

Machine learning 2

Support Vector Machines (SVM)

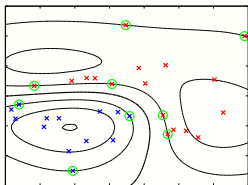
Pierre Chainais



1 SVM: Support Vector Machines

- SVM in brief
- Recall on linear classifiers
- Maximizing the margin : the linearly separable case
- Maximizing the margin: the non separable case
- SVM: the redescription space and the kernel
- Complements
- In summary...

- ▶ $\mathbf{x} \in \mathbb{R}^D$ et 2 classes $y \in \{-1, 1\}$
- ▶ Classification algorithm: **non-linear**
- ▶ Embeds inputs $\mathbf{x} \in \mathbb{R}^D$ to a higher dimensional space (potentially infinite) $\mathcal{F} : \mathbf{x} \Rightarrow \phi(\mathbf{x})$
- ▶ Applies a **linear** classification in this space \mathcal{F} ...
- ▶ ... without knowing explicitly functions $\phi = \text{"kernel trick"}$!

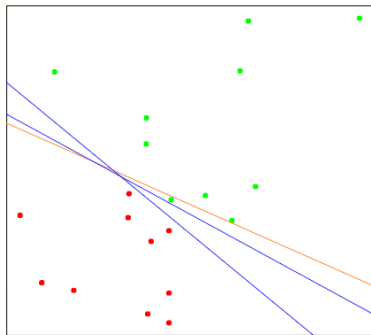


- ▶ $\mathbf{x} \in \mathbb{R}^D$ and 2 classes $y \in \{-1, 1\}$
- ▶ Training set (x_i, y_i) , $i = 1, \dots, N$
- ▶ First assume that the training set is **linearly separable** in \mathbb{R}^D
- ▶ Classification
= search for a hyperplane \mathcal{H}

$$h(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta} + \beta_0 = 0, \quad \boldsymbol{\beta} \in \mathbb{R}^D, \beta_0 \in \mathbb{R}$$

perfectly separating the 2 classes.

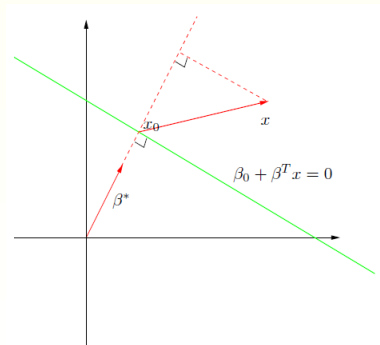
- ▶ If these classes are separable, \exists an infinite set of planes



- ▶ For all points $(\mathbf{x}_1, \mathbf{x}_2)$ on \mathcal{H} , $\beta^T(\mathbf{x}_1 - \mathbf{x}_2) = 0$
- ▶ $\beta^* = \beta / \|\beta\|$ is orthogonal to the hyperplane \mathcal{H} ,
- ▶ For all point \mathbf{x}_0 on \mathcal{H} , $\beta^T \mathbf{x}_0 = -\beta_0$

- ▶ The signed distance $d(\mathbf{x}, \mathcal{H})$ is

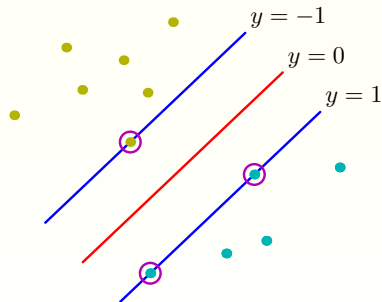
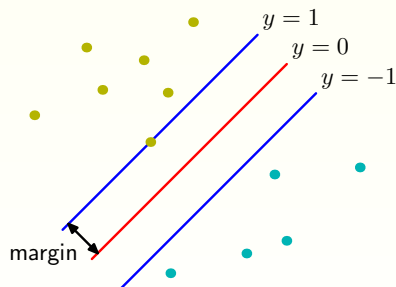
$$\beta^{*T}(\mathbf{x} - \mathbf{x}_0) = \frac{1}{\|\beta\|}(\beta^T \mathbf{x} + \beta_0)$$



Maximizing the margin : the linearly separable case

p.6

Notion of *margin*

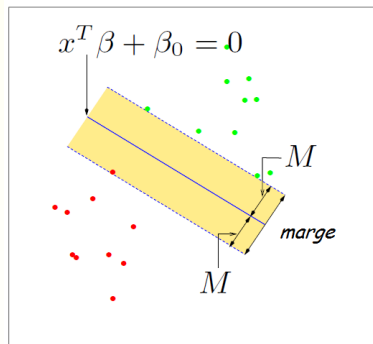


- ▶ One considers the optimisation problem

$$(\beta, \beta_0) = \operatorname{argmax}_{\beta, \beta_0, \|\beta\|=1} M$$

$$\text{s.t. } y_i(\mathbf{x}_i^T \beta + \beta_0) \geq M, i = 1, \dots, N$$

- ▶ These N conditions guarantee that all points are at least at a distance M of the hyperplane ($y_i = \pm 1$)
- ▶ We aim at maximizing the margin M

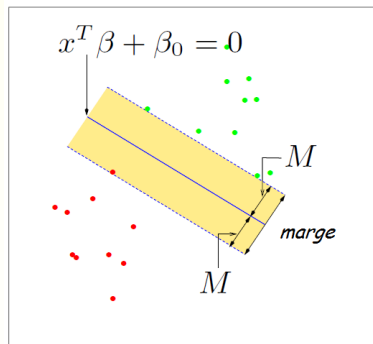


- ▶ One considers the optimisation problem

$$(\beta, \beta_0) = \operatorname{argmax}_{\beta, \beta_0, \|\beta\|=1} M$$

$$\text{s.c. } y_i(\mathbf{x}_i^T \beta + \beta_0) \geq M, i = 1, \dots, N$$

- ▶ $M = \min_{\mathbf{x}_i} \frac{1}{\|\beta\|} y_i(\beta^T \mathbf{x}_i + \beta_0)$
- ▶ Pb : $\operatorname{argmax}_{\beta, \beta_0} \frac{1}{\|\beta\|} \min_{\mathbf{x}_i} [y_i(\beta^T \mathbf{x}_i + \beta_0)]$



Support Vector Machine: primal problem

p.9

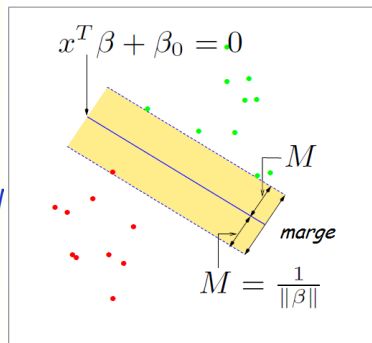
Formulation canonique

- ▶ Let $\min_{\mathbf{x}_i} [y_i(\beta^T \mathbf{x}_i + \beta_0)] = 1$, one gets the equivalent formulation [*primal problem*]

$$(\beta, \beta_0) = \operatorname{argmin}_{\beta, \beta_0} \frac{1}{2} \|\beta\|^2$$

such that $y_i(\mathbf{x}_i^T \beta + \beta_0) \geq 1, i = 1, \dots, l$

- ▶ Quadratic criterion with linear equality constraints
- ▶ **Convex** optimisation problem



- Lagrangian formulation

$$L(\beta, \beta_0, \alpha) = \frac{1}{2} \|\beta\|^2 - \sum_{i=1}^N \alpha_i [y_i (\mathbf{x}_i^T \beta + \beta_0) - 1]$$

où $\alpha_i \geq 0$, $i = 1, \dots, N$ sont les multiplicateurs de Lagrange

- Zeroes of the derivatives w.r.t. β et β_0 yield:

$$\beta = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i, \quad 0 = \sum_{i=1}^N \alpha_i y_i$$

- By substitution, on gets the dual formulation

$$L_{dual} = \sum_{i=1}^N \alpha_i - \underbrace{\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_i \alpha_j y_j \mathbf{x}_i^T \mathbf{x}_j}_{\|\beta\|^2}$$

- ▶ This is a quadratic programming problem of dimension n [*dual problem*]

$$\operatorname{argmax}_{\alpha} \left[\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_i \alpha_j y_j \mathbf{x}_i^T \mathbf{x}_j \right]$$

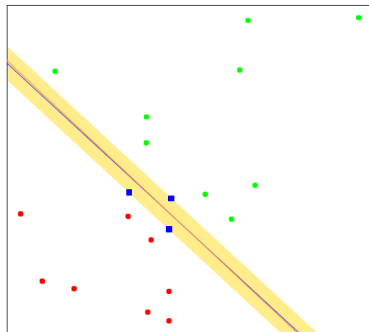
such that (constraints) $\forall i \quad \alpha_i \geq 0, \quad \sum_{i=1}^N \alpha_i y_i = 0$

- ▶ This is a convex optimization problem \Rightarrow Sequential Minimal (SMO) Optimization [Platt, 1999]
- ▶ Remark : depends on $\mathbf{x}_i^T \mathbf{x}_j$ only

- ▶ On a

$$\hat{\beta} = \sum_{i=1}^N \hat{\alpha}_i y_i \mathbf{x}_i$$

- ▶ Only few $\hat{\alpha}_i$, corresponding to points x_i on the margin, are non zero
- ▶ The corresponding x_i are called **support vectors**
- ▶ **Sparse** solution ("parcimonieuse" in french)



- One uses the function

$$h(\mathbf{x}) = \mathbf{x}^T \hat{\boldsymbol{\beta}} + \hat{\beta}_0 = \sum_{i=1}^N \hat{\alpha}_i y_i \mathbf{x}_i^T \mathbf{x}$$

to classify new elements

$$f(\mathbf{x}) = \hat{y}(\mathbf{x}) = \text{signe}[h(\mathbf{x})]$$

- Remark : depends on $\mathbf{x}_i^T \mathbf{x}$ only.

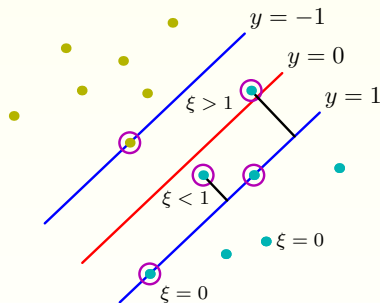
- ▶ One uses the function

$$h(\mathbf{x}) = \mathbf{x}^T \hat{\boldsymbol{\beta}} + \hat{\beta}_0 = \sum_{i=1}^N \hat{\alpha}_i y_i \mathbf{x}_i^T \mathbf{x}$$

to classify new elements

$$f(\mathbf{x}) = \hat{y}(\mathbf{x}) = \text{signe}[h(\mathbf{x})]$$

- ▶ Remark : depends on $\mathbf{x}_i^T \mathbf{x}$ only.



- ▶ Tolerance of the superposition of 2 classes
- ▶ Soft margin that tolerates small classification errors in the training set

- ▶ Introduction of **spring variables** $\xi_i > 0$, $i = 1, \dots, N$
- ▶ Primal problem where C is a constant ≥ 0

$$(\beta, \beta_0, \xi) = \operatorname{argmin}_{\beta, \beta_0, \xi} \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i$$

under conditions

$$\begin{cases} \xi_i \geq 0 \\ y_i(\mathbf{x}_i^T \beta + \beta_0) \geq 1 - \xi_i, i = 1, \dots, N \end{cases}$$

- Dual formulation

$$\max_{\alpha} \left[\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_i \alpha_j y_j \mathbf{x}_i^T \mathbf{x}_j \right]$$

with constraints

$$\forall i \quad 0 \leq \alpha_i \leq C, \quad \sum_{i=1}^N \alpha_i y_i = 0$$

- Difference : upper bound C on the α_i
- Problem depending on the $\mathbf{x}_i^T \mathbf{x}_j$ only...
- Convex problem - solution by quadratic programming
- Read Bishop pp. 331-334 for primal \Rightarrow dual

- Dual formulation

$$\max_{\alpha} \left[\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_i \alpha_j y_j \mathbf{x}_i^T \mathbf{x}_j \right]$$

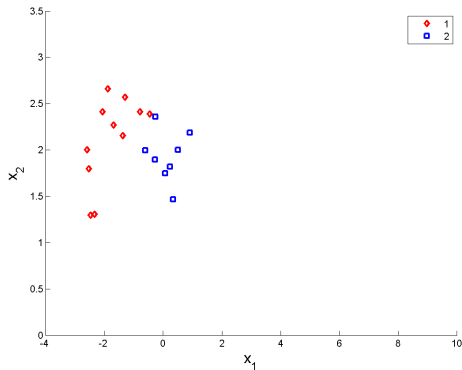
with constraints

$$\forall i \quad 0 \leq \alpha_i \leq C, \quad \sum_{i=1}^N \alpha_i y_i = 0$$

- Difference : upper bound C on the α_i
- Problem depending on the $\mathbf{x}_i^T \mathbf{x}_j$ only...
- Convex problem - solution by quadratic programming
- Read Bishop pp. 331-334 for primal \Rightarrow dual

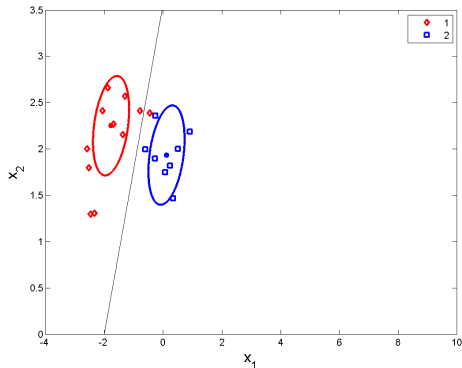
- ▶ SVM & LDA
= linear classification
- ▶ LDA = generative model (gaussian) for each class
- ▶ Atypical points, even far from the boundary \Rightarrow critical influence on the classification rule
- ▶ Points that are far from the margin \Rightarrow no influence on SVM

Linear discriminant analysis: influence of an outlier



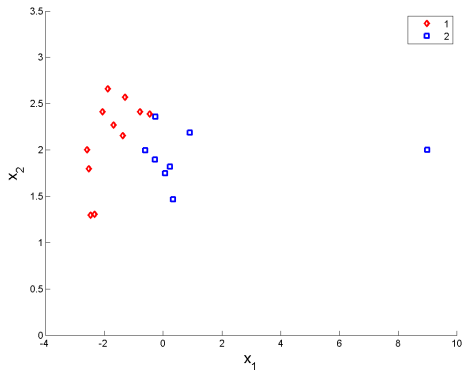
Training set without outlier

Linear discriminant analysis: influence of an outlier



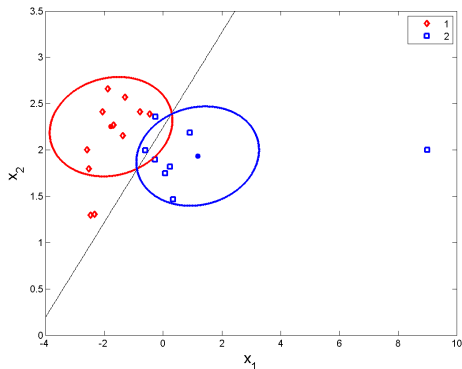
Linear discriminant analysis (without outlier)

Linear discriminant analysis: influence of an outlier



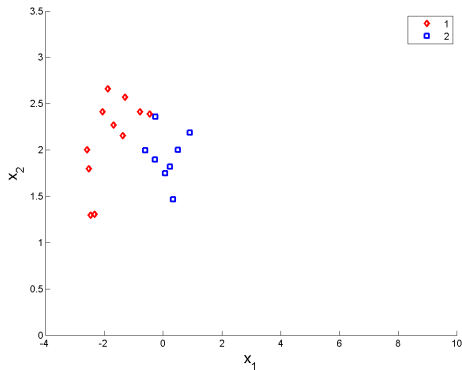
Training set with outlier

Linear discriminant analysis: influence of an outlier



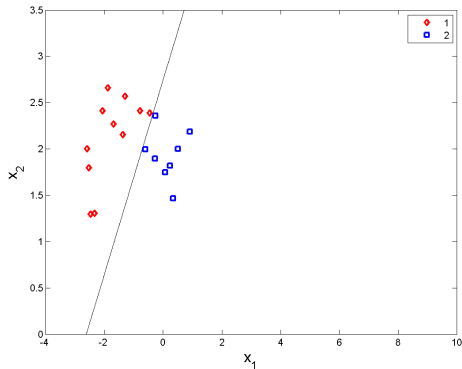
Linear discriminant analysis (with outlier)

SVM: no influence of an outlier



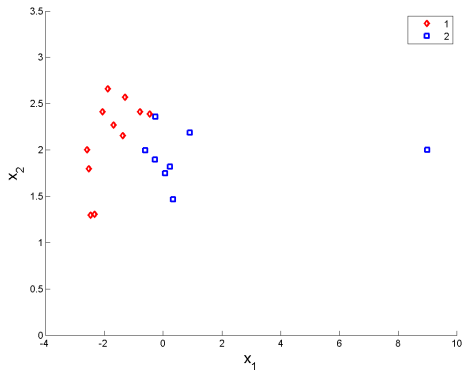
Training set without outlier

SVM: no influence of an outlier



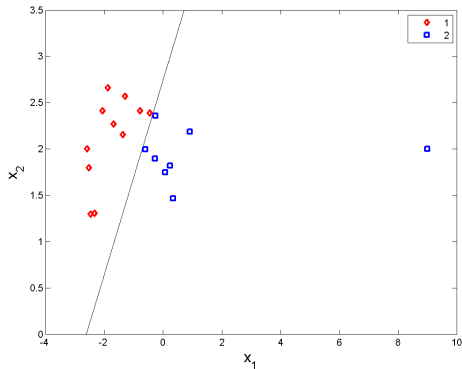
SVM (without outlier)

SVM: no influence of an outlier



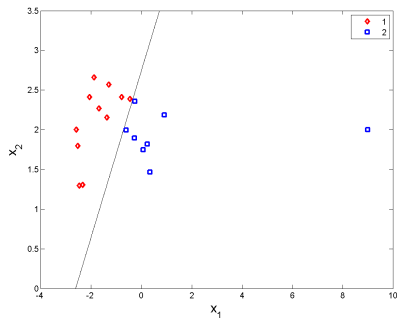
Training set with outlier

SVM: no influence of an outlier

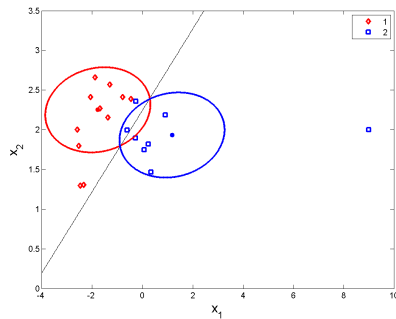


SVM (with outlier)

Comparison with linear discriminant analysis



SVM



Linear discriminant analysis

- ▶ In place of searching for a separating hyperplane in the inputs space, we embed the data in a redescription (feature space) of higher dimension

$$\begin{aligned}\phi : \mathbb{R}^D &\rightarrow \mathcal{F} \\ x &\rightarrow \phi(\mathbf{x})\end{aligned}$$

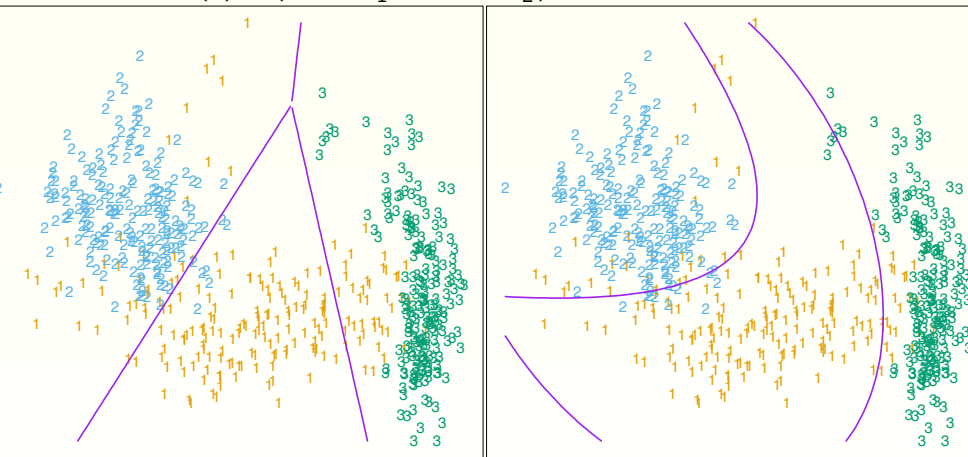
- ▶ Linear separation in \mathcal{F} yields a non-linear separation in the inputs space \mathbb{R}^D

Back to LDAe

p.30

Application to $\phi_j(\mathbf{x})$: example, $x_1, x_2, x_1^2, x_1x_2, x_2^2$

$$\mathbf{x} \in \mathbb{R}^2, \phi(\mathbf{x}) = (x_1, x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2)$$



Linear classification in the redescription

\iff non linear classification in the input space

- Dual formulation

$$\max_{\alpha} \left[\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_i \alpha_j y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \right]$$

with constraints

$$\forall i \quad 0 \leq \alpha_i \leq C, \quad \sum_{i=1}^N \alpha_i y_i = 0$$

- Solution given by

$$\begin{aligned} h(\mathbf{x}) &= \hat{\beta}^T \phi(\mathbf{x}) + \hat{\beta}_0 \\ &= \sum_{i=1}^N \hat{\alpha}_i y_i \underbrace{\phi(\mathbf{x}_i)^T \phi(\mathbf{x})}_{k(\mathbf{x}_i, \mathbf{x})} + \hat{\beta}_0 \end{aligned}$$

- ▶ The problem and its solution depend on the scalar product $\phi(\mathbf{x})^T \phi(\mathbf{x}')$ only
- ▶ One denotes by $k : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$ the **kernel function** defined by

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

- ▶ Kernel \Rightarrow computing in the input space
without making explicit the embedding $\mathbf{x} \rightarrow \phi(\mathbf{x})$
- ▶ This is the *kernel trick*)

- ▶ $\mathbf{x} \in \mathbb{R}^2$, $\phi(\mathbf{x}) = (x_1, x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2)$
- ▶ The scalar product in the redescription space is given by

$$\begin{aligned}\phi(\mathbf{x})^T \phi(\mathbf{x}') &= x_1x_1' + x_2x_2' + x_1^2x_1'^2 + 2x_1x_1'x_2x_2' + x_2^2x_2'^2 \\ &= (x_1x_1' + x_2x_2') + (x_1x_1' + x_2x_2')^2 \\ &= (x^T x') + (x^T x')^2 = k(x, x')\end{aligned}$$

- ▶ Can be computed without using ϕ (potentially unknown)

Theorem: Mercer's conditions

A symmetric function $k : \mathbb{R}^D \times \mathbb{R}^D$ is a kernel if for all \mathbf{x}_i , $(k(\mathbf{x}_i, \mathbf{x}_j))_{i,j}$ is a positive definite matrix. Then, there exists a space \mathcal{F} and a (vectorial) function ϕ such that

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

- ▶ This condition is difficult to check
- ▶ It does not permit to determine neither \mathcal{F} nor ϕ

- ▶ Linear

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

- ▶ Polynomials of order d

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^d \text{ ou } (1 + \mathbf{x}^T \mathbf{x}')^d$$

- ▶ Gaussian (radial basis)

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$$

- ▶ Neural network

$$k(\mathbf{x}, \mathbf{x}') = \tanh(\kappa_1 \mathbf{x}^T \mathbf{x}' + \kappa_2)$$

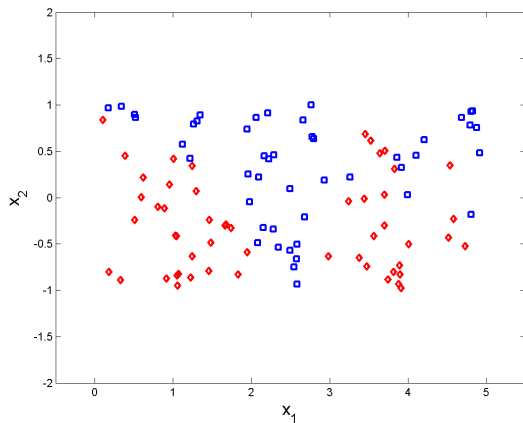
On reprend tout en remplaçant partout :

$$\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

\Rightarrow Gram matrix
and one gets:

$$\hat{\alpha} = \operatorname{argmax}_{\alpha} \left[\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_i \alpha_j y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \right]$$

$$h(\mathbf{x}) = \sum_{i=1}^N \hat{\alpha}_i y_i \kappa(\mathbf{x}_i, \mathbf{x}) + \hat{\beta}_0$$

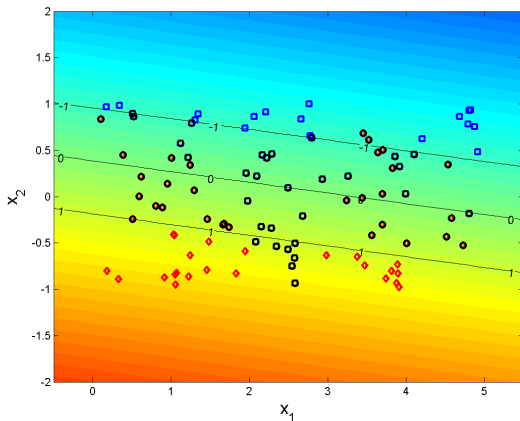


Training set

Support vector machines

Example

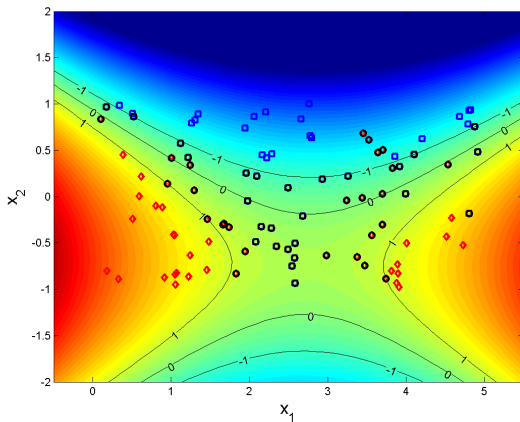
p.38



Linear kernel

Support vector machines

Example

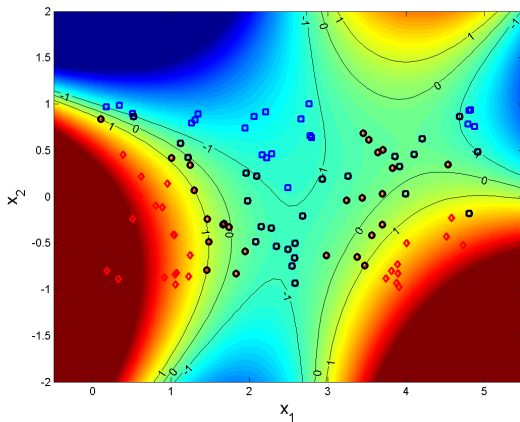


Polynomial kernel of order 2

Support vector machines

Example

p.40

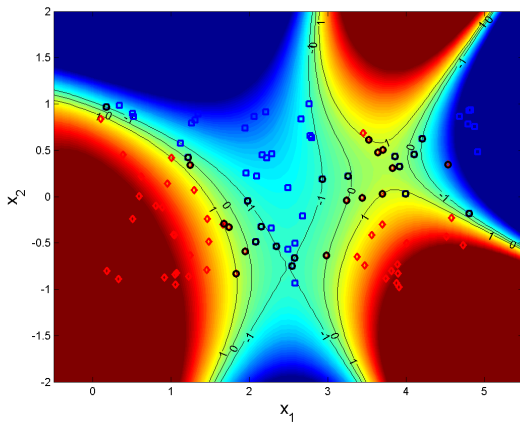


Polynomial kernel of order 3

Support vector machines

Example

p.41



Polynomial kernel of order 4

- ▶ Dimension of input data: D
- ▶ Size of the training set: N
- ▶ Between $O(DN^2)$ and $O(DN^3)$
- ▶ Therefore no too big N , but maybe high dim. D

- ▶ $y \in \{1, \dots, K\}$, K the number of classes
- ▶ Using several binary SVMs
 - One against all: K classifiers using one class vs others
 - One against one: $K(K - 1)/2$ classifiers using all possible pairs of classes
- ▶ Global objective function: K SVM simultaneously (expensive...)
- ▶ 1 class only: for anomaly detection

- ▶ One considers the following optimization problem

$$\underbrace{\min_{\beta, \beta_0} \sum_{i=1}^N [1 - y_i h(\mathbf{x}_i)]_+}_{\text{training error}} + \underbrace{\frac{\lambda}{2} \|\beta\|^2}_{\text{Regularisation}}$$

where $[x]_+ = \max(x, 0)$, $h(\mathbf{x}) = \phi(\mathbf{x})^T \beta + \beta_0$ and $\lambda = \frac{1}{C}$

- ▶ One searches for a solution of the form $\sum_i \alpha_i k(\mathbf{x}_i, \cdot)$
- ▶ Solution = SVM \Rightarrow representation $h(\mathbf{x}) = \sum_i \hat{\alpha}_i k(\mathbf{x}_i, \mathbf{x})$

Kernels for any kind of data

- ▶ Texts,
- ▶ Trees,
- ▶ Graphs,
- ▶ DNA sequences...

Kernel methods:

- ▶ Kernel PCA,
- ▶ Kernel FLD,
- ▶ Kernel clustering,
- ▶ One class SVM...

<http://rvlasveld.github.io/blog/2013/07/12/introduction-to-one-class-support-vector-machines/>

4 important ideas :

- ❶ data $\in \mathcal{X}$: **implicit non-linear proj.** on a vectorial space \mathcal{F} ,
- ❷ **linear classification** in \mathcal{F} (quadratic optimisation)
- ❸ algorithms use **the kernel** $k(\mathbf{x}_i, \mathbf{x}_j)$ **only**
- ❹ **support vectors** \Rightarrow a sparse representation of the classifier

Propriétés :

- ▶ No local minima,
- ▶ Kernel trick: computations in the input space,
- ▶ Control of the risk of overfitting (spring variables),
- ▶ Not too many parameters to tune (\neq neural networks),
- ▶ Robustness: support vectors,
- ▶ Very good results in general (often state of the art)
- ▶ Pb : computational cost if large N (big data)