

Lecture 9: Support Vector Machines II

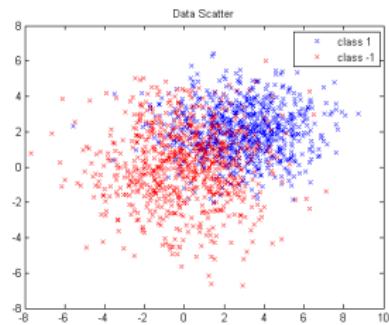
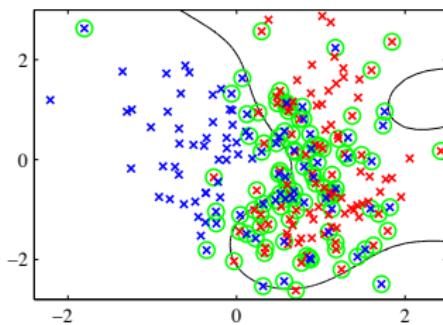
Reading: Section 12.3

GU4241/GR5241 Statistical Machine Learning

Linxi Liu
February 16, 2017

Motivation

More realistic data



Motivation: Kernels

sign $\langle V_h, X_i \rangle - c$, V_h is a p dimensional vector

This is a separating hyperplane for linear problem

$$X_i = (X_{i1}, \dots, X_{ip})$$

$$\sum V^h \cdot X_i \cdot J$$

$$Y = \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

Idea

- ▶ The SVM uses the scalar product $\langle x, \tilde{x}_i \rangle$ as a measure of similarity between x and \tilde{x}_i , and of distance to the hyperplane.
- ▶ Since the scalar product is linear, the SVM is a linear method.
- ▶ By using a *nonlinear* function instead, we can make the classifier nonlinear.

primal problem

$$\text{min: } \|V_h\|$$

$$\text{st } Y_i(\langle V_h, X_i \rangle - c) \geq 1$$

$$\text{max: } W(x) = \sum X_i - \sum a_i a_j y_i y_j \langle X_i, X_j \rangle$$

inner product us a function $\langle X_i, X_j \rangle$ for linear,
for non-linear: we use kernel

inner product is still linear classifier

Kernels in Detail

kernel will take all the linear features, and
take all the non-linear features to higher dimension
kernel help us map the non-linear data from low d to high d

- ▶ Scalar product can be regarded as a two-argument function

$$\langle \cdot, \cdot \rangle : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$

$$x_i \rightarrow (\Phi_1(x_i), \Phi_2(x_i))$$

- ▶ We will replace this function with a function $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ and substitute

$$k(\mathbf{x}, \mathbf{x}') \quad \text{for every occurrence of} \quad \langle \mathbf{x}, \mathbf{x}' \rangle$$

in the SVM formula.

- ▶ Under certain conditions on k , all optimization/classification results for the SVM still hold. Functions that satisfy these conditions are called **kernel functions**.

The Most Popular Kernel

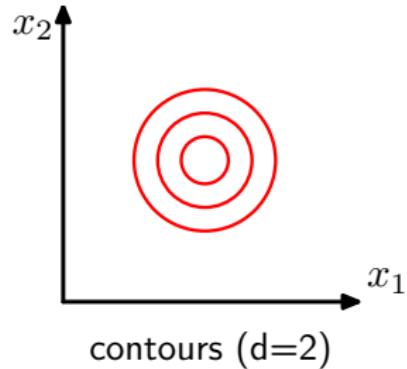
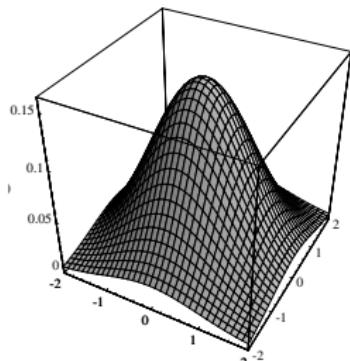
RBF Kernel

$$k_{\text{RBF}}(\mathbf{x}, \mathbf{x}') := \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2}\right) \quad \text{for some } \sigma \in \mathbb{R}_+$$

is called an **RBF kernel** (RBF = radial basis function). The parameter σ is called **bandwidth**.

Other names for k_{RBF} : Gaussian kernel, squared-exponential kernel.

If we fix \mathbf{x}' , the function $k_{\text{RBF}}(., \mathbf{x}')$ is (up to scaling) a spherical Gaussian density on \mathbb{R}^d , with mean \mathbf{x}' and standard deviation σ .



Choosing a kernel

Theory

To define a kernel:

- ▶ We have to define a function of two arguments and prove that it is a kernel.
- ▶ This is done by checking a set of necessary and sufficient conditions known as "Mercer's theorem".

Practice

The data analyst does not define a kernel, but tries some well-known standard kernels until one seems to work. Most common choices:

- ▶ The RBF kernel.
- ▶ The "linear kernel" $k_{SP}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle$, i.e. the standard, linear SVM.

Once kernel is chosen

- ▶ Classifier can be trained by solving the optimization problem using standard software.
- ▶ SVM software packages include implementations of most common kernels.

Which Functions work as Kernels?

Formal definition

A function $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is called a **kernel** on \mathbb{R}^d if there is *some* function $\phi : \mathbb{R}^d \rightarrow \mathcal{F}$ into *some* space \mathcal{F} with scalar product $\langle \cdot, \cdot \rangle_{\mathcal{F}}$ such that

$$k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{F}} \quad \text{for all } \mathbf{x}, \mathbf{x}' \in \mathbb{R}^d.$$

In other words

- ▶ k is a kernel if it can be interpreted as a scalar product on some other space.
- ▶ If we substitute $k(\mathbf{x}, \mathbf{x}')$ for $\langle \mathbf{x}, \mathbf{x}' \rangle$ in all SVM equations, we implicitly train a *linear* SVM on the space \mathcal{F} .
- ▶ The SVM still works: It still uses scalar products, just on another space.

The mapping ϕ

- ▶ ϕ has to transform the data into data on which a linear SVM works well.
- ▶ This is usually achieved by choosing \mathcal{F} as a higher-dimensional space than \mathbb{R}^d .

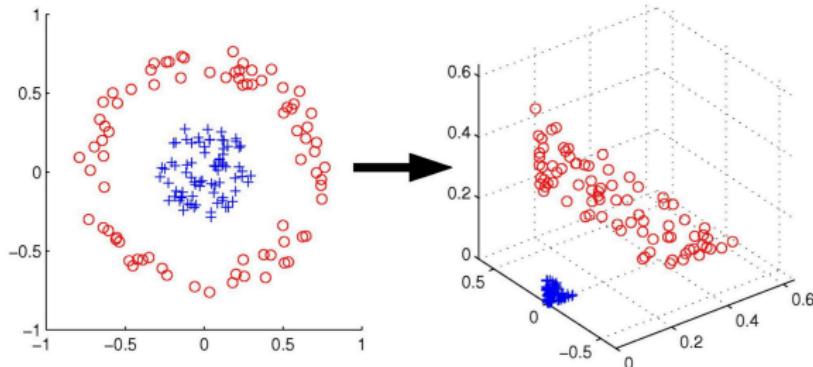
Mapping into Higher Dimensions

Example

How can a map into higher dimensions make class boundary (more) linear?

Consider

$$\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3 \quad \text{where} \quad \phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} := \begin{pmatrix} x_1^2 \\ 2x_1x_2 \\ x_2^2 \end{pmatrix}$$



Mapping into Higher Dimensions

Problem

In previous example: We have to know what the data looks like to choose ϕ !

Solution

- ▶ Choose high dimension h for \mathcal{F} .
- ▶ Choose components ϕ_i of $\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_h(\mathbf{x}))$ as different nonlinear mappings.
- ▶ If two points differ in \mathbb{R}^d , some of the nonlinear mappings will amplify differences.

The RBF kernel is an extreme case

- ▶ The function k_{RBF} can be shown to be a kernel, however:
- ▶ \mathcal{F} is infinite-dimensional for this kernel.

Determining whether k is a kernel

Mercer's theorem $\mathbf{x}_i \rightarrow (\Phi_1(\mathbf{x}), \Phi_2(\mathbf{x}))$

A mathematical result called *Mercer's theorem* states that, if the function k is positive, i.e.

$$\int_{\mathbb{R}^d \times \mathbb{R}^d} k(\mathbf{x}, \mathbf{x}') f(\mathbf{x}) f(\mathbf{x}') d\mathbf{x} d\mathbf{x}' \geq 0$$

all the kernel can be decomposed to lower dim

for all functions f , then it can be written as

$$k(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^{\infty} \lambda_j \phi_j(\mathbf{x}) \phi_j(\mathbf{x}')$$

The ϕ_j are functions $\mathbb{R}^d \rightarrow \mathbb{R}$ and $\lambda_i \geq 0$. This means the (possibly infinite) vector $\phi(\mathbf{x}) = (\sqrt{\lambda_1} \phi_1(\mathbf{x}), \sqrt{\lambda_2} \phi_2(\mathbf{x}), \dots)$ is a feature map.

Kernel arithmetic

Various functions of kernels are again kernels: If k_1 and k_2 are kernels, then e.g.

$$k_1 + k_2$$

$$k_1 \cdot k_2$$

$$\text{const.} \cdot k_1$$

are again kernels.

The Kernel Trick

Kernels in general

- ▶ Many linear machine learning and statistics algorithms can be "kernelized".
- ▶ The only conditions are:
 1. The algorithm uses a scalar product.
 2. In all relevant equations, the data (and all other elements of \mathbb{R}^d) appear *only inside a scalar product*.
- ▶ This approach to making algorithms non-linear is known as the "kernel trick".

primal problem: fixed by adding kernel, replace the inner product by kernel

min: $\|Vh\|$

st $Y_i(\langle Vh, X_i \rangle - c) \geq 1$

max: $W(x) = \sum X_i - \sum a_i^* a_j^* y_i^* y_j^* K \langle X_i, X_j \rangle$

Kernel SVM

Optimization problem

normally we use RBF kernel,
or polynomial kernel, which is easier
or gaussian kernel

$$\begin{aligned} \min_{\mathbf{v}_H, c} \quad & \|\mathbf{v}_H\|_{\mathcal{F}}^2 + \gamma \sum_{i=1}^n \xi_i^2 \\ \text{s.t.} \quad & \tilde{y}_i (\langle \mathbf{v}_H, \phi(\tilde{\mathbf{x}}_i) \rangle_{\mathcal{F}} - c) \geq 1 - \xi_i \quad \text{and } \xi_i \geq 0 \end{aligned}$$

Note: \mathbf{v}_H now lives in \mathcal{F} , and $\|\cdot\|_{\mathcal{F}}$ and $\langle \cdot, \cdot \rangle_{\mathcal{F}}$ are norm and scalar product on \mathcal{F} .

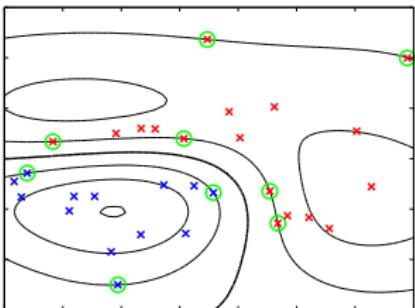
Dual optimization problem

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^n} \quad & W(\alpha) := \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j \tilde{y}_i \tilde{y}_j (\mathbf{k}(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j) + \frac{1}{\gamma} \mathbb{I}\{i = j\}) \\ \text{s.t.} \quad & \sum_{i=1}^n \tilde{y}_i \alpha_i = 0 \quad \text{and} \quad \alpha_i \geq 0 \end{aligned}$$

Classifier

$$f(\mathbf{x}) = \operatorname{sgn} \left(\sum_{i=1}^n \tilde{y}_i \alpha_i^* \mathbf{k}(\tilde{\mathbf{x}}_i, \mathbf{x}) - c \right)$$

SVM with RBF Kernel



we have the non-linear decision boundaries

$$f(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^n y_i \alpha_i^* k_{\text{RBF}}(\mathbf{x}_i, \mathbf{x}) \right)$$

this is the final classifier,
which is non-linear

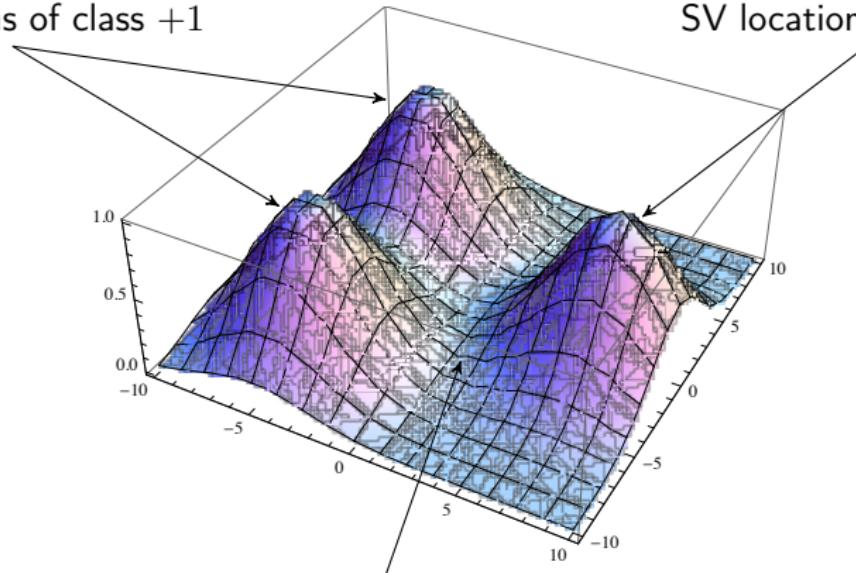
- ▶ Circled points are support vectors. The two contour lines running through support vectors are the nonlinear counterparts of the convex hulls.
- ▶ The thick black line is the classifier.
- ▶ Think of a Gaussian-shaped function $k_{\text{RBF}}(\cdot, \mathbf{x}')$ centered at each support vector \mathbf{x}' . These functions add up to a function surface over \mathbb{R}^2 .
- ▶ The lines in the image are contour lines of this surface. The classifier runs along the bottom of the "valley" between the two classes.
- ▶ Smoothness of the contours is controlled by σ

Decision Boundary with RBF Kernel

Text

SV locations of class +1

SV location of class -1



The decision boundary runs here.

The decision boundary of the classifier coincides with the set of points where the surfaces for class +1 and class -1 have equal value.

Summary: SVMs

Basic SVM

- ▶ Linear classifier for linearly separable data.
- ▶ Positions of affine hyperplane is determined by maximizing margin.
- ▶ Maximizing the margin is a convex optimization problem.

Full-fledged SVM

Ingredient	Purpose
Maximum margin	Good generalization properties
Slack variables	Overlapping classes
Kernel	Robustness against outliers Nonlinear decision boundary

Use in practice

- ▶ Software packages (e.g. libsvm, SVMLite)
- ▶ Choose a kernel function (e.g. RBF)
- ▶ Cross-validate margin parameter γ and kernel parameters (e.g. bandwidth)