

Machine Learning

MIRI Master

Lluís A. Belanche

belanche@cs.upc.edu



Soft Computing Research Group
Dept. de Ciències de la Computació (Computer Science)
Universitat Politècnica de Catalunya

Spring Semester 2016-2017

LECTURE 8: Learning with kernels (I) - The SVM

Learning with kernels (I): The SVM

Linear regression revisited

Problem: We wish to find a function $y(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + w_0$ which best models a data set $D = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\} \subset \mathbb{R}^d \times \mathbb{R}$

- We define $\mathbf{w} = (w_0, w_1, \dots, w_d)^\top$ and $\mathbf{x} = (1, x_1, \dots, x_d)^\top$
- Then we minimize the regularized (aka penalized) empirical error:

$$E_{\text{emp}}^\lambda(y) = \sum_{n=1}^N \left(t_n - y(\mathbf{x}_n)\right)^2 + \lambda \sum_{i=0}^d w_i^2 = \|\mathbf{t} - X\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2$$

The parameter $\lambda > 0$ defines a trade-off between the fit to the data and the complexity of the vector \mathbf{w}

Learning with kernels (I): The SVM

Linear regression revisited

Setting $\frac{\partial E_{\text{emp}}^\lambda(y)}{\partial \mathbf{w}} = 0$, we obtain the (regularized) normal equations:

$$-2X^\top(t - X\mathbf{w}) + 2\lambda\mathbf{w} = 0$$

with solution

$$\hat{\mathbf{w}} = (X^\top X + \lambda I_d)^{-1} X^\top t$$

and therefore

$$y(\mathbf{x}) = \hat{\mathbf{w}}^\top \mathbf{x}$$

Learning with kernels (I): The SVM

Linear regression revisited

It turns out that the regularized solution can be written as:

$$\hat{\mathbf{w}} = \sum_{n=1}^N \alpha_n \mathbf{x}_n, \quad \hat{\mathbf{w}} = \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \\ \vdots \\ \hat{w}_d \end{pmatrix}$$

$$y(\mathbf{x}) = \sum_{n=1}^N \alpha_n (\mathbf{x}_n^\top \mathbf{x}), \quad \boldsymbol{\alpha} = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{pmatrix}$$

The new vector of parameters is given by $\boldsymbol{\alpha} = (X X^\top + \lambda I_N)^{-1} \mathbf{t}$

Learning with kernels (I): The SVM

Linear regression revisited

So we have the **primal** and the **dual** forms for $y(\mathbf{x})$:

$$y(\mathbf{x}) = \hat{\mathbf{w}}^\top \mathbf{x} \quad \text{and} \quad y(\mathbf{x}) = \sum_{n=1}^N \alpha_n (\mathbf{x}_n^\top \mathbf{x})$$

The dual form is usually more convenient when $d \gg N$:

- the primal requires the computation & inversion of $X^\top X + \lambda I_d$, requiring $O(Nd^2 + d^3)$ operations
- the dual requires the computation & inversion of $XX^\top + \lambda I_N$, requiring $O(dN^2 + N^3)$ operations

Learning with kernels (I): The SVM

Key aspects of kernel methods

How can we achieve non-linear regression?

A **feature map** is a function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^M$:

$$\phi(\mathbf{x}) = \left(\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_M(\mathbf{x}) \right)^\top$$

- $\phi(\mathbf{x})$ is called the **feature vector**
- $\{\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^d\}$ is the **feature space** (FS), and typically $M \gg d$.

Learning with kernels (I): The SVM

Key aspects of kernel methods

- Define $\Phi_{N \times M}$ the matrix of the $\phi(\mathbf{x}_n)$ as

$$\phi_{nm} = \phi_m(\mathbf{x}_n), n = 1, \dots, N, m = 1, \dots, M.$$

- Suppose we perform ridge regression on the Φ matrix
- The new regression function has the **primal** representation:

$$y(\mathbf{x}) = \hat{\mathbf{w}}^\top \phi(\mathbf{x})$$

Note the primal now operates in feature space

Learning with kernels (I): The SVM

Key aspects of kernel methods

Given a feature map $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^M$, we define its associated **kernel function** $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ as:

$$k(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u})^\top \phi(\mathbf{v}), \quad \mathbf{u}, \mathbf{v} \in \mathbb{R}^d$$

- The (feature) space where k implicitly operates is \mathbb{R}^M
- For some feature maps, computing $k(\mathbf{u}, \mathbf{v})$ is independent of M (it depends only on d).

Learning with kernels (I): The SVM

Key aspects of kernel methods

Since $\hat{\mathbf{w}} = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)$, the new regression function has the **dual** representation:

$$y(\mathbf{x}) = \sum_{n=1}^N \alpha_n (\phi(\mathbf{x}_n)^\top \phi(\mathbf{x})) = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x})$$

The new vector of parameters is given by

$$\boldsymbol{\alpha} = (\mathbf{K} + \lambda I_N)^{-1} \mathbf{t}, \quad \text{where } \mathbf{K} = (k(\mathbf{x}_n, \mathbf{x}_m))$$

Kernel-Based Learning

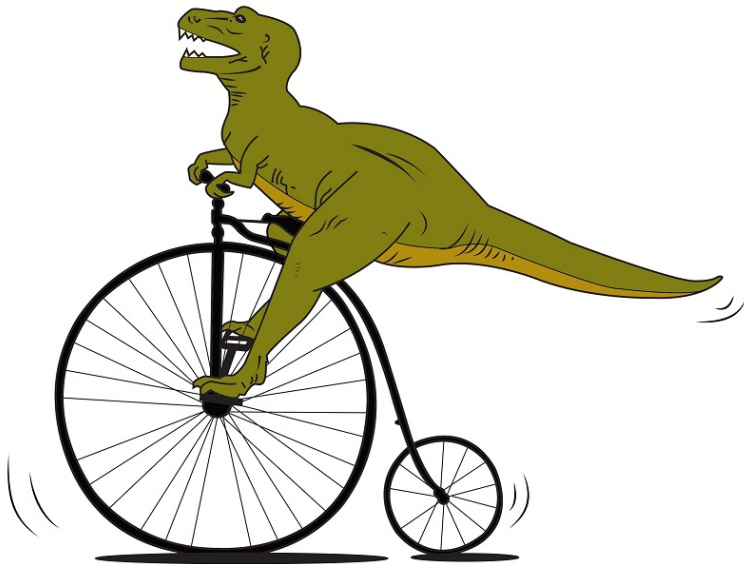
Example

What if we take the (simplest) choice $\phi(x) = x$?

\Rightarrow in this case $d = M$ and $k(u, v) = \langle u, v \rangle$

Then $\mathbf{K} = \mathbf{X}\mathbf{X}^\top$ and we are back to where we started

We thus have a “generalized linear model” (in a different sense than with GLMs)



Is this efficient or is it as out of date as a dinosaur riding a penny farthing?

Kernel-Based Learning

Key aspects of kernel methods

- A feature map is of the general form $\phi : \mathcal{X} \rightarrow \mathcal{H}$. The associated kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is $k(u, v) = \langle \phi(u), \phi(v) \rangle_{\mathcal{H}}$, $u, v \in \mathcal{X}$
- \mathcal{X} can be any space, \mathcal{H} is any **Hilbert space**:
 - An abstract complete vector space possessing the structure of an inner product
 - Examples would be \mathbb{R}^M or the l_2 space of square-summable sequences

In our previous discussion, $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{H} = \mathbb{R}^M$

Learning with kernels (I): The SVM

Key aspects of kernel methods

Kernel-based methods consist of two ingredients:

1. The kernel function (this is non-trivial)
2. The algorithm taking kernels as input
 - Data items are embedded into a vector space (feature space FS)
 - Linear relations are sought among the elements of the FS
 - The coordinates of these images are not needed: only their pairwise inner products
 - These inner products can sometimes be computed efficiently and implicitly in the input space (kernel function)
 - The solution vector is expressed as a linear combination of the kernel centered at the data

Learning with kernels (I): The SVM

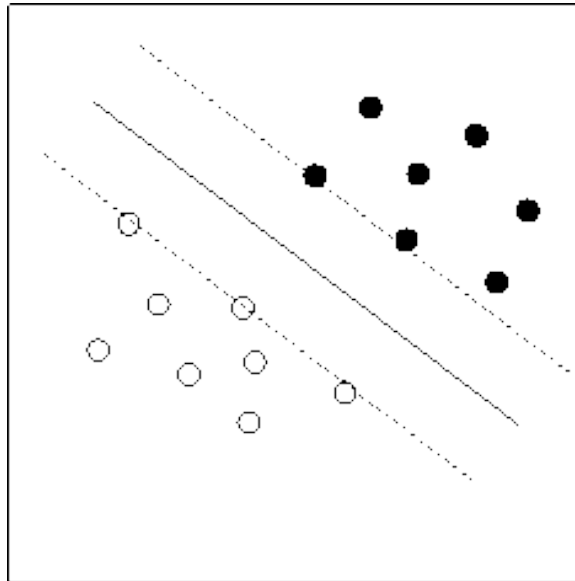
Key aspects of kernel methods

Many (classical and new) learning algorithms can be “kernelized”:

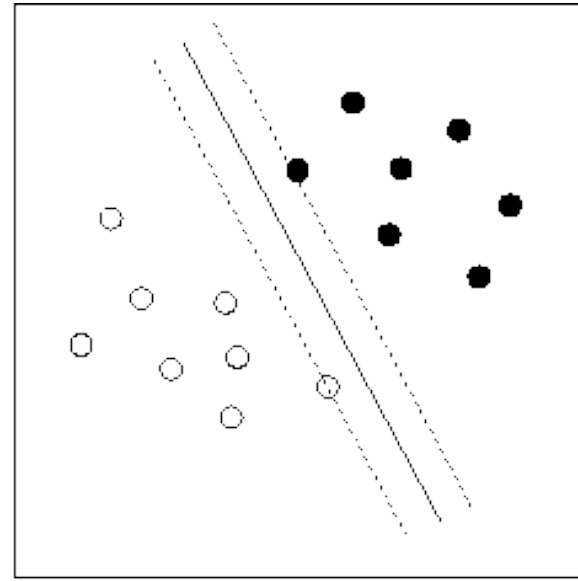
- The Support Vector Machine (SVM) and the Relevance Vector Machine (RVM)
- Fisher Discriminant Analysis (KFDA), Principal Components Analysis (KPCA), Canonical Correlation Analysis (KCCA), ...
- Kernel (regularized) linear regression
- Kernel k-means, kernel kNN
- (less known or very recent): PLS, Parzen Windows, logistic regression, statistical tests, ...

Support Vector Machines

Preliminaries



(a) Larger margin



(b) Smaller margin

Which solution is more likely to lead to better generalization?

Support Vector Machines

Preliminaries

Working Hypothesis (intuition):

The larger the margin, the better the generalization

- Criterion for building a two-class classifier:

Maximize the **width of the margin** between the classes

- **margin** = empty area around the decision boundary, defined by the distance to the nearest training examples

These examples will be called the **support vectors**

Goal: find the hyperplane (linear boundary) with the **largest margin**

Support Vector Machines

Formalisation

We have a data set $D = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$, with $\mathbf{x}_n \in \mathbb{R}^d$ and $t_n \in \{-1, +1\}$, describing a two class problem

We wish to find a function $y(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + w_0$ which best models D

- We would like to find \mathbf{w}, w_0 such that:

$$\langle \mathbf{w}, \mathbf{x}_n \rangle + w_0 > 0, \text{ when } t_n = +1$$

$$\langle \mathbf{w}, \mathbf{x}_n \rangle + w_0 < 0, \text{ when } t_n = -1$$

- In short, $t_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + w_0) > 0$, or $t_n y(\mathbf{x}_n) > 0$, $1 \leq n \leq N$

Support Vector Machines

Formalisation

- The quantity $t_n y(\mathbf{x}_n)$ is the **functional** margin of \mathbf{x}_n (there will be an “error” whenever $t_n y(\mathbf{x}_n) < 0$)
- Define the **loss** $L(t_n, \langle \mathbf{w}, \mathbf{x}_n \rangle) = \max(1 - t_n y(\mathbf{x}_n), 1)$
- Given the plane $\pi : y(\mathbf{x}) = 0$ or $\pi : \langle \mathbf{w}, \mathbf{x} \rangle + w_0 = 0$, the distance $d(\mathbf{x}, \pi) = \frac{|y(\mathbf{x})|}{\|\mathbf{w}\|}$ is called the **geometrical margin** of \mathbf{x}
- The optimal separating hyperplane (OSH) is the one that maximizes the geometrical margin for linearly separable data:

$$\max_{\mathbf{w}, b} \left\{ \min_{1 \leq n \leq N} d(\mathbf{x}_n, \pi) \right\} \quad \text{subject to } t_n (\langle \mathbf{w}, \mathbf{x}_n \rangle + w_0) > 0 \quad (1 \leq n \leq N)$$

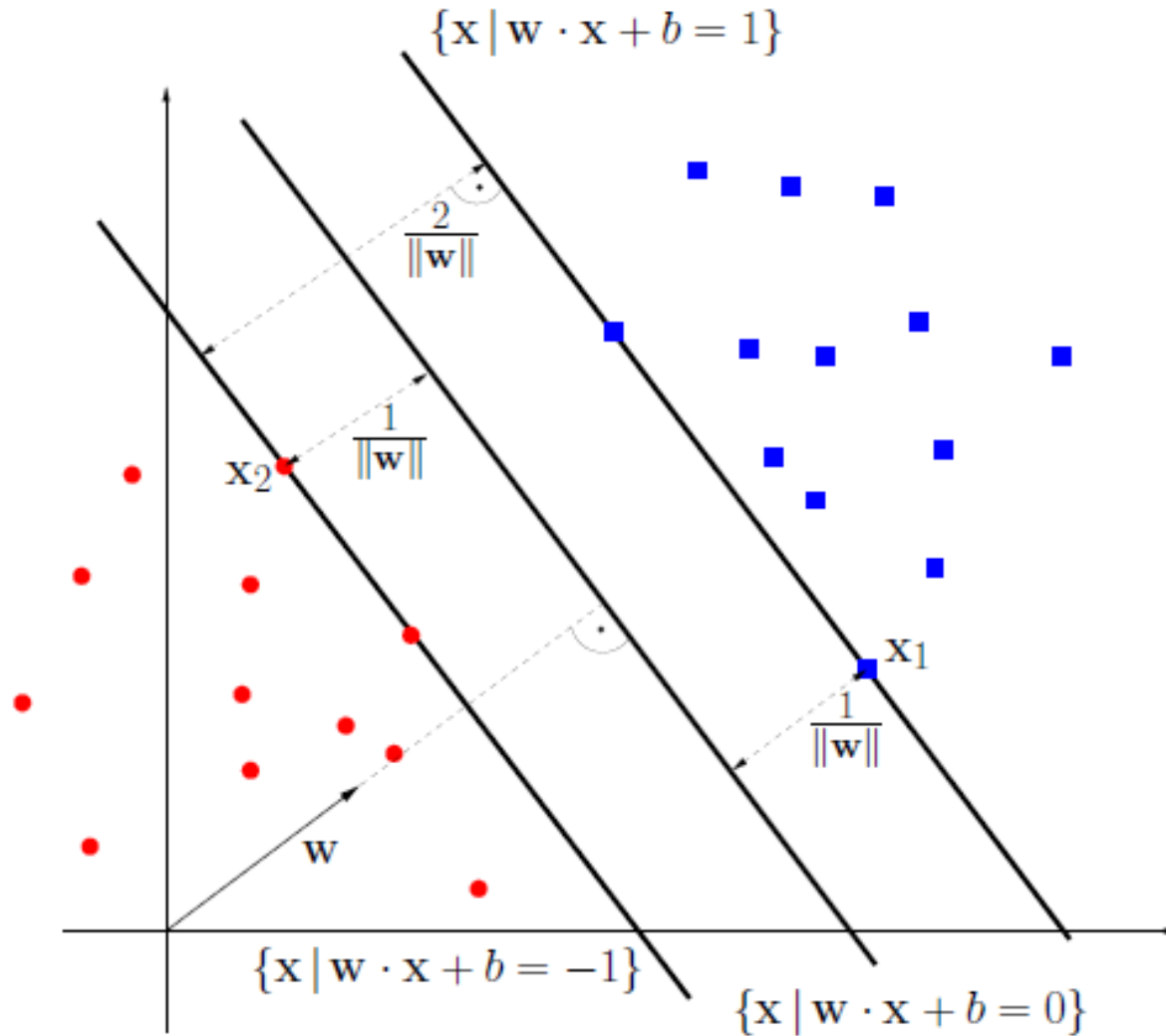
Support Vector Machines

Formalisation

- Rescaling \mathbf{w}, b such that $|\langle \mathbf{w}, \mathbf{x} \rangle + w_0| = 1$ for the points closest to the hyperplane, we obtain $|\langle \mathbf{w}, \mathbf{x} \rangle + w_0| \geq 1$
- The **support vectors** (SV) are the \mathbf{x}_n such that $|\langle \mathbf{w}, \mathbf{x}_n \rangle + w_0| = 1$
- The new loss is $\max(1 - t_n y(\mathbf{x}_n), 0) = (1 - t_n y(\mathbf{x}_n))_+$ (**hinge loss**)
- The **margin** of the OSH is twice the distance to any SV:
$$\text{margin(OSH)} = 2 d(\mathbf{x}_{SV}, \text{OSH}) = \frac{2}{\|\mathbf{w}\|}, \text{ since } |y(\mathbf{x}_{SV})| = 1$$

Support Vector Machines

Geometrical view



Support Vector Machines

A look on what's to come

1. We find the **canonical OSH** by solving

$$\max_{\mathbf{w}, b} \left\{ \frac{2}{\|\mathbf{w}\|} \quad / \quad t_n (\langle \mathbf{w}, \mathbf{x}_n \rangle + w_0) \geq 1, \quad 1 \leq n \leq N \right\}$$

2. The solution will be $\mathbf{w} = \sum_{n=1}^N t_n \alpha_n \mathbf{x}_n$, with $\alpha_n \geq 0$ (the **dual** form)
3. A fraction of the training \mathbf{x}_n will have $\alpha_n = 0$ (**sparsity**); the \mathbf{x}_n for which $\alpha_n > 0$ will coincide with the **support vectors**
4. The **discriminant function** will be written

$$y_{\text{SVM}}(\mathbf{x}) = \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle + w_0) = \text{sgn} \left(\sum_{n=1}^N t_n \alpha_n \langle \mathbf{x}, \mathbf{x}_n \rangle + w_0 \right)$$

Support Vector Machines

Using the VC dimension for two-class classification

For a two-class classifier, the **VC dimension** is the maximum number k of points that can be separated in all possible 2^k ways (*shattered*) by using functions representable by the classifier.

- Note it is *sufficient* that one set of k points exists that can be shattered for the VC dimension to be at least k
- If the VC dimension of a class is k , this means there is at least one set of k points that can be shattered by members of the class. It does not mean that every set of k points can be shattered.
- If no set of $k + 1$ points can be shattered by members of the class, then the VC dimension of the class is less than $k + 1$.

Support Vector Machines

Using the VC dimension for two-class classification

Theorem (Vapnik and Chervonenkis, 1974). Let D be an i.i.d data sample of size N and \mathcal{Y} a class of parametric binary classifiers. Let ϑ denote the VC dimension of \mathcal{Y} . Take $y \in \mathcal{Y}$ with empirical error $R_{\text{emp}}(y)$ on D . For all $\eta > 0$ it holds true that, with probability at least $1 - \eta$, the true error of y is bounded by:

$$R(y) \leq R_{\text{emp}}(y) + H(N, \vartheta, \eta)$$

where

$$H(N, \vartheta, \eta) = \sqrt{\frac{\vartheta(\ln(2N/\vartheta) + 1) - \ln(\eta/4)}{N}}$$

Support Vector Machines

More than an intuition

- Separating hyperplanes in \mathbb{R}^d have VC dimension $d + 1$
- When we use a feature map into a very high dimension $M \in (\mathbb{N} \cup \{\infty\})$, VC dimension will grow accordingly
- If we bound the margin of the hyperplanes, we limit VC dimension (therefore, we have an explicit control on complexity)

Support Vector Machines

More than an intuition

Theorem. Consider canonical hyperplanes $y(x) = \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle + w_0)$ w.r.t. a set of training patterns X (where $x_n \in \mathbb{R}^d$).

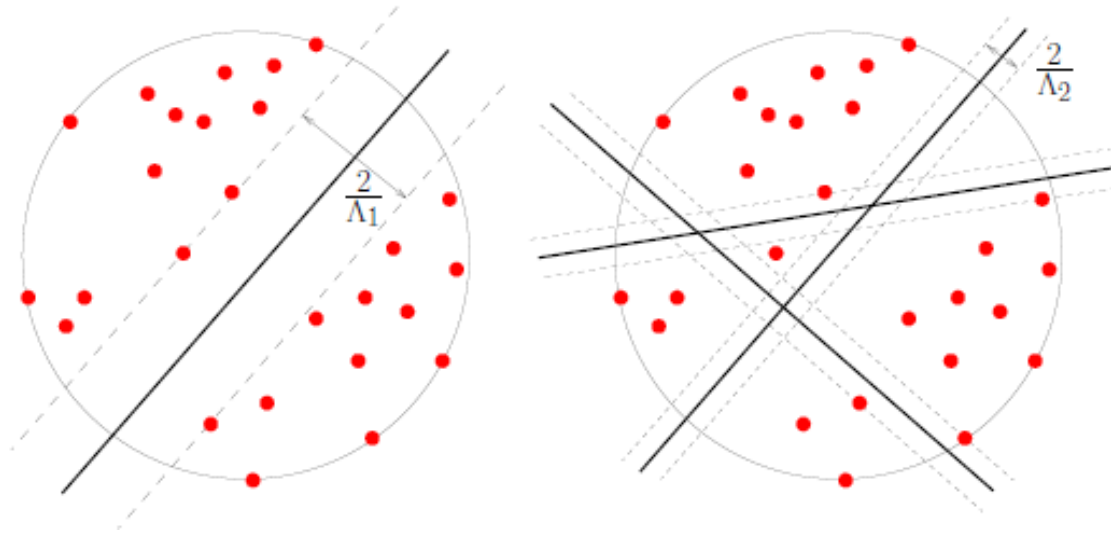
The subclass of linear classifiers with margin $m \geq m_0$ has VC dimension k bounded by

$$k \leq \min \left(\left\lceil \frac{R^2}{m_0^2} \right\rceil, d \right) + 1$$

where R is the radius of the smallest sphere centered at the origin containing X .

Support Vector Machines

More than an intuition



- Left: hyperplanes with a large margin have reduced chances to separate the data (the VC dimension is small)
- Right: smaller margins allow more separating hyperplanes (the VC dimension is large)

Support Vector Machines

Resolution

$$\begin{array}{ll}\text{minimize} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{subject to} & t_n (\langle \mathbf{w}, \mathbf{x}_n \rangle + w_0) \geq 1, \quad 1 \leq n \leq N\end{array}$$

This is solved (numerically) by QP techniques:

- Quadratic (therefore convex) function subject to linear constraints
- Unique solution (or set of equivalent ones); therefore, NO LOCAL MINIMA

Support Vector Machines

Margin violations

- In practice, we allow small margin violations ε_n , for each \mathbf{x}_n :

$$\begin{array}{ll} \text{minimize} & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \varepsilon_n \\ \text{subject to} & t_n (\langle \mathbf{w}, \mathbf{x}_n \rangle + w_0) \geq 1 - \varepsilon_n, \quad \varepsilon_n \geq 0, 1 \leq n \leq N \end{array}$$

- The ε_n are **slack** variables, leading to a **soft margin** ($\varepsilon_n > 0$ implying the functional margin $t_n y(\mathbf{x}_n) < 1$)

Support Vector Machines

Lagrangian form (primal)

- We first construct the Lagrangian:

$$\mathcal{L} = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N \alpha_n \left\{ t_n (\langle \mathbf{w}, \mathbf{x}_n \rangle + w_0) - 1 + \varepsilon_n \right\} + C \sum_{n=1}^N \varepsilon_n - \sum_{n=1}^N \mu_n \varepsilon_n$$

- The $\alpha_n, \mu_n \geq 0$ are the *Lagrange multipliers* (auxiliary variables to cope with the constraints)
- The solution is the saddle point of \mathcal{L} :
 1. the minimum of \mathcal{L} is taken with respect to \mathbf{w}, b
 2. the maximum of \mathcal{L} is taken with respect to the α_n

Support Vector Machines

Lagrangian form

The gradient of \mathcal{L} with respect to \boldsymbol{w}, b must vanish:

$$\frac{\delta \mathcal{L}}{\delta b} = \sum_{n=1}^N \alpha_n t_n = 0, \quad \frac{\delta \mathcal{L}}{\delta \boldsymbol{w}} = \boldsymbol{w} - \sum_{n=1}^N \alpha_n t_n \boldsymbol{x}_n = 0, \quad \frac{\delta \mathcal{L}}{\delta \varepsilon_n} = C - \alpha_n - \mu_n = 0$$

In addition, the so-called **KKT complementarity conditions** hold:

$$\alpha_n \left(t_n (\langle \boldsymbol{w}, \boldsymbol{x}_n \rangle + w_0) - 1 + \varepsilon_n \right) = 0, \quad 1 \leq n \leq N$$

Support Vector Machines

Dual formulation

The Lagrangian \mathcal{L} is convex; its optimization is equivalent to the maximization of its **dual problem** \mathcal{L}_D :

$$\begin{array}{ll} \text{maximize} & \mathcal{L}_D = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m \langle \mathbf{x}_n, \mathbf{x}_m \rangle \\ \\ \text{subject to} & 0 \leq \alpha_n \leq C, \quad 1 \leq n \leq N, \quad \text{and} \quad \sum_{n=1}^N \alpha_n t_n = 0 \end{array}$$

- Note how neither $\mu_n, \varepsilon_n, \mathbf{w}, w_0$ appear in the dual form
- Maximization is only with respect to the α_n

Support Vector Machines

Dual formulation

A closer look at the KKT complementarity conditions:

- $\alpha_n = 0$ implies $t_n y(\mathbf{x}_n) > 1$ and $\varepsilon_n = 0$ (\mathbf{x}_n is **not a SV**)
- $\alpha_n \in (0, C)$ implies $t_n y(\mathbf{x}_n) = 1$ and $\varepsilon_n = 0$ (\mathbf{x}_n is a **non-bound SV**)
- $\alpha_n = C$ implies $t_n y(\mathbf{x}_n) < 1$ and $\varepsilon_n > 0$ (\mathbf{x}_n is a **bound SV**)

Support Vector Machines

Result

- The discriminant function is:

$$y_{\text{SVM}}(\mathbf{x}) = \text{sgn} \left(\sum_{n=1}^N \alpha_n t_n \langle \mathbf{x}_n, \mathbf{x} \rangle + w_0 \right)$$

with $\alpha_n > 0$ *only* for the support vectors (for the rest, $\alpha_n = 0$).

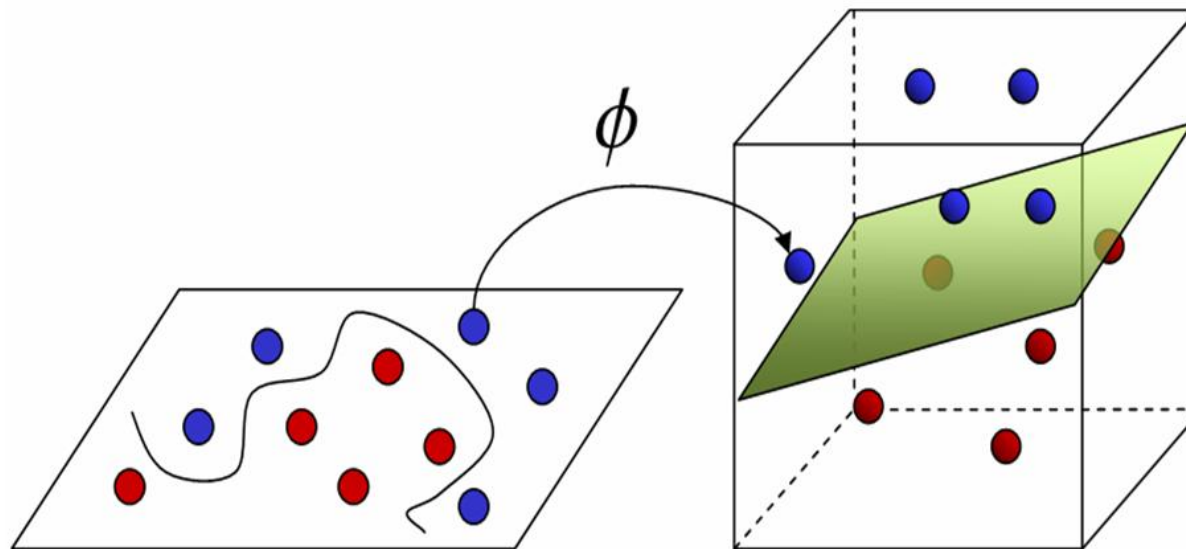
- By setting $C = \infty$ we are in the *hard* margin case
 $\Rightarrow C$ represents a tradeoff between **margin width** and **training error**
- $w_0 = - \left(\sum_{n=1}^N \sum_{m=1}^N \alpha_n t_n \alpha_m t_m \langle \mathbf{x}_n, \mathbf{x}_m \rangle \right) / \left(\sum_{n=1}^N \alpha_n \right)$

Support Vector Machines

General feature maps

Recall the idea of mapping input data into some Hilbert space (called the *feature space*) via a non-linear mapping $\phi : \mathcal{X} \rightarrow \mathcal{H}$

The associated kernel function is $k(u, v) = \langle \phi(u), \phi(v) \rangle$, $u, v \in \mathcal{X}$



Support Vector Machines

SVM training ... back to the OSH

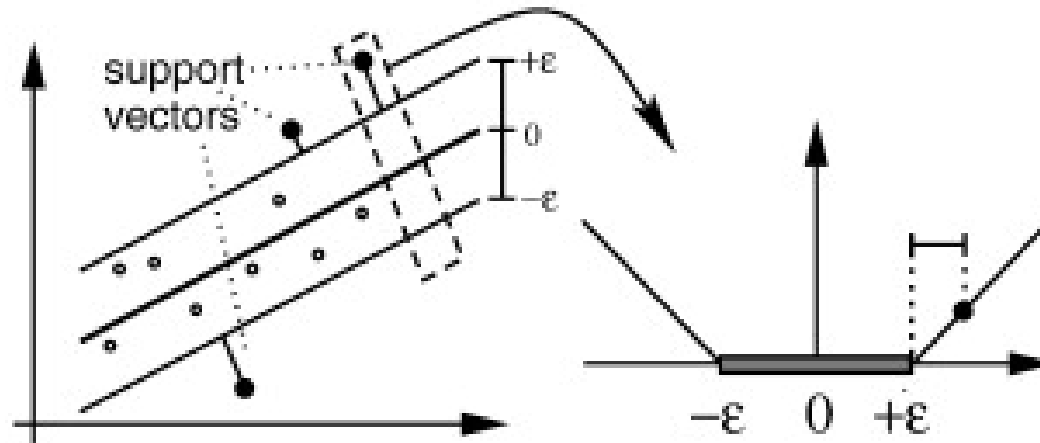
- We now substitute \mathbf{x}_n by $\phi(\mathbf{x}_n)$, then build the OSH in \mathcal{H}
- The discriminant function becomes:

$$y_{\text{SVM}}(\mathbf{x}) = \text{sgn} \left(\sum_{n=1}^N \alpha_n t_n k(\mathbf{x}, \mathbf{x}_n) + w_0 \right)$$

- The dual of the new QP problem is formulated exactly as before, replacing $\langle \mathbf{x}_n, \mathbf{x}_m \rangle$ with $k(\mathbf{x}_n, \mathbf{x}_m)$

Support Vector Machines

SVMs for regression



We choose the loss: $|t - y(x)|_{\epsilon} = \max(|t - y(x)| - \epsilon, 0)$

Support Vector Machines

SVMs for regression

$$\text{minimize} \quad \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N (\epsilon_n + \epsilon_n^*)$$

$$\text{subject to} \quad \langle w, x_n \rangle + w_0 - t_n \leq \epsilon + \epsilon_n \quad (1)$$

$$t_n - \langle w, x_n \rangle + w_0 \leq \epsilon + \epsilon_n^* \quad (2)$$

$$\epsilon_n, \epsilon_n^* \geq 0 \quad (3)$$

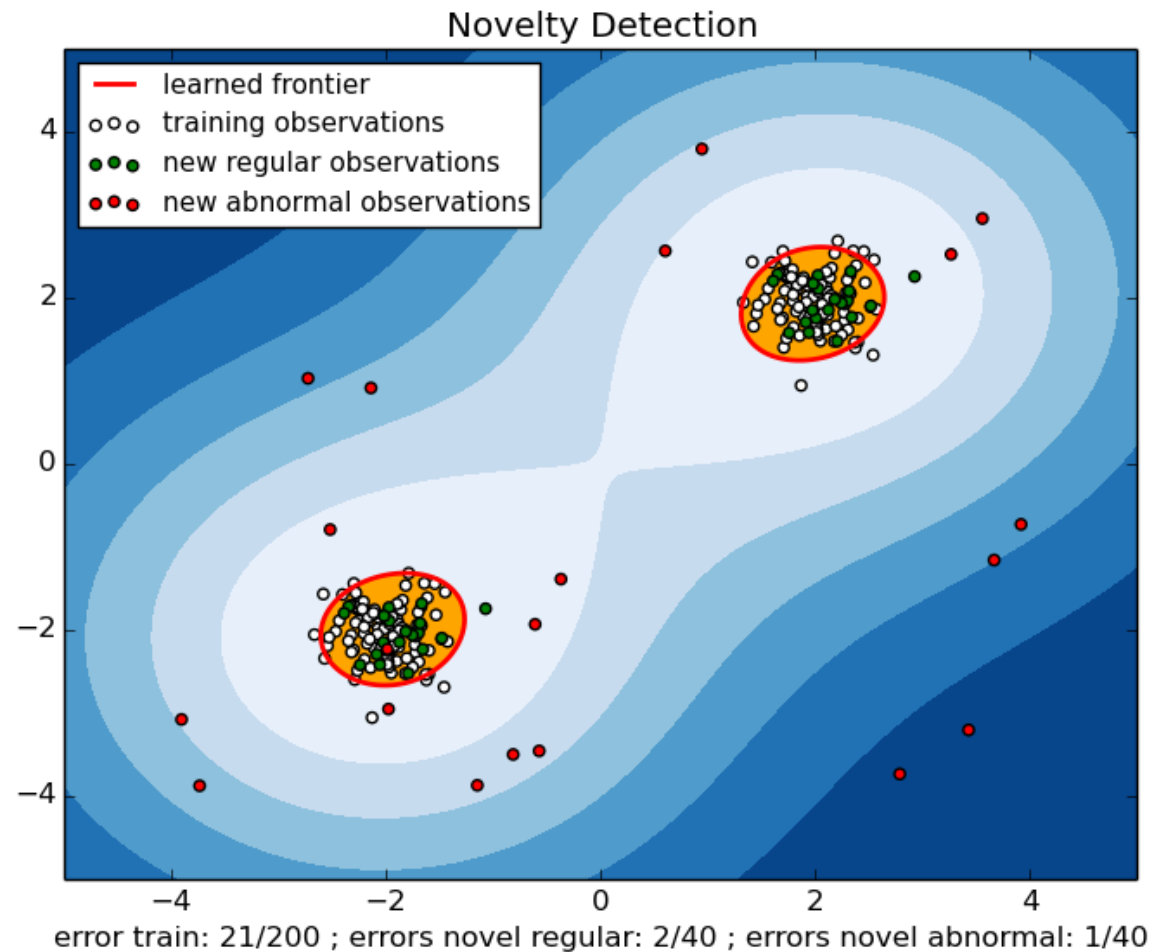
Support Vector Machines

SVMs for novelty detection

- You are given a dataset drawn from a pdf P ; you want to estimate a “simple” subset S of input space s.t. the probability that a test point drawn from P lies outside of S equals some a priori specified $\nu \in (0, 1)$
- The SVM approach to this problem is by estimating a function f which is positive on S and negative on its complement
- The ν parameter characterizes both the fraction of SVs and the fraction of outliers

Support Vector Machines

SVMs for novelty detection



Support Vector Machines

SVMs for novelty detection

USPS dataset of handwritten digits: 9,298 digit images of size 16×16

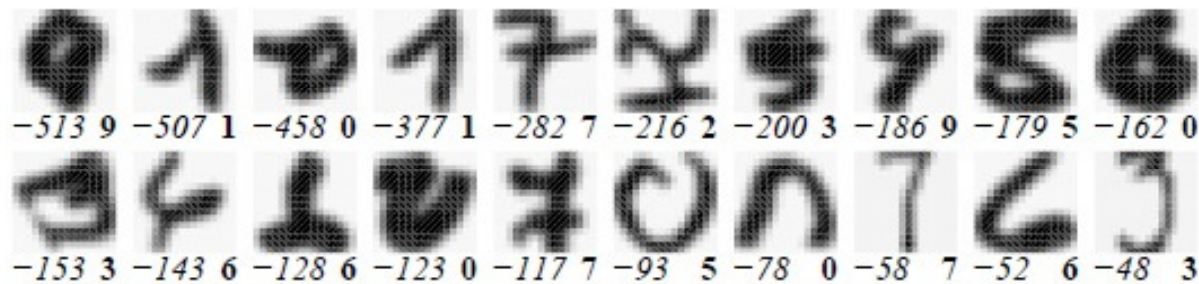


Figure 2: Outliers identified by the proposed algorithm, ranked by the negative output of the SVM (the argument of the sgn in the decision function). The outputs (for convenience in units of 10^{-5}) are written underneath each image in italics, the (alleged) class labels are given in bold face. Note that most of the examples are “difficult” in that they are either atypical or even mislabelled.

The 20 worst outliers for the USPS test set (here $\nu = 0,05$)

(from Schoelkopf et al, *Support Vector Method for Novelty Detection*, NIPS 2000)

Support Vector Machines

In conclusion ...

- (↑) No local minima, no initial conditions, few parameters to set
 - (↑) Provide with a compact description of original data
 - (↑) Decision surfaces can be RBFs, polynomials, multilayer perceptrons, ...
 - (↑) Accept any data type as input
 - (↑) Explicit complexity control via C , ϵ and the kernel parameters
 - (↑) Less affected by large input dimensions than other methods
 - (↑) Excellent practical performance: OCR, text categorization, face detection, ...
-
- (↓) Training requires solving a QP problem
 - (↓) Choice of best kernel is an open issue; **kernel design** is an active area of research
 - (↓) Sometimes fraction of SVs very high (indicating a poor model)
 - (↓) Performance usually depends on a careful choice of the parameters

Machine Learning

Syllabus

1. Introduction to Machine Learning
2. Theoretical issues (I): regression
3. Linear regression and beyond
4. Theoretical issues (II): classification
5. Generative classifiers
6. Discriminative classifiers

7. Clustering
8. Learning with kernels (I): The SVM
9. Learning with kernels (II): Kernel functions
10. Learning with kernels (III): Other kernel methods
11. Artificial neural networks (I): the MLP
12. Artificial neural networks (II): the RBF
13. Ensemble methods: Random Forests
14. Advanced topics and frontiers