf{w}}^T \mathbf{x}$$

- once we plug in the definition of  $\hat{\mathbf{w}}$  we have

$$\hat{y}(\mathbf{x}) = \hat{w}_0 + \sum_i \alpha_i \mathbf{x}_i^T \mathbf{x}$$

and we can **kernelize** it by replacing the inner product with the kernel function

$$\hat{y}(\mathbf{x}) = \hat{w}_0 + \sum_i \alpha_i \kappa(\mathbf{x}, \mathbf{x}_i)$$

# Outline

- 1 Intro
  - Why Kernels?
  - Kernel Characterization
- 2 Kernels Functions
  - RBF Kernels
  - Kernels for Comparing Documents
  - Mercer Kernels
- 3 Kernel-based Models
  - Kernel Machines
  - Kernel Trick
  - Kernelized 1NN Classification
  - Kernelized Ridge Regression
- 4 Support Vector Machine (SVM)
  - Loss Functions
  - SVM for Regression
  - SVM for Classification
  - Large Margin Principle

- let's consider a logistic regression problem where  $y_i \in \{-1, +1\}$
- the objective function is

$$J(\mathbf{w}, \lambda) = \sum_{i=1}^N L(y_i, \hat{y}_i) + \lambda \|\mathbf{w}\|^2$$

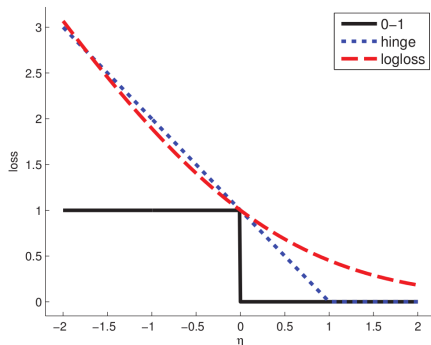
with a **log-loss**  $L(y_i, \eta_i) = -\log p(y_i | \mathbf{x}_i, \mathbf{w}_i) = -\log(\text{sigm}(y_i \eta_i)) = \log(1 + e^{-y_i \eta_i})$   
where  $\eta_i = f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + w_0$  represents our “confidence” in choosing label  $\hat{y}_i = 1$

- in principle, we could use the **0-1 loss**  $L(y_i, \eta_i) = \mathbb{I}(y_i \neq \eta_i) = \mathbb{I}(y_i \eta_i < 0)$  so as to minimize the **misclassification error**
- unfortunately, the 0-1 risk is a very non-smooth objective and hence is hard to optimize
- the SVM algorithm replaces the log-loss with the **hinge loss**

$$L_{\text{hinge}}(y_i, \eta_i) \triangleq \max(0, 1 - y_i \eta_i) = (1 - y_i \eta_i)_+$$

where  $(v)_+ \triangleq \max(0, v)$

- $\eta = f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$  represents our “confidence” in choosing label  $y = 1$
- **log-loss**  $L(y, \eta) = -\log p(y|\mathbf{x}, \mathbf{w}) = \log(1 + e^{-y\eta})$
- **0-1 loss**  $L(y, \eta) = \mathbb{I}(y \neq \eta) = \mathbb{I}(y\eta < 0)$
- **hinge loss**  $L_{\text{hinge}}(y, \eta) \triangleq \max(0, 1 - y\eta) = (1 - y\eta)_+$
- the hinge loss and log-loss represent smooth convex upper bounds on the 0-1 loss
- in the figure below, the horizontal axis is the margin  $y\eta$ , the vertical axis is the loss



- the objective function is

$$J = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N (1 - y_i \eta_i)_+$$

with  $\eta_i = f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + w_0$

- since the function  $(1 - y_i \eta_i)_+$  is not differentiable we can introduce the **slack variable**  $\xi_i \geq 0$  and rewrite the objective as

$$J = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \quad \text{s.t.} \quad \xi_i \geq 0, \quad y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1 - \xi, \quad i = 1 : N$$

- the objective

$$J = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \quad \text{s.t.} \quad \xi_i \geq 0, \quad y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1 - \xi_i, \quad i = 1 : N$$

- it is possible to show that the **optimal solution** has the form

$$\hat{\mathbf{w}} = \sum_i \alpha_i \mathbf{x}_i \quad \text{with} \quad \alpha_i = \lambda_i y_i$$

and where  $\alpha$  is sparse (because of the hinge loss)

- the  $\mathbf{x}_i$  for which  $\alpha_i > 0$  are called **support vectors**: these are points which are either incorrectly classified, or are classified correctly but are on or inside the margin



- at test time we evaluate the  $y$  function as

$$\hat{y}(\mathbf{x}) = \text{sign}(f(\mathbf{x})) = \text{sign}(\hat{w}_0 + \hat{\mathbf{w}}^T \mathbf{x})$$

- once we plug in the definition of  $\hat{\mathbf{w}}$  we have

$$\hat{y}(\mathbf{x}) = \text{sign}(\hat{w}_0 + \sum_{i=1}^N \alpha_i \mathbf{x}_i^T \mathbf{x})$$

and we can **kernelize** by replacing the inner product with the kernel function

$$\hat{y}(\mathbf{x}) = \text{sign}(\hat{w}_0 + \sum_{i=1}^N \alpha_i \kappa(\mathbf{x}, \mathbf{x}_i))$$

# Outline

- 1 Intro
  - Why Kernels?
  - Kernel Characterization
- 2 Kernels Functions
  - RBF Kernels
  - Kernels for Comparing Documents
  - Mercer Kernels
- 3 Kernel-based Models
  - Kernel Machines
  - Kernel Trick
  - Kernelized 1NN Classification
  - Kernelized Ridge Regression
- 4 Support Vector Machine (SVM)
  - Loss Functions
  - SVM for Regression
  - SVM for Classification
  - Large Margin Principle

# SVM Classification

## Large Margin Principle

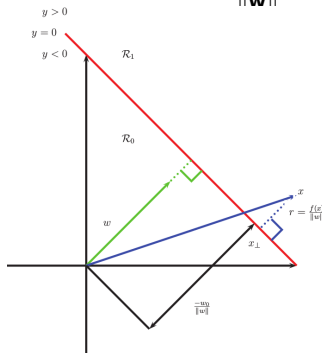
- now let's revise the previous concepts under a geometrical point of view
- $f(\mathbf{x})$  is the **discriminant function** which will be **linear** in the **feature space implied by the choice of the kernel**
- $f(\mathbf{x}) = 0$  is the **decision boundary**
- now, for simplicity, let's assume that  $\mathbf{x}$  belongs to the kernel induced space  $\phi(\chi)$

# SVM Classification

## Large Margin Principle

- for each point we have  $\mathbf{x} = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$  where  $r$  is the distance of  $\mathbf{x}$  from the decision boundary  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = 0$  (which is an hyperplane whose normal vector is  $\mathbf{w}$ ), and  $\mathbf{x}_\perp$  is the orthogonal projection of  $\mathbf{x}$  onto this boundary (hence  $\mathbf{w}^T \mathbf{x}_\perp + w_0 = 0$ )
- if we plug the decomposition of  $\mathbf{x}$  in  $f(\mathbf{x})$ , we have

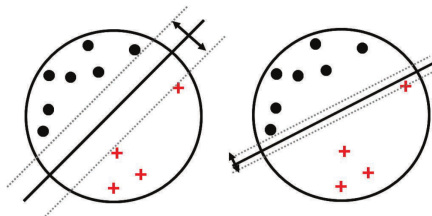
$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = (\mathbf{w}^T \mathbf{x}_\perp + w_0) + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} = r \|\mathbf{w}\| \implies r = \frac{f(\mathbf{x})}{\|\mathbf{w}\|}$$



# SVM Classification

## Large Margin Principle

- for each point  $\mathbf{x}_i$  we would like to make this perpendicular distance  $r = \frac{f(\mathbf{x}_i)}{\|\mathbf{w}\|}$  as large as possible
- in particular, there might be many lines that perfectly separate the training data
- the best one to pick is the one that **maximizes the margin**, i.e., the **perpendicular distance to the closest point**
- in addition, we want to ensure each point is on the correct side of the boundary, hence we want  $f(\mathbf{x}_i)y_i > 0$



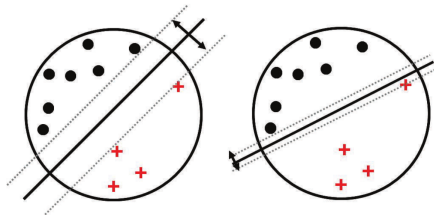
# SVM Classification

## Large Margin Principle

- our objective becomes

$$\max_{\mathbf{w}, w_0} \min_{i=1}^N \frac{y_i(\mathbf{w}^T \mathbf{x}_i + w_0)}{\|\mathbf{w}\|}$$

- rescaling the parameters using  $\mathbf{w} \rightarrow k\mathbf{w}$  and  $w_0 \rightarrow kw_0$ , we do not change the distance of any point to the boundary, since the  $k$  factor cancels out when we divide by  $\|\mathbf{w}\|$



# SVM Classification

## Large Margin Principle

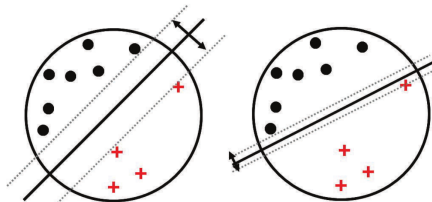
- our objective is

$$\max_{\mathbf{w}, w_0} \min_{i=1}^N \frac{y_i(\mathbf{w}^T \mathbf{x}_i + w_0)}{\|\mathbf{w}\|}$$

- therefore let us define the scale factor such that  $y_i f_i = 1$  for the point that is **closest** to the decision boundary
- note that maximizing  $1/\|\mathbf{w}\|$  is equivalent to minimizing  $\|\mathbf{w}\|^2$
- thus we get the new objective

$$\min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1, \quad i = 1 : N$$

- the constraint says that we want all points to be on the correct side of the decision boundary with a margin of at least 1

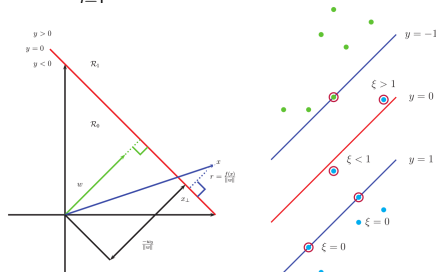


# SVM Classification

## Large Margin Principle

- if the data is **not linearly separable** (even after using the kernel trick), there will be no feasible solution in which  $y_i f_i \geq 1$  for all  $i$
- we therefore introduce **slack variables**  $\xi \geq 0$  such that  $\xi_i = 0$  if the point is on or inside the correct margin boundary, and  $\xi_i = |y_i - f_i|$  otherwise
- we replace the **hard constraints** that  $y_i f_i \geq 1$  with the **soft margin constraints** that  $y_i f_i \geq 1 - \xi_i$
- the new objective becomes

$$\min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \quad \text{s.t.} \quad \xi_i \geq 0, \quad y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1 - \xi_i, \quad i = 1 : N$$





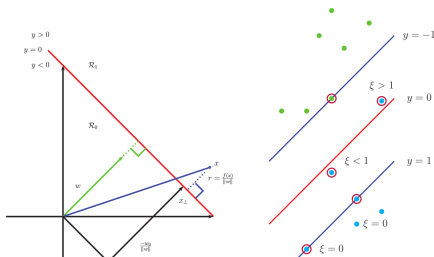
# SVM Classification

## Large Margin Principle

- the new objective

$$\min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \quad \text{s.t.} \quad \xi_i \geq 0, \quad y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1 - \xi_i, \quad i = 1 : N$$

- since  $\xi_i > 1$  means point  $i$  is misclassified, we can interpret  $\sum_i \xi_i$  as an upper bound on the **number of misclassified points**
- the parameter  $C$  is a regularization parameter that controls the number of errors we are willing to tolerate on the training set
- it is common to define this using  $C = 1/(\nu N)$  where  $0 < \nu \leq 1$  controls the fraction of misclassified points that we allow during the training phase



- Kevin Murphy's book