

MSc Big Data for Business - *MAP 534*

Introduction to supervised learning

Supervised classification

Linear/Quadratic discriminant analysis (LDA/QDA) & Support Vector Machines (SVM)

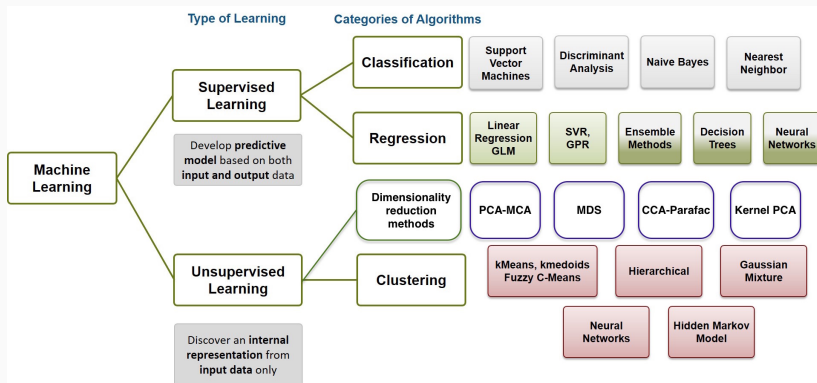
Unsupervised learning

- Dimensionality reduction methods (PCA).
- Clustering (k -means, Hierarchical clustering).
- Handling missing values/ matrix completion.
⇒ Data visualization - exploratory data analysis.

Supervised learning

- Theoretical framework - Bayes risk.
- Logistic regression.
- Optimization.
- Nonparametric methods.

Machine Learning



Unsupervised Learning

- Goal: Discover a structure within a set of individuals (\mathbf{X}_i).
- Data: Learning set (\mathbf{X}_i)

Supervised Learning:

- Goal: Learn a function f predicting a variable Y from an individual \mathbf{X} .
- Data: Learning set (\mathbf{X}_i, Y_i)

Second case is better posed.

Introduction

Statistical Supervised Learning

Parametric models

- Linear/quadratic discriminant analysis

- Support Vector Machine

Introduction

Statistical Supervised Learning

Parametric models

Linear/quadratic discriminant analysis

Support Vector Machine

Use data to extract a prediction function

Search engine, text-mining.

Diagnosis, Fault detection.

Business analytics.

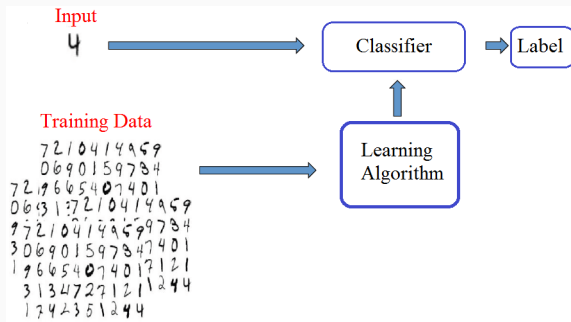
Social network analytics.

Data-mining.

Link prediction.

Robotics.

A definition of Machine Learning

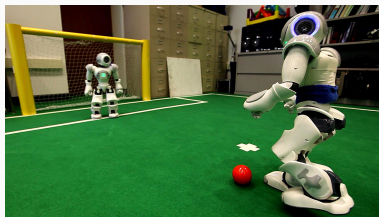


A definition by Tom Mitchell (<http://www.cs.cmu.edu/~tom/>)

A computer program is said to learn from **experience E** with respect to some **class of tasks T** and **performance measure P**, if its performance at tasks in T, as measured by P, improves with experience E.

A robot that learns

A robot endowed with a set of sensors and an online learning algorithm:

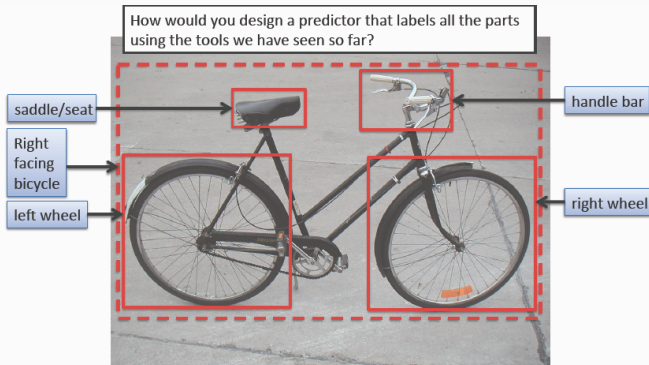


Task: play football.

Performance: score.

Experience: current environment and outcome, past games...

Object recognition in an image

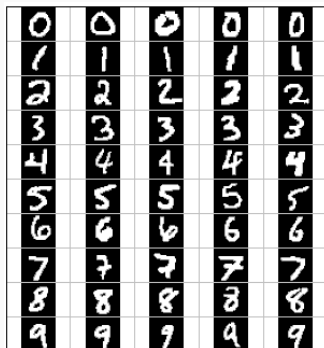


Task: say if an object is present or not in the image.

Performance: number of errors.

Experience: set of previously seen labeled image.

Number



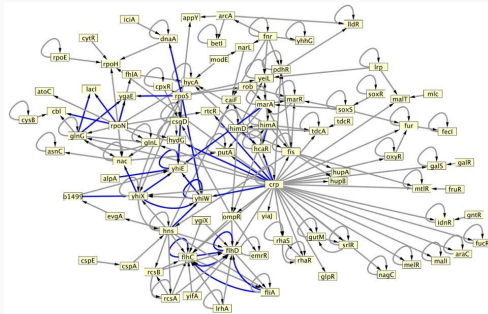
Task: Read a ZIP code from an envelop.

Goal: give a number from an image.

X = image.

Y = corresponding number.

Applications in biology



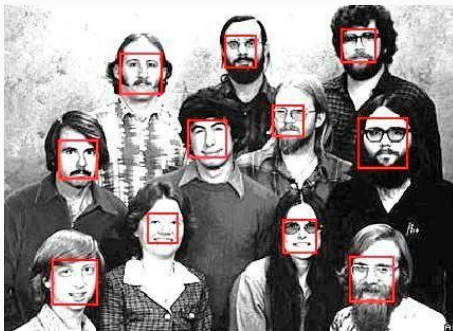
Task: protein interaction network prediction.

Goal: predict (unknown) interactions between proteins.

X = pair of proteins.

Y = existence or no of interaction.

Numerous similar questions in bio(informatics), genomic, etc.



Goal: detect the position of faces in an image.

Different setting?... Reformulation as a supervised learning problem.

X = sub window in the image

Y = presence or no of a face...

Lots of answer in a single image. Post processing required...

Support Vector Machine

Linear Discriminant Analysis

Logistic Regression

Trees/ Random Forests

Kernel methods

Neural Networks

Many more...

Introduction

Statistical Supervised Learning

Parametric models

Linear/quadratic discriminant analysis

Support Vector Machine

Supervised Learning Framework

Input measurement $\mathbf{X} \in \mathcal{X}$ (often $\mathcal{X} \subset \mathbb{R}^d$).

Output measurement $Y \in \mathcal{Y}$.

The joint distribution of (\mathbf{X}, Y) is p with p unknown.

Training data : $\mathcal{D}_n = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}$ (i.i.d. $\sim p$).

$Y \in \{-1, 1\}$ (classification) or $Y \in \mathbb{R}$ (regression).

A **predictor** is a measurable function in $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$.

Goal

Construct a **good** predictor \hat{f}_n from the training data.

Need to specify the meaning of good.

Loss function for a generic predictor

- **Loss function:** $\ell(Y, f(\mathbf{X}))$ measures the goodness of the prediction of Y by $f(\mathbf{X})$.
- Examples:
 - Prediction loss: $\ell(Y, f(\mathbf{X})) = \mathbf{1}_{Y \neq f(\mathbf{X})}$.
 - Quadratic loss: $\ell(Y, \mathbf{X}) = |Y - f(\mathbf{X})|^2$.

Risk function

- Risk measured as the average loss for a new couple:

$$\mathcal{R}(f) = \mathbb{E}_{(X, Y) \sim p} [\ell(Y, f(\mathbf{X}))].$$

- Examples:
 - Prediction loss: $\mathbb{E} [\ell(Y, f(\mathbf{X}))] = \mathbb{P} \{Y \neq f(\mathbf{X})\}$.
 - Quadratic loss: $\mathbb{E} [\ell(Y, f(\mathbf{X}))] = \mathbb{E} [|Y - f(\mathbf{X})|^2]$.

- **Beware:** As \hat{f}_n depends on \mathcal{D}_n , $\mathcal{R}(\hat{f}_n)$ is a random variable!

The best solution f^* (which is independent of \mathcal{D}_n) is

$$f^* = \arg \min_{f \in \mathcal{F}} R(f) = \arg \min_{f \in \mathcal{F}} \mathbb{E} [\ell(Y, f(\mathbf{X}))] .$$

Bayes Predictor (explicit solution)

In binary classification with 0 – 1 loss:

$$f^*(\mathbf{X}) = \begin{cases} +1 & \text{if } \mathbb{P}\{Y = 1|\mathbf{X}\} \geq \mathbb{P}\{Y = -1|\mathbf{X}\} \\ & \Leftrightarrow \mathbb{P}\{Y = 1|\mathbf{X}\} \geq 1/2, \\ -1 & \text{otherwise.} \end{cases}$$

In regression with the quadratic loss

$$f^*(\mathbf{X}) = \mathbb{E}[Y|\mathbf{X}] .$$

Issue: the explicit solution requires to know $\mathbb{E}[Y|\mathbf{X}]$ for all values of \mathbf{X} !

Estimate $\eta(\mathbf{x}) = \mathbb{P}[Y = 1 | \mathbf{X} = \mathbf{x}]$ by $\hat{\eta}_n(\mathbf{x})$ using the training dataset and plug it the Bayes classifier.

Plugin Bayes Classifier

In binary classification with 0 – 1 loss:

$$\hat{f}_n(\mathbf{X}) = \begin{cases} +1 & \text{if } \hat{\eta}_n(\mathbf{x}) \geq 1/2, \\ -1 & \text{otherwise.} \end{cases}$$

Input: a data set \mathcal{D}_n .

Learn $Y|\mathbf{X}$ or equivalently $\mathbb{P}\{Y = k|\mathbf{X}\}$ (using the data set) and plug this estimate in the Bayes classifier.

Output: a classifier $\hat{f}_n : \mathbb{R}^d \rightarrow \{-1, 1\}$

$$\hat{f}_n(\mathbf{X}) = \begin{cases} +1 & \text{if } \hat{\eta}_n(\mathbf{X}) \geq 1/2, \\ -1 & \text{otherwise.} \end{cases}$$

Can we certify that the classifier is good if $Y|\mathbf{X}$ is well estimated?

The missclassification error satisfies:

$$0 \leq \mathbb{P}(\hat{f}_n(\mathbf{X}) \neq Y) - L^* \leq 2\mathbb{E} [|\eta(\mathbf{X}) - \hat{\eta}_n(\mathbf{X})|^2]^{1/2},$$

where

$$L^* = \mathbb{P}[f^*(\mathbf{X}) \neq Y]$$

and $\hat{\eta}_n(\mathbf{x})$ is an empirical estimate based on the training dataset of

$$\eta(\mathbf{x}) = \mathbb{P}[Y = 1 | \mathbf{X} = \mathbf{x}].$$

Instantiation: how to estimate $Y|X$?

Fully generative modeling.

Estimate the law of (\mathbf{X}, Y) and use the **Bayes formula** to deduce an estimate of $Y|\mathbf{X}$: *LDA/QDA, Naive Bayes, Gaussian Processes...*

Parametric conditional modeling.

Estimate the law of $Y|\mathbf{X}$ by a **parametric** law $\mathcal{L}_\theta(\mathbf{X})$: *linear regression, logistic regression...*

Nonparametric conditional modeling.

Estimate the law of $Y|\mathbf{X}$ by a **non parametric** estimate: *kernel methods, nearest neighbors...*

Introduction

Statistical Supervised Learning

Parametric models

- Linear/quadratic discriminant analysis

- Support Vector Machine

Introduction

Statistical Supervised Learning

Parametric models

Linear/quadratic discriminant analysis

Support Vector Machine

If the law of (\mathbf{X}, Y) is **known** everything can be easy!

Bayes formula

With a slight abuse of notation,

$$\mathbb{P}\{Y|\mathbf{X}\} = \frac{\mathbb{P}\{(\mathbf{X}, Y)\}}{\mathbb{P}\{\mathbf{X}\}} = \frac{\mathbb{P}\{\mathbf{X}|Y\} \mathbb{P}\{Y\}}{\mathbb{P}\{\mathbf{X}\}}$$

Generative Modeling

Propose a model for (\mathbf{X}, Y) (or equivalently $\mathbf{X}|Y$ and Y).

Estimate it as a density estimation problem.

Plug the estimate in the Bayes formula.

Plug the conditional estimate in the Bayes *classifier*.

Remark: require to estimate (\mathbf{X}, Y) rather than only $Y|\mathbf{X}$!

Great flexibility in the model design but may lead to complex computation.

Simplest setting in classification!

Bayes formula

$$\mathbb{P}\{Y = k|\mathbf{X}\} = \frac{\mathbb{P}\{\mathbf{X}|Y = k\} \mathbb{P}\{Y = k\}}{\mathbb{P}\{\mathbf{X}\}}.$$

Binary Bayes classifier (the best solution).

$$f^*(\mathbf{X}) = \begin{cases} +1 & \text{if } \mathbb{P}\{Y = 1|\mathbf{X}\} \geq \mathbb{P}\{Y = -1|\mathbf{X}\}, \\ -1 & \text{otherwise.} \end{cases}$$

Heuristic: estimate those quantities and plug the estimations. Choosing different models/estimators for $\mathbb{P}\{\mathbf{X}|Y\}$ leads to different classifiers.

Remark: no need to renormalize by $\mathbb{P}\{\mathbf{X}\}$ to take the decision!

Naive Bayes

Classical algorithm using a crude modeling for $\mathbb{P}\{\mathbf{X}|Y\}$:

- Feature **independence** assumption:

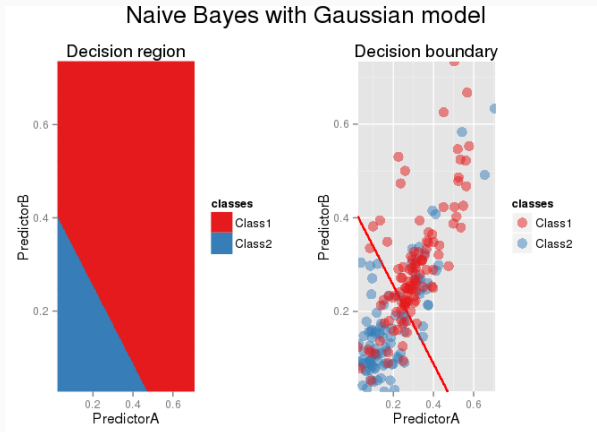
$$\mathbb{P}\{\mathbf{X}|Y\} = \prod_{i=1}^d \mathbb{P}\{X^{(i)}|Y\} .$$

- Simple featurewise model: binomial if binary, multinomial if finite and Gaussian if continuous.

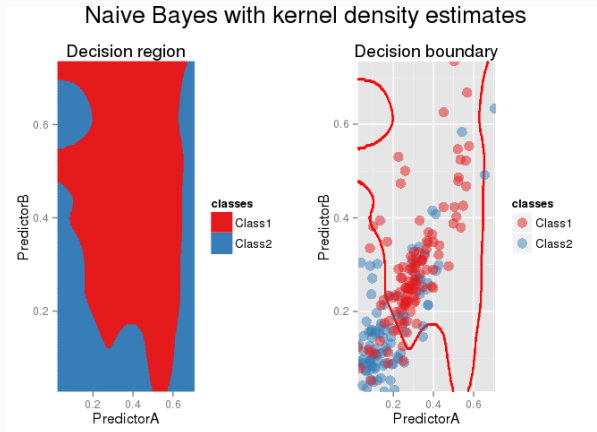
If all features are continuous, similar to the previous Gaussian but with a **diagonal covariance matrix**!

Very simple learning even in **very high dimension**!

Example: Naive Bayes



Example: Naive Bayes



Discriminant Analysis (Gaussian model)

The densities are modeled as multivariate normal, i.e., for all class k , conditionnally on $\{Y = k\}$,

$$\mathbf{X} \sim \mathcal{N}(\mu_k, \Sigma_k).$$

Discriminant functions:

$$g_k(\mathbf{X}) = \ln(\mathbb{P}\{\mathbf{X}|Y = k\}) + \ln(\mathbb{P}\{Y = k\}).$$

In a two-classes problem, the optimal classifier is

$$f^* : x \mapsto 2\mathbb{1}\{g_1(x) > g_{-1}(x)\} - 1.$$

QDA (different Σ_k in each class) and LDA ($\Sigma_k = \Sigma$ for all k)

Beware: this model can be false but the methodology remains valid!

Estimation

In practice, we will need to estimate μ_k , Σ_k and $\pi_k := \mathbb{P}\{Y = k\}$.

Estimated proportions $\hat{\pi}_k = \frac{n_k}{n} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{Y_i=k\}}$.

Maximum likelihood estimate of $\hat{\mu}_k$ and $\hat{\Sigma}_k$ (explicit formulas).

The DA classifier then becomes

$$\hat{f}_G(\mathbf{X}) = \begin{cases} +1 & \text{if } \hat{g}_{+1}(\mathbf{X}) \geq \hat{g}_{-1}(\mathbf{X}), \\ -1 & \text{otherwise.} \end{cases}$$

If $\Sigma_{-1} = \Sigma_1 = \Sigma$ then the **decision boundary is an affine hyperplane**.

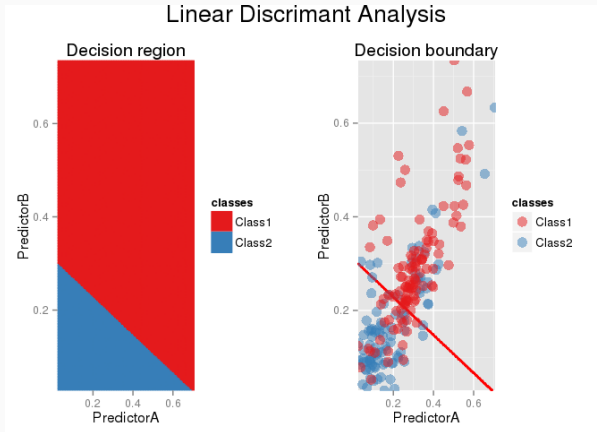
The loglikelihood of the observations is given by

$$\begin{aligned}\log \mathbb{P}_{\theta} (X_{1:n}, Y_{1:n}) &= \sum_{i=1}^n \log \mathbb{P}_{\theta} (X_i, Y_i) , \\ &= -\frac{nd}{2} \log(2\pi) - \frac{n}{2} \log \det(\Sigma) + \left(\sum_{i=1}^n \mathbb{1}_{Y_i=1} \right) \log \pi_1 + \left(\sum_{i=1}^n \mathbb{1}_{Y_i=-1} \right) \log(1 - \pi_1) \\ &\quad - \frac{1}{2} \sum_{i=1}^n \mathbb{1}_{Y_i=1} (X_i - \mu_1)' \Sigma^{-1} (X_i - \mu_1) - \frac{1}{2} \sum_{i=1}^n \mathbb{1}_{Y_i=-1} (X_i - \mu_{-1})' \Sigma^{-1} (X_i - \mu_{-1}) .\end{aligned}$$

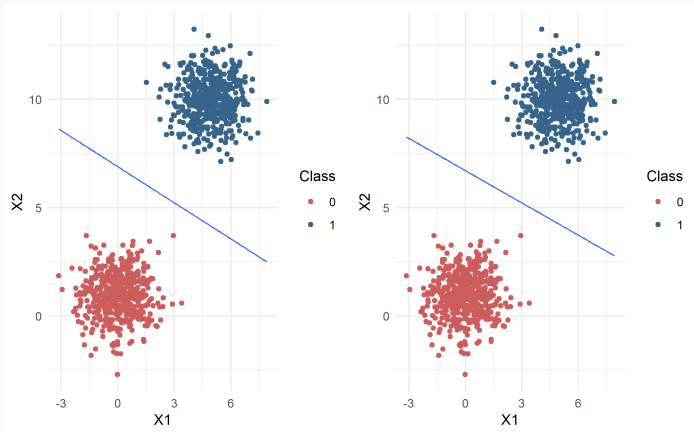
This yields, for $k \in \{-1, 1\}$,

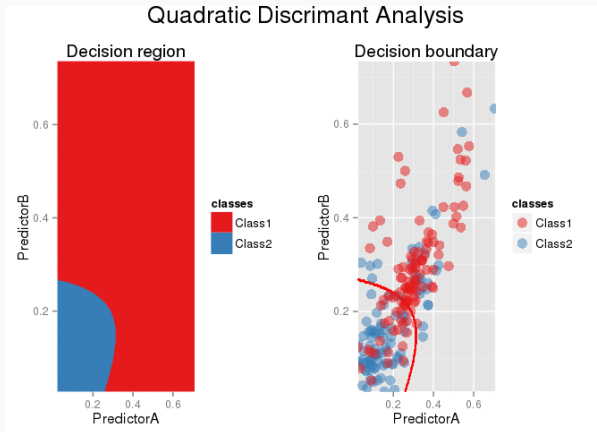
$$\begin{aligned}\hat{\pi}_k^n &= \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i=k} , \\ \hat{\mu}_k^n &= \frac{1}{\sum_{i=1}^n \mathbb{1}_{Y_i=k}} \sum_{i=1}^n \mathbb{1}_{Y_i=k} X_i , \\ \hat{\Sigma}^n &= \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}_{Y_i}^n) (X_i - \hat{\mu}_{Y_i}^n)' .\end{aligned}$$

Remains to plug these estimates in the classification boundary.

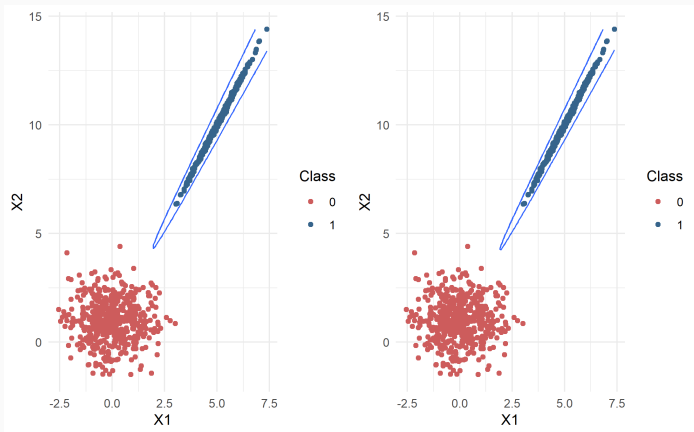


Example: LDA





Example: QDA



Function `svm` in package `e1071`.

Function `lda` and `qda` in package `MASS`.

Function `naive_bayes` in package `naivebayes`.

Introduction

Statistical Supervised Learning

Parametric models

Linear/quadratic discriminant analysis

Support Vector Machine

Setting

You have past/historical data from your platform, containing data about individuals $i = 1, \dots, n$.

You have a **features** vector $\mathbf{x}_i \in \mathbb{R}^d$ for each individual i .

For each i , you know if he/she clicked ($y_i = 1$) or not ($y_i = -1$).

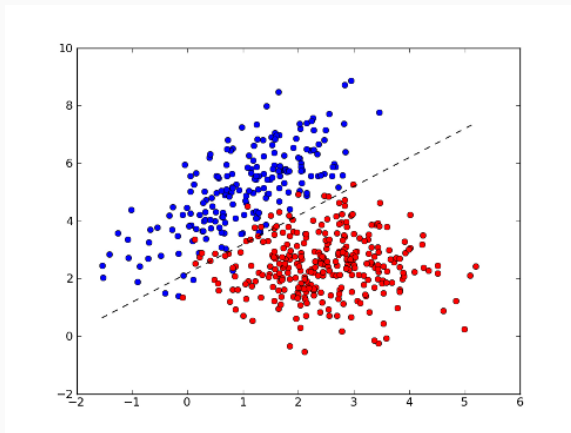
We call $y_i \in \{-1, 1\}$ the **label** of i .

Aim

Given a features vector \mathbf{x} (with no corresponding label), predict a label $\hat{y} \in \{-1, 1\}$.

Use data $\mathcal{D}_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ to construct a **classifier**.

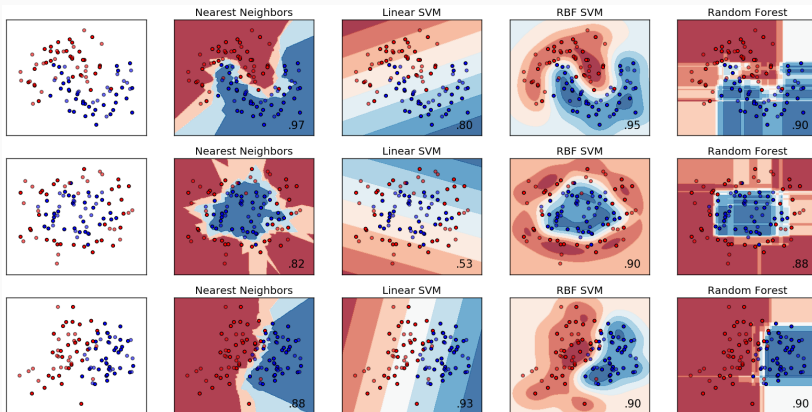
Geometrically



Learn a way to separate two “groups” of points

Classification

But, many ways to separate points!



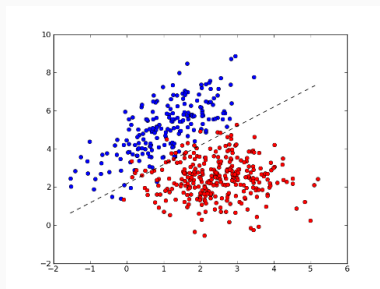
Linear classification

Simple to interpret and to implement.

On very large datasets (n is large, say $n \geq 10^6$), no other choice (training complexity).

Big data paradigm: lots of data \Rightarrow simple methods are enough.

A linear classifier



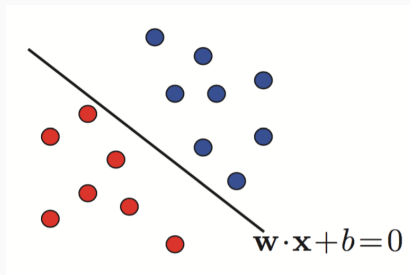
Learn $\hat{\mathbf{w}} \in \mathbb{R}^d$ and \hat{b} such that
$$\hat{y} = \text{sign}(\langle \mathbf{x}, \hat{\mathbf{w}} \rangle + \hat{b})$$

is a good classifier

Linearly separable data

A dataset is **linearly separable** if we can find an hyperplane H (linear classification rule) that puts

- Points $\mathbf{x}_i \in \mathbb{R}^d$ such that $y_i = 1$ on one side of the hyperplane
- Points $\mathbf{x}_i \in \mathbb{R}^d$ such that $y_i = -1$ on the other
- H do not pass through a point \mathbf{x}_i



A hyperplane

$$H_{\mathbf{w},\mathbf{x}} = \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{w}, \mathbf{x} \rangle + b = 0\}$$

is a translation of a set of vectors orthogonal to \mathbf{w} .

- $\mathbf{w} \in \mathbb{R}^d$ is a **non-zero vector normal** to the hyperplane.
- $b \in \mathbb{R}$ is a scalar.

Following for instance the results obtained for linear discriminant analysis and logistic regression, a **hyperplane $H_{\mathbf{w},b}$ may be used as a classifier** by defining

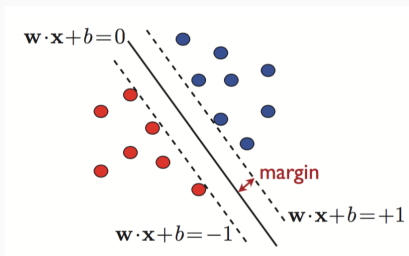
$$h_{\mathbf{w},b} : \mathbf{x} \mapsto \begin{cases} 1 & \text{if } \langle \mathbf{w}; \mathbf{x} \rangle + b > 0, \\ -1 & \text{otherwise.} \end{cases}$$

Some geometry

Definition of H is invariant by multiplication of \mathbf{w} and b by a non-zero scalar

If H does not pass through any sample point \mathbf{x}_i , we can scale \mathbf{w} and b so that

$$\min_{(\mathbf{x}, y) \in \mathcal{D}_n} |\langle \mathbf{w}, \mathbf{x} \rangle + b| = 1$$



For such \mathbf{w} and b , we call H the *canonical* hyperplane

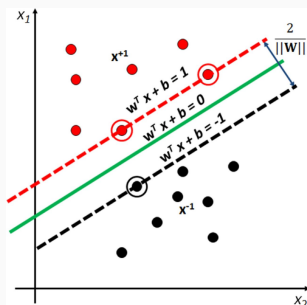
Some geometry

The distance of any point $\mathbf{x}' \in \mathbb{R}^d$ to H is given by

$$\frac{|\langle \mathbf{w}, \mathbf{x}' \rangle + b|}{\|\mathbf{w}\|}$$

So, if H is a canonical hyperplane, its **margin** is given by

$$\min_{(\mathbf{x}, y) \in \mathcal{D}_n} \frac{|\langle \mathbf{w}, \mathbf{x} \rangle + b|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}.$$



Summary

If \mathcal{D}_n is strictly linearly separable, we can find a canonical separating hyperplane

$$H = \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{w}, \mathbf{x} \rangle + b = 0\}.$$

that satisfies

$$|\langle \mathbf{w}, \mathbf{x}_i \rangle + b| \geq 1 \text{ for any } i = 1, \dots, n,$$

which entails that a point \mathbf{x}_i is correctly classified if

$$y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) \geq 1.$$

The margin of H is equal to $1/\|\mathbf{w}\|$.

Hard Support Vector Machines is a classification procedure which aims at building a linear classifier with the largest possible margin, i.e. the largest minimal distance between a point in the training set and the hyperplane.

The hyperplane which **correctly separates all training data sets with the largest margin** is obtained with:

$$(\hat{w}_n, \hat{b}_n) \in \underset{\substack{(w,b) \in \mathbb{R}^d \times \mathbb{R}^d; \|w\|=1, \\ \forall i \in \{1, \dots, n\}, Y_i(\langle w; X_i \rangle + b) > 0}}{\operatorname{argmax}} \left\{ \min_{1 \leq i \leq n} |\langle w; X_i \rangle + b| \right\}.$$

The **hard Support Vector Machines** procedure is equivalent to solving the following optimization problem:

$$(\hat{w}_n, \hat{b}_n) \in \operatorname{argmax}_{(w,b) \in \mathbb{R}^d \times \mathbb{R}; \|w\|=1} \left\{ \min_{1 \leq i \leq n} Y_i (\langle w; X_i \rangle + b) \right\},$$

A solution to the hard Support Vector Machines optimization problem is obtained by setting $(\hat{w}_n, \hat{b}_n) = (w_\star / \|w_\star\|, b_\star / \|w_\star\|)$ where

$$(w_\star, b_\star) \in \operatorname{argmin}_{\substack{(w,b) \in \mathbb{R}^d \times \mathbb{R} \\ \forall i \in \{1, \dots, n\}, Y_i (\langle w; X_i \rangle + b) \geq 1}} \|w\|^2.$$

A way of classifying \mathcal{D}_n with maximum margin is to solve the following problem:

$$(w_*, b_*) \in \underset{\substack{(w,b) \in \mathbb{R}^d \times \mathbb{R} \\ \forall i \in \{1, \dots, n\}, Y_i(\langle w; X_i \rangle + b) \geq 1}}{\operatorname{argmin}} \|w\|^2.$$

This problem admits a **unique** solution

It is a **quadratic programming** problem, which is easy to solve numerically.

Dedicated optimization algorithms can solve this on a large scale very efficiently.

Linear SVM: separable case

The optimization problem is solved using [Karush-Kuhn-Tucker's theorem](#)).

There are $\alpha_i \geq 0$, $i = 1, \dots, n$, called [dual variables](#), such that the solution (\mathbf{w}, b) of this problem satisfies:

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \quad \text{and} \quad \alpha_i ((y_i \langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1) = 0 \quad \text{for } i = 1, \dots, n.$$

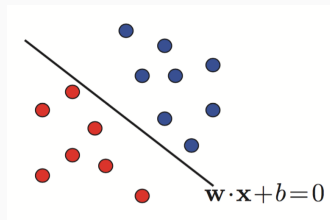
$\alpha_i \neq 0$ iff $y_i \langle \mathbf{w}, \mathbf{x}_i \rangle + b = 1$, meaning that \mathbf{x}_i is on the marginal hyperplane.

Weights vector \mathbf{w} is a linear combination of the vectors \mathbf{x}_i that belong to a marginal hyperplane.

Such points \mathbf{x}_i are called [support vectors](#).

Linear SVM: non-separable case

Have you ever seen a dataset that looks that this?



Restricting the problem to linearly separable training data sets is a **somehow strong assumption**.

Inequality constraints in the quadratic optimization problem **can be relaxed**.

Introduction of nonnegative variables $(\xi_i)_{1 \leq i \leq n}$ which **quantify the nonfeasability of the constraint** $y_i(\langle \mathbf{w}; \mathbf{x}_i \rangle + b) \geq 1$.

$$y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i .$$

Replace the constraint

$$y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 \quad \text{for all } i = 1, \dots, n,$$

which is too strong by the **relaxed** one

$$y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - s_i \quad \text{for all } i = 1, \dots, n,$$

for **slack variables** $s_1, \dots, s_n \geq 0$.

Linear SVM: non-separable case

The original problem is then replaced by

$$(w_*, b_*, \xi_*) \in \underset{(w, b, \xi) \in \mathbb{R}^d \times \mathbb{R} \times \mathbb{R}_+^n}{\operatorname{argmin}} \left\{ \lambda \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \xi_i \right\},$$
$$\forall i \in \{1, \dots, n\}, Y_i(\langle w; X_i \rangle + b) \geq 1 - \xi_i$$

where $\lambda > 0$.

The **soft Support Vector Machines** algorithm simultaneously the margin of the linear classifier and the average value of these slack variables

Note that, if (w_*, b_*) is solution to

$$(w_*, b_*) \in \underset{(w, b) \in \mathbb{R}^d \times \mathbb{R}}{\operatorname{argmin}} \left\{ \lambda \|w\|^2 + \frac{1}{n} \sum_{i=1}^n (1 - Y_i(\langle w; X_i \rangle + b))_+ \right\},$$

then $(w_*/\|w_*\|, b_*, \xi_*/\|w_*\|)$ is solution to the soft SVM problem.

The **soft SVM problem boils down** to computing:

$$(w_*, b_*, \xi_*) \in \underset{\substack{(w, b, \xi) \in \mathbb{R}^d \times \mathbb{R} \times \mathbb{R}_+^n \\ \forall i \in \{1, \dots, n\}, Y_i(\langle w; X_i \rangle + b) \geq 1 - \xi_i}}{\operatorname{argmin}} \left\{ \lambda \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \xi_i \right\},$$

where $\lambda > 0$.

This problem admits a **unique** solution.

It is a **quadratic programming** problem, which is easy to solve numerically.

Dedicated optimization algorithms can solve this on a large scale very efficiently.

The hinge loss

This problem can be reformulated as follows

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \max \left(0, 1 - y_i (\langle \mathbf{x}_i, \mathbf{w} \rangle + b) \right),$$

or even better, by introducing the [hinge loss](#)

$$\ell(y, y') = \max(0, 1 - yy') = (1 - yy')_+,$$

the problem can be written as

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \ell(y_i, \langle \mathbf{x}_i, \mathbf{w} \rangle + b).$$

We'll see later more loss functions ℓ and explain the general setting hidden behind this.

The joint law of (X, Y) is **not assumed to belong to any parametric or semiparametric** family of models.

The classification risk **cannot be computed nor minimized**, it is instead estimated by the empirical classification risk defined as

$$\hat{L}_{\text{miss}}^n(f) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i \neq f(\mathbf{x}_i)},$$

where $(X_i, Y_i)_{1 \leq i \leq n}$ are independent with the same distribution as (X, Y) .

The classification problem then boils down to solving

$$\hat{f}^n \in \operatorname{argmin}_{f \in \mathcal{F}} \hat{L}_{\text{miss}}^n(f),$$

for a chosen class \mathcal{F} of classifiers.

Introduction to nonparametric classification

Nonparametric classification based on the empirical risk minimization may seem appealing

It cannot be used to derive efficient practical classifiers due to **the computational cost of the optimization problem**.

The target loss function \hat{L}_{miss}^n is replaced by **a convex surrogate and its minimization is constrained to a convex set of classifiers**.

For any convex function $\phi : \mathcal{X} \rightarrow \mathbb{R}$, it is possible to build a classifier f given by $f_\phi = \text{sign}(\phi)$. The associated empirical classification is then

$$\hat{L}_{\text{miss}}^n(\phi) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i \neq f_\phi(X_i)} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i \phi(X_i) < 0}.$$

Replacing the indicator function by **any convex loss function ℓ yields a convex surrogate**:

$$\hat{L}_{\text{miss}}^{n,\text{conv}}(\phi) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i \phi(X_i)).$$

Penalizing the smoothness of the function ϕ is penalized, $\hat{L}_{\text{miss}}^{n,\text{conv}}$ may be replaced by

$$\hat{L}_{\text{miss}}^{n,\text{conv}}(\phi) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i \phi(X_i)) + \lambda \|\phi\|^2,$$

where $\lambda > 0$ and $\|\cdot\|$ is a norm on the space \mathcal{H} .

The soft Support Vector Machines algorithm fits this framework with the affine base function $\phi : \mathbf{x} \mapsto \langle \mathbf{w} ; \mathbf{x} \rangle + b$ and ℓ chosen as the hinge loss $\ell : x \mapsto (1 - x)_+$ when the target function is penalized by its margin $\|\mathbf{w}\|^2$.

A useful case in practice consists in choosing \mathcal{H} as a **Reproducing Kernel Hilbert Space** with positive definite kernel k on $\mathcal{X} \times \mathcal{X}$.

A function k on $\mathcal{X} \times \mathcal{X}$ is said to be a **positive definite kernel** if and only if it is symmetric and if for all $n \geq 1$, $(x_1, \dots, x_n) \in \mathcal{X}^n$ and all $(a_1, \dots, a_n) \in \mathbb{R}^n$,

$$\sum_{1 \leq i, j \leq n} a_i a_j k(x_i, x_j) \geq 0.$$

The **Reproducing Kernel Hilbert Space** with kernel k is the only Hilbert space $\mathcal{H} \subset \mathbb{R}^{\mathcal{X}}$ such that for all $x \in \mathcal{X}$, $k(x, \cdot) \in \mathcal{H}$ and for all $x \in \mathcal{X}$ and all $f \in \mathcal{H}$,

$$f(x) = \langle f; k(x, \cdot) \rangle_{\mathcal{H}}.$$

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ a positive definite kernel and \mathcal{H} the RKHS with kernel k .

$$\hat{\phi}_{\mathcal{H}}^n \in \operatorname{argmin}_{\phi \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i \phi(X_i)) + \lambda \|\phi\|_{\mathcal{H}}^2,$$

where $\|\phi\|_{\mathcal{H}}^2 = \langle \phi; \phi \rangle$, is given by

$$\hat{\phi}_{\mathcal{H}}^n : x \mapsto \sum_{i=1}^n \hat{\alpha}_i k(X_i, x),$$

with

$$\hat{\alpha} \in \operatorname{argmin}_{\alpha \in \mathbb{R}^n} \left\{ \frac{1}{n} \sum_{i=1}^n \ell \left(\sum_{j=1}^n \alpha_j Y_j k(X_j, X_i) \right) + \lambda \sum_{1 \leq i, j \leq n} \alpha_i \alpha_j k(X_i, X_j) \right\}.$$