

Support vector machines

Pavlo Mozharovskyi¹

¹LTCI, Télécom Paris, Institut Polytechnique de Paris

Machine learning

Paris, March 12, 2022

Today

Vapnik-Chervonenkis theory, simplest case

The support vector machine

- Optimal margin classifier

- Introducing kernels (1992)

- Allowing for misclassification: soft margin (1995)

Implementation

Literature

Learning materials include but are not limited to:

- ▶ Hastie, T., Tibshirani, R., and Friedman, J. (2009).
The Elements of Statistics Learning: Data Mining, Inference, and Prediction (Second Edition).
Springer.
 - ▶ Section 12.{1,2,3.1,3.2}.
- ▶ Slides of the lecture.
- ▶ Boser, B. E., Guyon, I. M., and V. N. Vapnik (1992).
A training algorithm for optimal margin classifiers.
In: *Proceedings of the Fifth Annual Workshop of Computational Learning Theory*, Pittsburgh, ACM, 5, 144–152.
- ▶ Cortes, C. and Vapnik, V. (1995).
Support-vector networks.
Machine learning, 20, 273–297.
- ▶ Vapnik, V. N. (1998).
Statistical Learning Theory.
John Wiley & Sons.

Binary supervised classification (reminder)

Notation:

- ▶ **Given:** for the random pair (X, Y) in $\mathbb{R}^d \times \{-1, 1\}$ consisting of a random observation X and its random binary label Y (class), a sample of n i.i.d.: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$.
- ▶ **Goal:** predict the label of the new (unseen before) observation \mathbf{x} .
- ▶ **Method:** construct a classification rule:

$$g : \mathbb{R}^d \rightarrow \{-1, 1\}, \mathbf{x} \mapsto g(\mathbf{x}),$$

so $g(\mathbf{x})$ is the prediction of the label for observation \mathbf{x} .

- ▶ **Criterion:** of the performance of g is the **error probability**:

$$R(g) = \mathbb{P}[g(X) \neq Y] = \mathbb{E}[\mathbb{1}(g(X) \neq Y)].$$

- ▶ **The best solution:** is to know the distribution of (X, Y) :

$$g(\mathbf{x}) = \text{sign}(2\mathbb{E}[Y|X = \mathbf{x}] - 1 > 0).$$

Contents

Vapnik-Chervonenkis theory, simplest case

The support vector machine

- Optimal margin classifier

- Introducing kernels (1992)

- Allowing for misclassification: soft margin (1995)

Implementation

Glivenko-Cantelli theorem

Consider the classification rule for a new observation \mathbf{x} given the weights vector \mathbf{w} :

$$g(\mathbf{x}, \mathbf{w}) = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} > 0, \\ -1 & \text{otherwise.} \end{cases}$$

What can be said about the **error probability**, i.e. about the relationship between

$$\mathbb{P}(g(X, \mathbf{w}) \neq Y) = \int_{\mathbb{R}^d} \mathbb{1}(g(\mathbf{x}, \mathbf{w}) \neq Y) dF_X \quad \text{and} \quad \frac{1}{n} \sum_{i=1}^n \mathbb{1}(g(\mathbf{x}_i, \mathbf{w}) \neq y_i) ?$$

Let X_1, \dots, X_n be a random sample on \mathbb{R} . The **empirical distribution function** is defined as

$$\mathbb{F}_n(t) = \frac{1}{n} \sum \mathbb{1}(X_i \leq t).$$

Theorem (Glivenko-Cantelli)

If X_1, X_2, \dots are i.i.d. random variables with distribution function F , then

$$\|\mathbb{F}_n - F\|_\infty = \sup_{x \in \mathbb{R}} |\mathbb{F}_n(x) - F(x)| \xrightarrow{\text{a.s.}} 0.$$

Uniform one-sided convergence

Under *additional* conditions, for $g(\mathbf{x}, \mathbf{w})$ and a probability measure F_X , for any $\epsilon > 0$ it holds

$$\mathbb{P} \left\{ \underbrace{\sup_{\mathbf{w}} \left(\mathbb{P}(g(X, \mathbf{w}) \neq Y) \right)}_{L(g(\cdot, \mathbf{w}))} - \underbrace{\frac{1}{n} \sum_{i=1}^n \mathbb{1}(g(X_i, \mathbf{w}) \neq Y_i)}_{L_{emp}(g(\cdot, \mathbf{w}))} > \epsilon \right\} \xrightarrow{n \rightarrow \infty} 0.$$

What can be said about the **rate of convergence**?

Regard **finite set of classification rules** $g(\mathbf{x}, \mathbf{w}_k)$, $k = 1, \dots, N$. The restriction is naturally posed by the finite number of elements in the training set.

$$\begin{aligned} & \mathbb{P} \left\{ \sup_{k \in \{1, \dots, N\}} \left(L(g(\cdot, \mathbf{w}_k)) - L_{emp}(g(\cdot, \mathbf{w}_k)) \right) > \epsilon \right\} \\ & \leq \sum_{k=1}^N \mathbb{P} \left\{ \left(L(g(\cdot, \mathbf{w}_k)) - L_{emp}(g(\cdot, \mathbf{w}_k)) \right) > \epsilon \right\} \end{aligned}$$

Uniform one-sided convergence

Theorem (Chernoff-Hoeffding, Bernoulli scheme)

If X_1, \dots, X_n are i.i.d. random variables taking values in $\{0, 1\}$, then for any $\epsilon > 0$ it holds

$$\mathbb{P}\left(\mathbb{E}[X_i] - \frac{1}{n} \sum_{i=1}^n X_i > \epsilon\right) < e^{-2\epsilon^2 n}.$$

This allows for:

$$\begin{aligned} & \sum_{k=1}^N \mathbb{P}\left\{\left(L(g(\cdot, \mathbf{w}_k)) - L_{emp}(g(\cdot, \mathbf{w}_k))\right) > \epsilon\right\} \\ = & \sum_{k=1}^N \mathbb{P}\left\{\left(\underbrace{\mathbb{P}(g(X, \mathbf{w}_k) \neq Y)}_{\mathbb{E}[\mathbb{1}(g(X, \mathbf{w}_k) \neq Y)]} - \frac{1}{n} \sum_{i=1}^n \mathbb{1}(g(X_i, \mathbf{w}_k) \neq Y_i)\right) > \epsilon\right\} \\ \leq & Ne^{-2\epsilon^2 n}. \end{aligned}$$

Vapnik-Chervonenkis inequality

So:

$$\mathbb{P}\left\{ \sup_{k \in \{1, \dots, N\}} \left(L(g(\cdot, \mathbf{w}_k)) - L_{emp}(g(\cdot, \mathbf{w}_k)) \right) > \epsilon \right\} \leq N e^{-2\epsilon^2 n}.$$

Let us fix this probability having chosen $0 < \eta \leq 1$, by that maintaining reliability $1 - \eta$:

$$N e^{-2\epsilon^2 n} = \eta \quad \text{or equivalently} \quad \epsilon = \sqrt{\frac{\log N - \log \eta}{2n}}.$$

This allows for the following result:

Theorem (Vapnik-Chervonenkis, 1974)

If from a set consisting of N classification rules a rule $g(\cdot, \mathbf{w})$ is chosen, which delivers empirical risk $L_{emp}(g(\cdot, \mathbf{w}))$, then with reliability $1 - \eta$ one can state that the error probability $L(g(\cdot, \mathbf{w}))$ is bounded from above as follows

$$L(g(\cdot, \mathbf{w})) \leq L_{emp}(g(\cdot, \mathbf{w})) + \sqrt{\frac{\log N - \log \eta}{2n}}.$$

Particular case: linear rule

Let us try to estimate N for the linear classification rule.

The number $\Phi(d, n)$ of all possible separations of n points in \mathbb{R}^d by a hyperplane via the origin is computed as

$$\Phi(d, n) = \begin{cases} 2 \sum_{l=0}^{d-1} \binom{n-1}{l} & \text{if } d \leq n, \\ 2^n & \text{otherwise.} \end{cases}$$

For $d \leq n$, one can approximate it from above using:

$$\Phi(d, n) \leq 3 \frac{n^{d-1}}{(d-1)!} \leq n^d.$$

Plugging this into the Vapnik-Chervonenkis inequality gives:

$$L(g(\cdot, \mathbf{w})) \leq L_{emp}(g(\cdot, \mathbf{w})) + \sqrt{\frac{d \log n - \log \eta}{2n}}.$$

Contents

Vapnik-Chervonenkis theory, simplest case

The support vector machine

- Optimal margin classifier

- Introducing kernels (1992)

- Allowing for misclassification: soft margin (1995)

Implementation

The principle

- ▶ The conservative upper **bound of Vapnik and Chervonenkis** is very **pessimistic**, as even for a linear classification rule a very large training data set is required to guarantee meaningfulness of the achieved empirical risk.
- ▶ As an example, consider the case of two **linearly separable** training classes. Even in this case, only **little can be said** about probability of **points from one class inside the other** one.
- ▶ Sticking to this “trivial” case, the safest **separating hyperplane** would be the one having **maximal** and equal **margin** to each of the classes.
- ▶ Finding such a hyperplane in a systematic way constitutes the main idea of the **optimal margin hyperplane** algorithm.

Contents

Vapnik-Chervonenkis theory, simplest case

The support vector machine

Optimal margin classifier

Introducing kernels (1992)

Allowing for misclassification: soft margin (1995)

Implementation

Optimal margin hyperplane

- ▶ Let the training sample consist of n pairs $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$ taking values in $\mathbb{R}^d \times \{-1, 1\}$.
- ▶ This set is said to be **linearly separable** if there exist a non-zero vector $\boldsymbol{\psi} \in \mathbb{R}^d$ and a scalar $b \in \mathbb{R}$ such that the n following inequalities hold:

$$\begin{aligned}\boldsymbol{\psi}^T \mathbf{x}_i + b &\geq 0 && \text{if } y_i = 1, \\ \boldsymbol{\psi}^T \mathbf{x}_i + b &\leq 0 && \text{if } y_i = -1.\end{aligned}$$

- ▶ Instead of simply requiring separation (the parts “ ≥ 0 ” and “ ≤ 0 ” in the above inequality) one can **introduce margin** $M > 0$, *i.e.*, require the distance between any two points stemming from different classes — in projection onto $\boldsymbol{\psi}$ — be at least $2M$.
- ▶ Involving the output (in this notation corresponding to the sign) allows for rewriting the above (restricting) inequalities in the following way:

$$\frac{y_i(\boldsymbol{\psi}^T \mathbf{x}_i + b)}{\|\boldsymbol{\psi}\|} \geq M, \quad i = 1, \dots, n.$$

Optimal margin hyperplane

- ▶ The objective of the training algorithm is then to find the parameter vector ψ that maximizes M :

$$M^* = \max M$$

$$\text{subject to } \|\psi\| = 1,$$

$$y_i(\psi^T \mathbf{x}_i + b) \geq M, \quad i = 1, \dots, n.$$

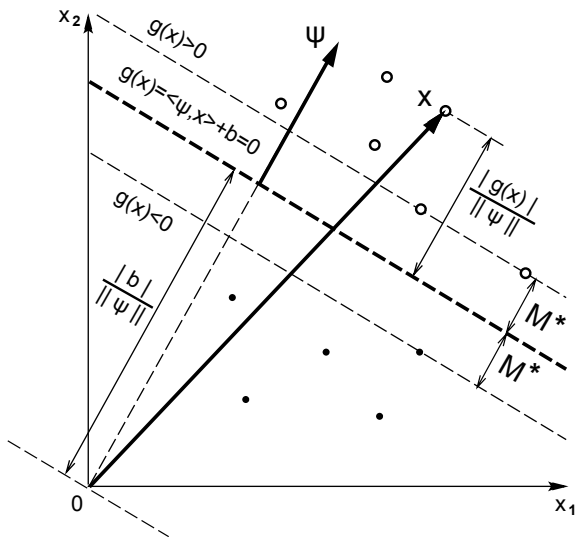
- ▶ The (last) bound is attained for those patterns satisfying

$$\min_{i \in \{1, \dots, n\}} y_i(\psi^T \mathbf{x}_i + b) = M^*.$$

- ▶ These patterns are called the **support vectors** of the decision boundary.
- ▶ Thus, the problem of finding a hyperplane with maximum margin can be seen as a **minimax** problem:

$$\max_{\psi \in \mathbb{R}^d, \|\psi\|=1} \min_{i \in \{1, \dots, n\}} y_i(\psi^T \mathbf{x}_i + b).$$

Optimal margin hyperplane : illustration



Optimal margin hyperplane

- ▶ Instead of fixing the norm of ψ , the product of the margin M and the norm of ψ can be fixed, e.g. by:

$$M\|\psi\| = 1.$$

- ▶ Now, **maximizing** the **margin** M is **equivalent to minimizing** the **norm** $\|\psi\|$.
- ▶ Then the problem of finding a maximum margin separating hyperplane, characterized by ψ , reduces to solving the following **quadratic optimization problem**:

$$\min \frac{1}{2}\|\psi\|^2$$

$$\text{subject to } y_i(\psi^T \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, n.$$

- ▶ The maximum margin is:

$$M^* = \frac{1}{\|\psi^*\|}.$$

- ▶ This approach is **impractical**:
 - if the **dimension** d is large or **infinite**,
 - because no information about **support vectors** is gained.

Optimal margin hyperplane: Lagrangian

- ▶ Construct a Lagrangian:

$$L(\boldsymbol{\psi}, b, \boldsymbol{\Lambda}) = \frac{1}{2} \boldsymbol{\psi}^T \boldsymbol{\psi} - \sum_{i=1}^n \alpha_i (y_i (\boldsymbol{\psi}^T \mathbf{x}_i + b) - 1)$$

with $\boldsymbol{\Lambda} = (\alpha_1, \dots, \alpha_n)^T$ being the vector of non-negative **Lagrange multipliers** corresponding to the inequality constraints.

- ▶ The solution to the optimization problem is determined by the saddle point of this Lagrangian in the $(d + 1 + n)$ -dimensional space of $\boldsymbol{\psi}$, b , and $\boldsymbol{\Lambda}$.
- ▶ The **minimum** should be taken w.r.t. the parameters $\boldsymbol{\psi}$ and b , the **maximum** should be taken w.r.t. the Lagrange multipliers $\boldsymbol{\Lambda}$.

Optimal margin hyperplane: Lagrangian

- ▶ At the point of minimum (w.r.t. ψ and b) one obtains:

$$\left. \frac{\partial L(\psi, b, \Lambda)}{\partial \psi} \right|_{\psi=\psi^*} = \left(\psi^* - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \right) = 0,$$

$$\left. \frac{\partial L(\psi, b, \Lambda)}{\partial b} \right|_{b=b^*} = \sum_{i=1}^n y_i \alpha_i = 0.$$

- ▶ From the upper equality one can derive:

$$\psi^* = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i.$$

- ▶ This means that the optimal hyperplane can be written as a **linear combination of training observations**.
- ▶ Only training observations \mathbf{x}_i with (strictly) positive Lagrange multipliers (*i.e.* with $\alpha_i > 0$) have an efficient contribution to the sum — the **support vectors**.

Optimal margin hyperplane: Lagrangian

- ▶ Substitution of the minimum conditions into the Lagrangian yields the following optimization problem:

$$\max W(\boldsymbol{\Lambda}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\text{subject to } \sum_{i=1}^n \alpha_i y_i = 0,$$

$$\alpha_i \geq 0, \quad i = 1, \dots, n.$$

- ▶ Usually it is written in the matrix form:

$$\max W(\boldsymbol{\Lambda}) = \boldsymbol{\Lambda}^T \mathbf{1} - \frac{1}{2} \boldsymbol{\Lambda}^T \mathbf{D} \boldsymbol{\Lambda}$$

$$\text{subject to } \boldsymbol{\Lambda}^T \mathbf{Y} = 0,$$

$$\boldsymbol{\Lambda} \geq \mathbf{0}$$

with \mathbf{D} being a $(n \times n)$ -dimensional matrix with entries $D_{ij} = y_i y_j \mathbf{x}_i^T \mathbf{x}_j$, $\mathbf{Y} = (y_1, \dots, y_n)^T$, and $\mathbf{0}$ and $\mathbf{1}$ standing for n -dimensional vectors of zeros and ones.

Optimal margin hyperplane: classification

- ▶ After the optimal pair $(\boldsymbol{\psi}^*, b^*)$ is obtained, classification of an observation $\mathbf{x} \in \mathbb{R}^d$ reduces to determining its position in the projection onto $\boldsymbol{\psi}^*$:

$$\begin{aligned}g(\mathbf{x}) &= \text{sign}(\boldsymbol{\psi}^{*T} \mathbf{x} + b^*) \\ &= \text{sign}\left(\sum_{i=1}^n y_i \alpha_i^* \mathbf{x}_i^T \mathbf{x} + b^*\right).\end{aligned}$$

- ▶ From this it becomes clear how to calculate b^* : it should position the separating hyperplane exactly in the middle between two support vectors from different classes, in the projection onto $\boldsymbol{\psi}^*$:

$$\begin{aligned}b^* &= -\frac{1}{2}(\boldsymbol{\psi}^{*T} \mathbf{x}_A + \boldsymbol{\psi}^{*T} \mathbf{x}_B) \\ &= -\frac{1}{2} \sum_{i=1}^n y_i \alpha_i^* (\mathbf{x}_i^T \mathbf{x}_A + \mathbf{x}_i^T \mathbf{x}_B).\end{aligned}$$

with $\mathbf{x}_A \in \{\mathbf{x}_i : y_i = 1, \alpha_i^* > 0, i = 1, \dots, n\}$ and $\mathbf{x}_B \in \{\mathbf{x}_i : y_i = -1, \alpha_i^* > 0, i = 1, \dots, n\}$.

- ▶ Only **support vectors** influence the classification rule.
(Analogy with a mine field on the front line between two enemies.)

Optimal margin classifier (algorithm)

Finding the optimal margin hyperplane (training)

Input: Training sample $((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)) \subset \mathbb{R}^d \times \{-1, 1\}$.

1. Solve the constraint quadratic optimization problem to obtain $\mathbf{\Lambda}^* = (\alpha_1^*, \dots, \alpha_n^*)^T$:

$$\begin{aligned} \max \quad & \mathbf{\Lambda}^T \mathbf{1} - \frac{1}{2} \mathbf{\Lambda}^T \mathbf{D} \mathbf{\Lambda} \\ \text{subject to} \quad & \mathbf{\Lambda}^T \mathbf{Y} = 0, \\ & \mathbf{\Lambda} \geq \mathbf{0}. \end{aligned}$$

2. Taking any two support vectors from opposite classes $\mathbf{x}_A \in \{\mathbf{x}_i : y_i = 1, \alpha_i^* > 0\}$ and $\mathbf{x}_B \in \{\mathbf{x}_i : y_i = -1, \alpha_i^* > 0\}$, calculate the threshold:

$$b^* = -\frac{1}{2} \sum_{i=1}^n y_i \alpha_i^* (\mathbf{x}_i^T \mathbf{x}_A + \mathbf{x}_i^T \mathbf{x}_B).$$

Output: The classifier: $g(\mathbf{x}) = \text{sign}\left(\sum_{i=1}^n y_i \alpha_i^* \mathbf{x}_i^T \mathbf{x} + b^*\right)$.

Optimal margin classifier: some comments

- ▶ The training phase is reduced to solving a problem of **quadratic optimization**, which is usually **computationally tractable**.
- ▶ The time of the training algorithm depends on dimension d only when calculating the matrix of quadratic coefficients D , the **dimension** of the original space is **irrelevant for the optimization time**.
- ▶ As **only** the **support vectors** are relevant, only these **should be stored** for the classification rule.
- ▶ The problem **can be solved iteratively by chunks**, as in each (previous) chunk only support vectors are important (for further chunks).

But:

- ▶ **Linear classification rule has poor approximation performance.**
- ▶ **Misclassification is not allowed.**

Contents

Vapnik-Chervonenkis theory, simplest case

The support vector machine

Optimal margin classifier

Introducing kernels (1992)

Allowing for misclassification: soft margin (1995)

Implementation

Convolution of the inner product

- ▶ The algorithm described above constructs a hyperplane (defining by that linear classification rule) in the input space \mathbb{R}^d .
- ▶ **Idea:** To increase the approximation power of the classification rule but keep its algorithmic linearity, one maps the input space to a feature space, *i.e.* transforms the d -dimensional input vector \mathbf{x} into a D -dimensional feature space through a choice of a D -dimensional vector function ϕ :

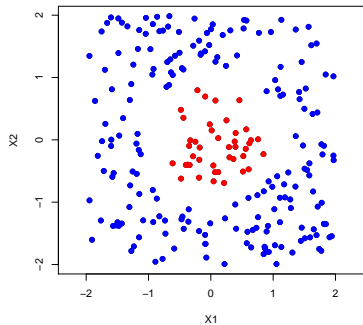
$$\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D .$$

- ▶ Then, a D -dimensional linear separator $(\psi, b) \in \mathbb{R}^D \times \mathbb{R}$ is constructed for the set of transformed vectors:

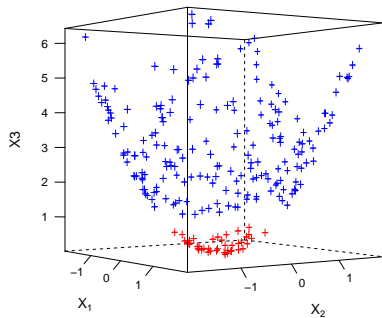
$$\phi(\mathbf{x}_i) = (\phi_1(\mathbf{x}_i), \phi_2(\mathbf{x}_i), \dots, \phi_D(\mathbf{x}_i)) \in \mathbb{R}^D, \quad i = 1, \dots, n.$$

- ▶ Note: \mathbb{R}^D can be of infinite dimension.

Convolution of the inner product



$$(x_1, x_2)$$



$$\begin{aligned} &\mapsto (\phi(x_1, x_2)) \\ &= (x_1, x_2, x_1^2 + x_2^2) \end{aligned}$$

Convolution of the inner product

- ▶ The classification of an unknown vector \mathbf{x} is done by first transforming it into the feature space

$$\mathbf{x} \mapsto \phi(\mathbf{x}),$$

and then classifying the featured vector with

$$g(\mathbf{x}) = \text{sign}(\psi^{*T} \phi(\mathbf{x}) + b^*).$$

- ▶ According to the properties of the classifier, it can be written as a linear combination of the support vectors (in the feature space):

$$\psi^* = \sum_{i=1}^n y_i \alpha_i^* \phi(\mathbf{x}_i).$$

- ▶ The linearity of the inner product implies that the **classifier** $g(\mathbf{x})$ **only depends on the inner products**:

$$g(\mathbf{x}) = \text{sign}(\psi^{*T} \phi(\mathbf{x}) + b^*) = \text{sign}\left(\sum_{i=1}^n y_i \alpha_i^* \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) + b^*\right).$$

- ▶ The **quadratic problem depends only inner products** as well.

Convolution of the inner product

- ▶ Consider the general form of the inner product in a Hilbert space:

$$\phi(\mathbf{u})^T \phi(\mathbf{v}) = K(\mathbf{u}, \mathbf{v}).$$

- ▶ According to Hilbert-Schmidt Theory any symmetric function $K(\mathbf{u}, \mathbf{v})$, with $K(\mathbf{u}, \mathbf{v}) \in L_2$, can be expanded in the form:

$$K(\mathbf{u}, \mathbf{v}) = \sum_{j=1}^{\infty} \lambda_j \phi_j(\mathbf{u}) \phi_j(\mathbf{v})$$

with $\lambda_i \in \mathbb{R}$ and ϕ_i being eigenvalues and eigenfunctions of the integral operator defined by the kernel $K(\mathbf{u}, \mathbf{v})$, *i.e.*

$$\int K(\mathbf{u}, \mathbf{v}) \phi_j(\mathbf{u}) d\mathbf{u} = \lambda_j \phi_j(\mathbf{v}).$$

- ▶ A sufficient condition to ensure that $K(\mathbf{u}, \mathbf{v})$ defines an inner product in the feature space is that all the eigenvalues λ_i are positive.

Convolution of the inner product

- ▶ According to Mercer's theorem, for λ_i s to be positive, it is necessary and sufficient that

$$\int \int K(\mathbf{u}, \mathbf{v}) h(\mathbf{u}) h(\mathbf{v}) d\mathbf{u} d\mathbf{v} > 0$$

holds for all h such that

$$\int h^2(\mathbf{u}) d\mathbf{u} < \infty.$$

- ▶ Thus, **functions that satisfy the Mercer's theorem can be used as inner products in the feature space.**

Examples of such functions:

- ▶ **Gaussian kernel = potential function = radial basis function:**

$$K(\mathbf{u}, \mathbf{v}) = e^{-\frac{\|\mathbf{u}-\mathbf{v}\|^2}{2\sigma^2}} = e^{-\gamma\|\mathbf{u}-\mathbf{v}\|^2}.$$

- ▶ **Polynomial kernel:**

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v} + 1)^\beta.$$

Convolution of the inner product

- ▶ Using **different kernel functions** $K(\mathbf{u}, \mathbf{v})$ as inner products (**with different parameters**, e.g., σ, γ, β) one can construct different learning machines with arbitrary types of decision surfaces.
- ▶ To find the optimal coefficient vector $\mathbf{\Lambda}^* = (\alpha_1^*, \dots, \alpha_n^*)$, threshold b^* , and support vectors \mathbf{x}_i s, one follows the same solution scheme as for the original optimal margin classifier by solving the quadratic optimization problem.
- ▶ The only difference consists in using the matrix \mathbf{D} with entries:

$$D_{ij} = y_i y_j K(\mathbf{x}_i, \mathbf{x}_j), \quad i, j = 1, \dots, n.$$

- ▶ The **decision rule** has then form:

$$g(\mathbf{x}) = \sum_{i=1}^n y_i \alpha_i^* K(\mathbf{x}, \mathbf{x}_i) + b^* .$$

where one can only restrict to support vectors \mathbf{x}_i and their coefficients $\alpha_i^* > 0$.

Contents

Vapnik-Chervonenkis theory, simplest case

The support vector machine

Optimal margin classifier

Introducing kernels (1992)

Allowing for misclassification: soft margin (1995)

Implementation

Soft margin classifier

- ▶ The **kernel trick** allows to “ignore” the transform on the algorithmic level.
- ▶ Consider the case where the **training data cannot be separated** without error.
- ▶ In this case one may want to **separate** the training set **with a minimal number of errors**.
- ▶ Let us introduce non-negative variables $\xi_i \geq 0$, $i = 1, \dots, n$.
- ▶ We can then **minimize** the functional

$$\sum_{i=1}^n \xi_i^\sigma$$

for some small $\sigma > 0$ subject to constraints

$$\begin{aligned} y_i(\boldsymbol{\psi}^T \mathbf{x}_i + b) &\geq 1 - \xi_i, & i = 1, \dots, n, \\ \xi_i &\geq 0, & i = 1, \dots, n. \end{aligned}$$

- ▶ For sufficiently small σ the minimized functional describes the number of errors on the training set.

Soft margin classifier

- ▶ In the minimum, strictly positive $\xi_j > 0$, $j = 1, \dots, k$ will identify the **minimal subset of training errors**:

$$(\mathbf{x}_{i_1}, y_{i_1}), (\mathbf{x}_{i_2}, y_{i_2}), \dots, (\mathbf{x}_{i_k}, y_{i_k}).$$

- ▶ After these data are **excluded**, one can separate the remaining part of the training set without errors using the usual optimal separating hyperplane.
- ▶ Formally this can be expressed as:

$$\min \quad \frac{1}{2} \|\boldsymbol{\psi}\|^2 + CF \left(\sum_{i=1}^n \xi_i^\sigma \right)$$

$$\begin{aligned} \text{subject to} \quad & y_i(\boldsymbol{\psi}^T \mathbf{x}_i + b) \geq 1 - \xi_i, & i = 1, \dots, n, \\ & \xi_i \geq 0, & i = 1, \dots, n. \end{aligned}$$

with $F(u)$ being a monotonic convex function and C being a positive constant.

- ▶ For sufficiently large C and sufficiently small σ , the pair $(\boldsymbol{\psi}^*, b^*)$ minimizing this functional will determine the **hyperplane minimizing the number of errors and separating the rest with maximum margin**.

Soft margin classifier

- ▶ However, the problem of finding a hyperplane minimizing number of errors is **NP-complete**.
- ▶ For the reasons of computational tractability, we consider the (most commonly used) case:

$$F(u) = u,$$
$$\sigma = 1,$$

and choose appropriate value for the regularizing constant C .

- ▶ The problem then becomes:

$$\begin{aligned} \min \quad & \frac{1}{2} \|\psi\|^2 + C \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & y_i(\psi^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, n, \\ & \xi_i \geq 0, \quad i = 1, \dots, n. \end{aligned}$$

Soft margin classifier

- ▶ The corresponding Lagrangian is:

$$L(\boldsymbol{\psi}, b, \boldsymbol{\xi}, \boldsymbol{\Lambda}, \mathbf{r}) = \frac{1}{2} \boldsymbol{\psi}^T \boldsymbol{\psi} + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i (y_i (\boldsymbol{\psi}^T \mathbf{x}_i + b) - 1 + \xi_i) - \sum_{i=1}^n r_i \xi_i$$

with $\boldsymbol{\Lambda} = (\alpha_1, \dots, \alpha_n)^T$ and $\mathbf{r} = (r_1, \dots, r_n)^T$ being the vectors of non-negative **Lagrange multipliers** corresponding to the two groups of inequality constraints.

- ▶ The solution to the optimization problem is determined by the saddle point of this Lagrangian in the $(d + 1 + n + n + n)$ -dimensional space of $\boldsymbol{\psi}$, b , $\boldsymbol{\xi}$, $\boldsymbol{\Lambda}$, and \mathbf{r} .
- ▶ The **minimum** should be taken w.r.t. the parameters $\boldsymbol{\psi}$, b , and $\boldsymbol{\xi}$, the **maximum** should be taken w.r.t. the Lagrange multipliers $\boldsymbol{\Lambda}$ and \mathbf{r} .

Soft margin classifier

- ▶ At the point of minimum (w.r.t. ψ , b , and ξ) one obtains:

$$\left. \frac{\partial L(\psi, b, \xi, \mathbf{\Lambda}, \mathbf{r})}{\partial \psi} \right|_{\psi=\psi^*} = \left(\psi^* - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \right) = 0,$$

$$\left. \frac{\partial L(\psi, b, \xi, \mathbf{\Lambda}, \mathbf{r})}{\partial b} \right|_{b=b^*} = \sum_{i=1}^n y_i \alpha_i = 0,$$

$$\left. \frac{\partial L(\psi, b, \xi, \mathbf{\Lambda}, \mathbf{r})}{\partial \xi_i} \right|_{\xi_i=\xi_i^*} = C - \alpha_i - r_i = 0, \quad i = 1, \dots, n.$$

- ▶ This leads to the following quadratic problem:

$$\max W(\mathbf{\Lambda}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\text{subject to} \quad \sum_{i=1}^n \alpha_i y_i = 0,$$

$$0 \leq \alpha_i \leq C, \quad i = 1, \dots, n.$$

Support vector machine (SVM)

- ▶ The **training** phase:

$$\begin{aligned} \max \quad & \boldsymbol{\Lambda}^T \mathbf{1} - \frac{1}{2} \boldsymbol{\Lambda}^T \mathbf{D} \boldsymbol{\Lambda} \\ \text{subject to} \quad & \boldsymbol{\Lambda}^T \mathbf{Y} = 0, \\ & \mathbf{0} \leq \boldsymbol{\Lambda} \leq C \mathbf{1}, \end{aligned}$$

with $\mathbf{Y} = (y_1, \dots, y_n)^T$, $\mathbf{0}$ and $\mathbf{1}$ standing for n -dimensional vectors of zeros and ones, C being a properly chosen constant, and \mathbf{D} being a $(n \times n)$ -dimensional matrix with entries

$$D_{ij} = y_i y_j K(\mathbf{x}_i, \mathbf{x}_j), \quad i, j = 1, \dots, n,$$

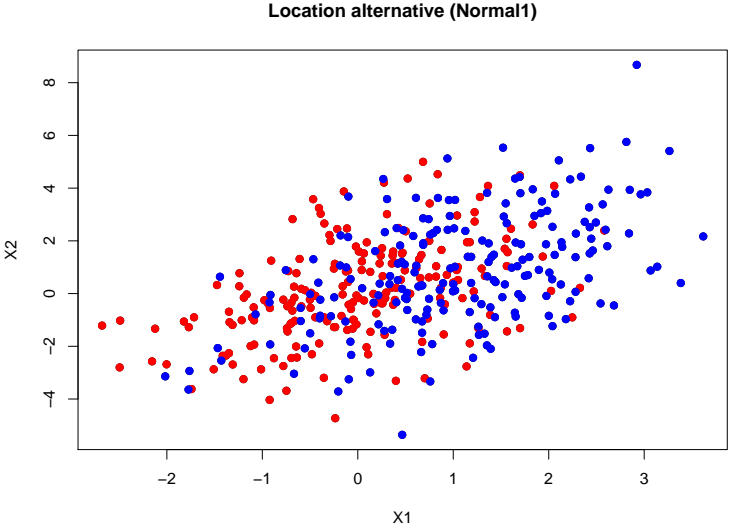
where $K(\mathbf{u}, \mathbf{v})$ is a properly chosen kernel function.

The result is the optimal vector $\boldsymbol{\Lambda}^* = (\alpha_1^*, \dots, \alpha_n^*)^T$.

Then, taking any two support vectors \mathbf{x}_{i_A} and \mathbf{x}_{i_B} from opposite classes, *i.e.* with $i_A \in \arg \max_{j: y_j=1, \alpha_j^* > 0} \sum_{i=1}^n y_i \alpha_i^* K(\mathbf{x}_j, \mathbf{x}_i)$ and $i_B \in \arg \min_{j: y_j=-1, \alpha_j^* > 0} \sum_{i=1}^n y_i \alpha_i^* K(\mathbf{x}_j, \mathbf{x}_i)$, calculate threshold:

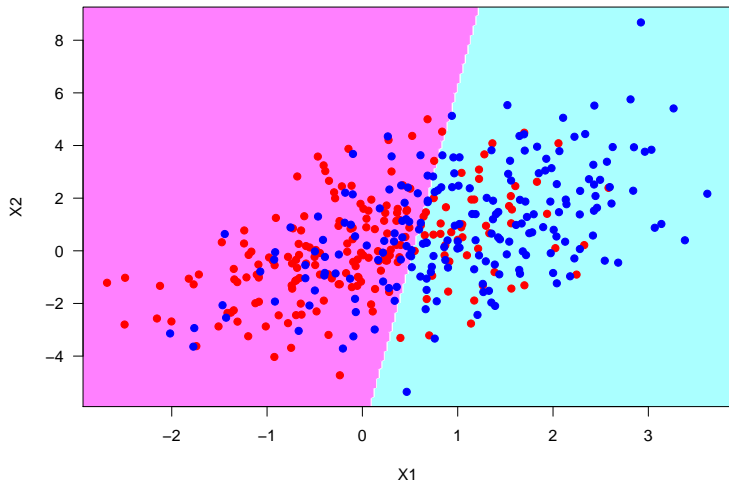
$$b^* = -\frac{1}{2} \sum_{i=1}^n y_i \alpha_i^* (K(\mathbf{x}_i, \mathbf{x}_{i_A}) + K(\mathbf{x}_i^T \mathbf{x}_{i_B})).$$

Normal location alternative

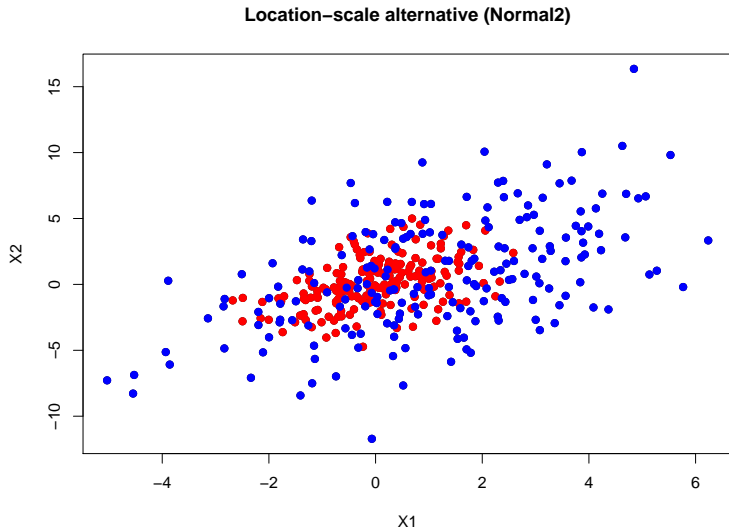


SVM: normal location alternative

SVM (linear kernel) for Normal1 data

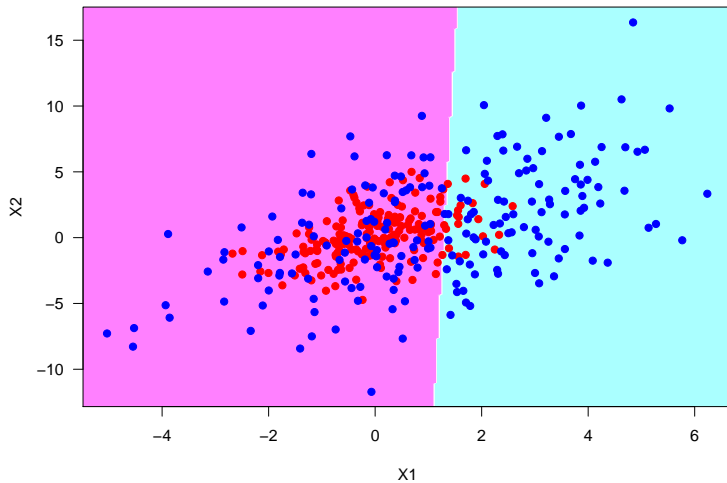


Normal location-scale alternative



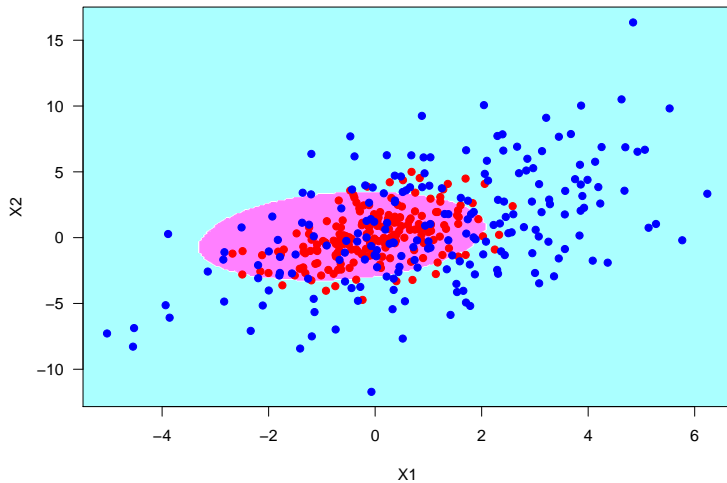
SVM: normal location-scale alternative

SVM (linear kernel) for Normal2 data



SVM: normal location-scale alternative

SVM (radial kernel) for Normal2 data



Contents

Vapnik-Chervonenkis theory, simplest case

The support vector machine

- Optimal margin classifier

- Introducing kernels (1992)

- Allowing for misclassification: soft margin (1995)

Implementation

Implementing SVM

When implementing and applying SVM, its parameters have to be chosen:

- ▶ **kernel** function,
- ▶ **kernel parameter**,
- ▶ **regularization constant** (=box constraint).

In practice, these parameters are usually chosen by cross-validation. This process is called **tuning of the SVM**. The SVM possesses certain degree of insensitivity w.r.t. parameters, which can be limited depending on the application of interest.

For R-software, SVM is implemented in such packages as, e.g., `e1071`, `kernlab`, `klaR`, `svmpath`.

For an overview, see, e.g.:

- ▶ Karatzoglou, A., Meyer, D., and Hornik, K. (2006).
Support vector machines in R.
Journal of Statistical Software, 15(9).

Thank you for your attention!

Thank you for your attention!

And some references

- ▶ Hastie, T., Tibshirani, R., and Friedman, J. (2009).
The Elements of Statistics Learning: Data Mining, Inference, and Prediction (Second Edition).
Springer.
- ▶ Devroye, L., Gyöfri, L., Lugosi, G. (1996).
A Probabilistic Theory of Pattern Recognition.
Springer.
- ▶ Vapnik, V. N. (1998).
Statistical Learning Theory.
John Wiley & Sons.
- ▶ Haykin, S. (2009).
Neural Networks and Learning Machines (Third Edition).
Pearson.