

Sapienza University of Rome

Master in Artificial Intelligence and Robotics

Machine Learning

A.Y. 2024/2025

Prof. Luca Iocchi

7. Linear models for classification

Luca Iocchi

Overview

- Linearly separable data
- Linear models
- Least squares
- Perceptron
- Fisher's linear discriminant
- Support Vector Machines

References

C. Bishop. Pattern Recognition and Machine Learning. Sect. 4.1, 7.1

T. Mitchell. Machine Learning. Section 4.4

Linear Models for Classification

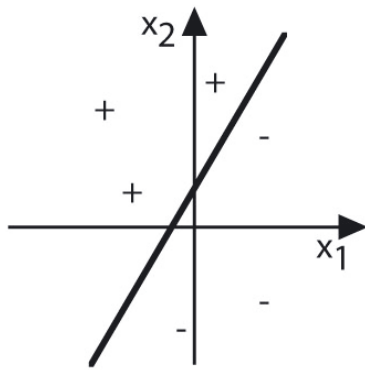
Learning a function $f : X \rightarrow Y$, with ...

- $X \subseteq \mathbb{R}^d$
- $Y = \{C_1, \dots, C_k\}$

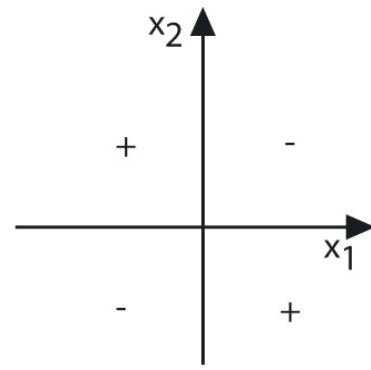
assuming *linearly separable* data.

Linearly separable data

Instances in a data set are *linearly separable* iff there exists a hyperplane that separates the instance space into two regions, such that differently classified instances are separated



(a)



(b)

Linear discriminant functions

Linear discriminant function

$$y : X \rightarrow \{C_1, \dots, C_K\}$$

Two classes:

$$y(x) = w^T x + w_0$$

K -classes:

$$y_1(x) = w_1^T x + w_{10}$$

...

$$y_K(x) = w_K^T x + w_{K0}$$

Compact notation

Two classes:

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}, \text{ with:}$$

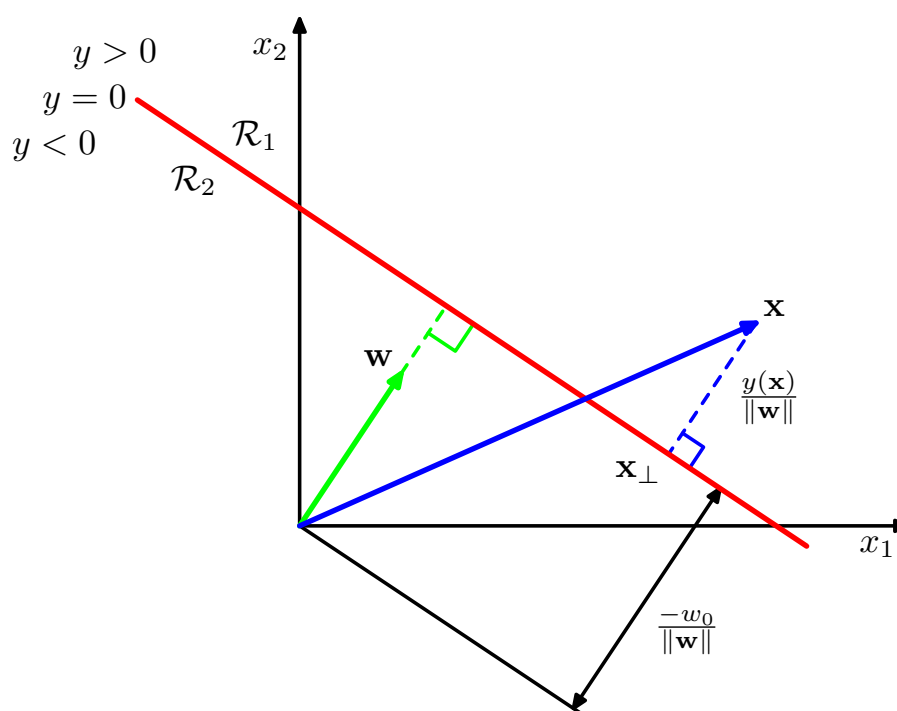
$$\tilde{\mathbf{w}} = \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix}, \tilde{\mathbf{x}} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

K classes:

$$y(\mathbf{x}) = \begin{pmatrix} y_1(\mathbf{x}) \\ \vdots \\ y_K(\mathbf{x}) \end{pmatrix} = \begin{pmatrix} \mathbf{w}_1^T \mathbf{x} + w_{10} \\ \vdots \\ \mathbf{w}_K^T \mathbf{x} + w_{K0} \end{pmatrix} = \begin{pmatrix} \tilde{\mathbf{w}}_1^T \\ \vdots \\ \tilde{\mathbf{w}}_K^T \end{pmatrix} \tilde{\mathbf{x}} = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}}, \text{ with:}$$

$$\tilde{\mathbf{W}}^T = \begin{pmatrix} \tilde{\mathbf{w}}_1^T \\ \vdots \\ \tilde{\mathbf{w}}_K^T \end{pmatrix}, \text{ i.e.: } \tilde{\mathbf{W}} = (\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_K)$$

Linear discriminant functions



Multiple classes

K -class discriminant comprising K linear functions

$$y(x) = \begin{pmatrix} y_1(x) \\ \vdots \\ y_K(x) \end{pmatrix} = \begin{pmatrix} \tilde{w}_1^T \tilde{x} \\ \vdots \\ \tilde{w}_K^T \tilde{x} \end{pmatrix} = \tilde{W}^T \tilde{x}$$

Classify x as C_k with $k = \operatorname{argmax}_{j=1,\dots,K} \{y_j(x)\}$

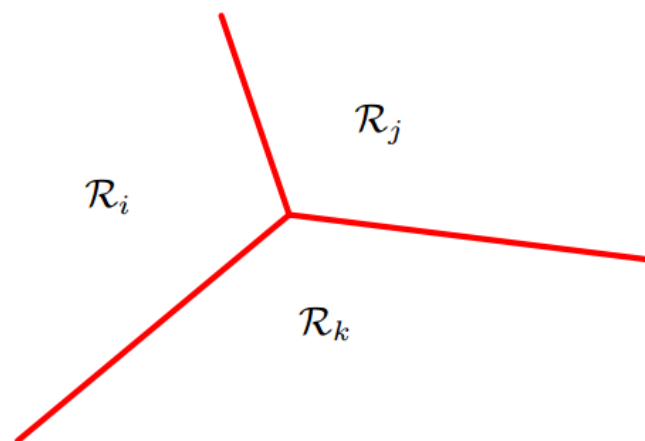
Decision boundary between C_k and C_j (hyperplane in \mathbb{R}^{D-1}):

$$(\tilde{w}_k - \tilde{w}_j)^T \tilde{x} = 0$$

Multiple classes

when we have n -dimensional space

Example of K -class discriminant



Learning linear discriminants

Given a multi-class classification problem and data set D with linearly separable data,

determine \tilde{W} such that $y(x) = \tilde{W}^T \tilde{x}$ is the K -class discriminant.

Approaches to learn linear discriminants

- Least squares
- Perceptron
- Fisher's linear discriminant
- Support Vector Machines

Least squares

Given $D = \{(x_n, t_n)_{n=1}^N\}$, find the linear discriminant

$$y(x) = \tilde{W}^T \tilde{x}$$

1-of-K coding scheme for t : $x \in C_k \rightarrow t_k = 1, t_j = 0$ for all $j \neq k$.
 E.g., $t_n = (0, \dots, 1, \dots, 0)^T$

$$\tilde{X} = \begin{pmatrix} \tilde{x}_1^T \\ \dots \\ \tilde{x}_N^T \end{pmatrix} \quad T = \begin{pmatrix} t_1^T \\ \dots \\ t_N^T \end{pmatrix}$$

↳ complete representation of dataset

Least squares

Minimize sum-of-squares error function

$$E(\tilde{W}) = \frac{1}{2} \text{Tr} \left\{ (\tilde{X}\tilde{W} - T)^T (\tilde{X}\tilde{W} - T) \right\}$$

Closed-form solution:

$$\tilde{W} = \underbrace{(\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T}_\tilde{X}^\dagger T$$

↳ pseudoinverse

↙ solution of minimization can be written in closed form

Learned model:

$$y(X) = \tilde{W}^T \tilde{X} = T^T (\tilde{X}^\dagger)^T \tilde{X}$$

Least squares

Classification of new instance $\mathbf{x}' \notin D$

Use learnt \tilde{W} to compute:

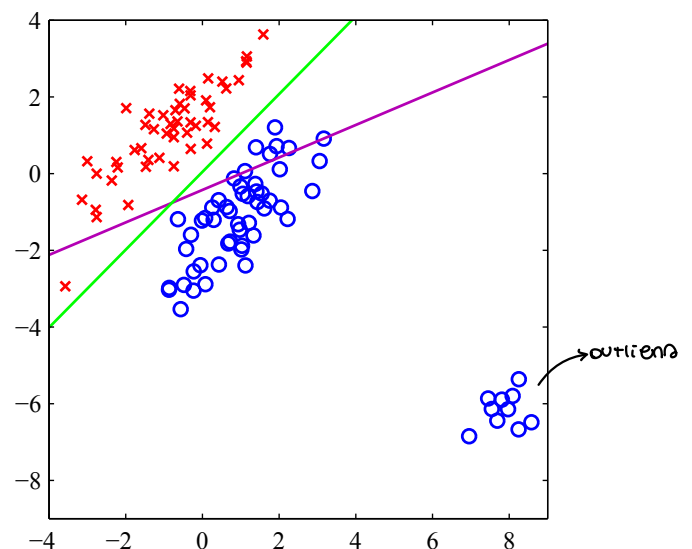
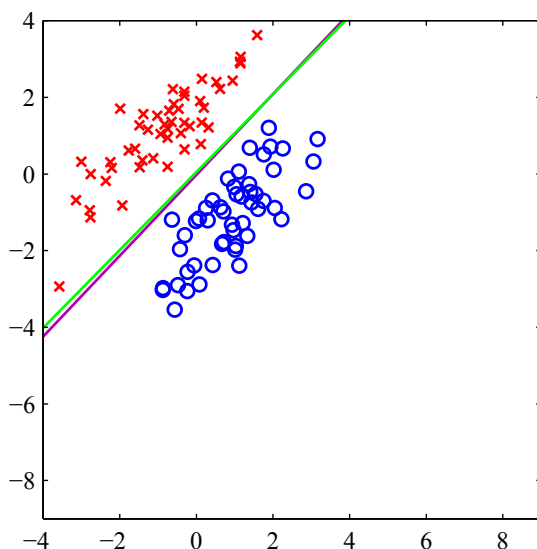
$$y(\mathbf{x}') = \tilde{W}^T \tilde{\mathbf{x}}' = \begin{pmatrix} y_1(\mathbf{x}') \\ \vdots \\ y_K(\mathbf{x}') \end{pmatrix}$$

Predict C_k with:

$$k = \operatorname{argmax}_{j \in \{1, \dots, K\}} \{y_j(\mathbf{x}')\}$$

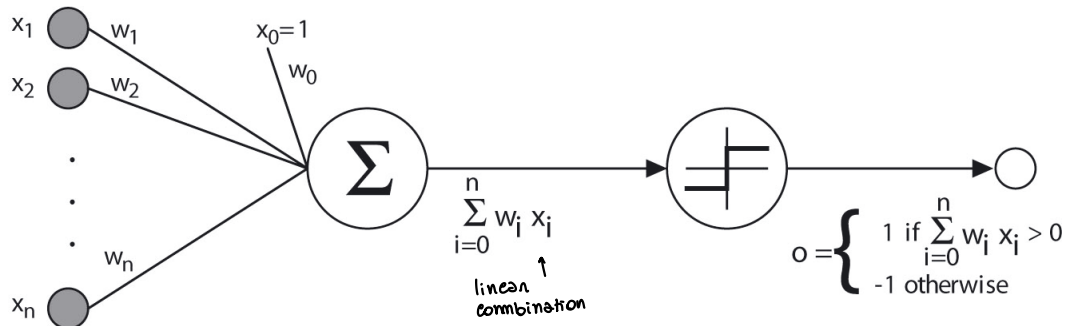
Issues with least squares

Assume Gaussian conditional distributions. Not robust to outliers!



Perceptron

Learning $f : \mathbb{R}^d \rightarrow \{-1, +1\}$



$$o(x_1, \dots, x_d) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_d x_d > 0 \\ -1 & \text{otherwise.} \end{cases}$$

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

Perceptron training rule

Consider the *unthresholded linear unit*, where

$$o = w_0 + w_1 x_1 + \dots + w_d x_d = w^T x$$

Let's learn w_i from training examples $D = \{(x_n, t_n)_{n=1}^N\}$ that minimize the squared error (*loss function*)

$$E(w) \equiv \underbrace{\frac{1}{2} \sum_{n=1}^N (t_n - o_n)^2}_{\substack{\text{o when model predict} \\ \text{exactly the label of} \\ \text{sample n of dataset}}} = \frac{1}{2} \sum_{n=1}^N (t_n - w^T x_n)^2$$

Perceptron training rule

$$\begin{aligned}
 \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{n=1}^N (t_n - w^T x_n)^2 = \frac{1}{2} \sum_{n=1}^N \frac{\partial}{\partial w_i} (t_n - w^T x_n)^2 \\
 &= \frac{1}{2} \sum_{n=1}^N 2(t_n - w^T x_n) \frac{\partial}{\partial w_i} (t_n - w^T x_n) \\
 &= \sum_{n=1}^N (t_n - w^T x_n) \frac{\partial}{\partial w_i} (t_n - w^T x_n) \\
 &= \sum_{n=1}^N (t_n - w^T x_n) (-x_{i,n})
 \end{aligned}$$

Perceptron training rule

Unthresholded unit:

Update of weights w

$$\begin{aligned}
 w_i &\leftarrow w_i + \Delta w_i \\
 \Delta w_i &= -\eta \frac{\partial E}{\partial w_i} = \eta \sum_{n=1}^N (t_n - w^T x_n) x_{i,n}
 \end{aligned}$$

η is a small constant (e.g., 0.05) called *learning rate*

Perceptron training rule

Thresholded unit:

Update of weights w

$$w_i \leftarrow w_i + \Delta w_i$$

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} = \eta \sum_{n=1}^N (t_n - \text{sign}(w^T x_n)) x_{i,n}$$

Perceptron algorithm

Given perceptron model $o(x) = \text{sign}(w^T x)$ and data set D , determine weights w .

- 1 Initialize \hat{w} (e.g. small random values)
- 2 Repeat until termination condition
 - $\hat{w}_i \leftarrow \hat{w}_i + \Delta w_i$
- 3 Output \hat{w}

Perceptron algorithm

Batch mode: Consider all dataset D

$$\Delta w_i = \eta \sum_{(x,t) \in D} (t - o(x)) x_i \quad \Rightarrow \text{high computational}$$

Mini-Batch mode: Choose a small subset $S \subset D$ \Rightarrow better choice

$$\Delta w_i = \eta \sum_{(x,t) \in S} (t - o(x)) x_i$$

Incremental mode: Choose one sample $(x, t) \in D$ \Rightarrow high variance

$$\Delta w_i = \eta (t - o(x)) x_i$$

$o(x) = w^T x$ for unthresholded, $o(x) = \text{sign}(w^T x)$ for thresholded

Incremental and mini-batch modes speed up convergence and are less sensitive to local minima.

Perceptron algorithm

Termination conditions

- Predefined number of iterations
- Threshold on changes in the loss function $E(w)$

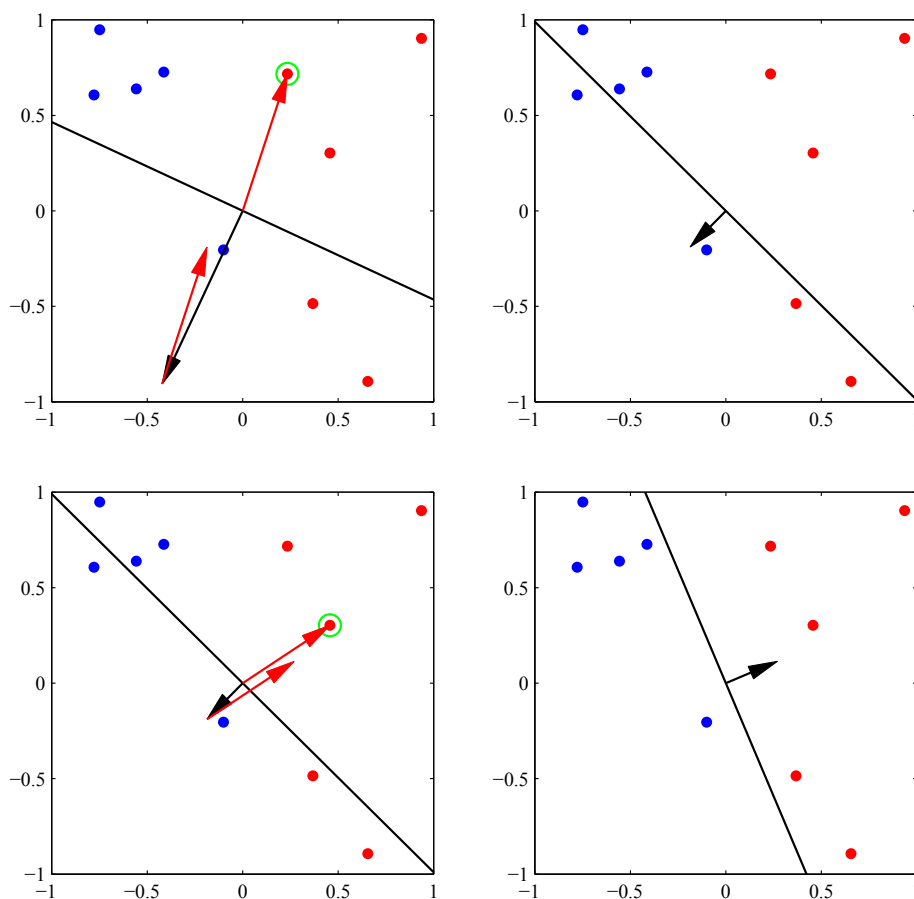
Perceptron training rule

Example:

$$\eta = 0.1, x_i = 0.8$$

- if $t = 1$ and $o = -1$ then $\Delta w_i = 0.16$
- if $t = -1$ and $o = 1$ then $\Delta w_i = -0.16$

Perceptron training rule



Perceptron training rule

Can prove it will converge:

- if training data is linearly separable
- and η sufficiently small

Small $\eta \rightarrow$ slow convergence.

Perceptron: Prediction

Classification of new instance $x' \notin D$:

Predict C_k , with $k = \text{sign}(w^T x')$, using learnt w

Fisher's linear discriminant

Consider two classes case.

Determine $y = w^T x$
and classify $x \in C_1$ if $y \geq -w_0$, $x \in C_2$ otherwise.

Corresponding to the projection on a line determined by w .

Fisher's linear discriminant

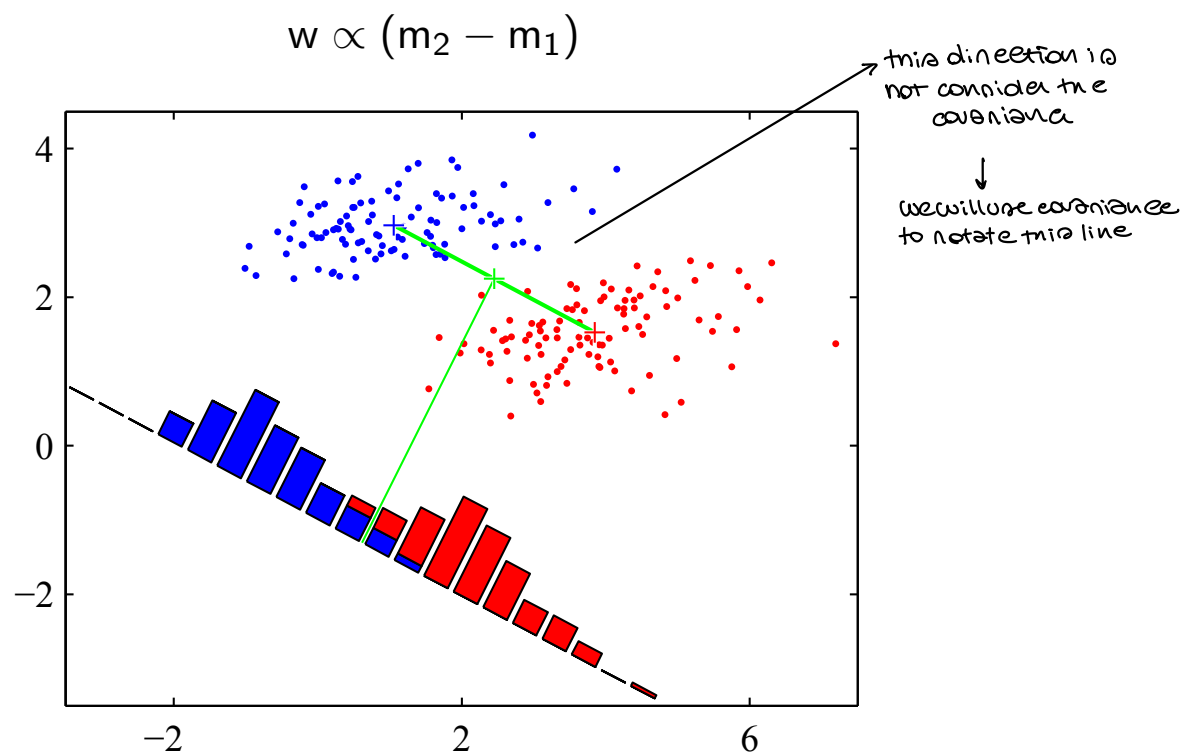
Adjusting w to find a direction that maximizes class separation.

Consider a data set with N_1 points in C_1 and N_2 points in C_2

$$m_1 = \frac{1}{N_1} \sum_{n \in C_1} x_n \quad m_2 = \frac{1}{N_2} \sum_{n \in C_2} x_n$$

Choose w that maximizes $J(w) = w^T (m_2 - m_1)$, subject to $\|w\| = 1$.

Fisher's linear discriminant



Fisher's linear discriminant

Fisher criterion

$$J(w) = \frac{w^T S_B w}{w^T S_W w}$$

with

$$S_B = (m_2 - m_1)(m_2 - m_1)^T$$

Between class scatter

$$S_W = \sum_{n \in C_1} (x_n - m_1)(x_n - m_1)^T + \sum_{n \in C_2} (x_n - m_2)(x_n - m_2)^T$$

Within class scatter

Choose w that maximizes $J(w)$.

Fisher's linear discriminant

Find w that maximizes

$$J(w) = \frac{w^T S_B w}{w^T S_W w}$$

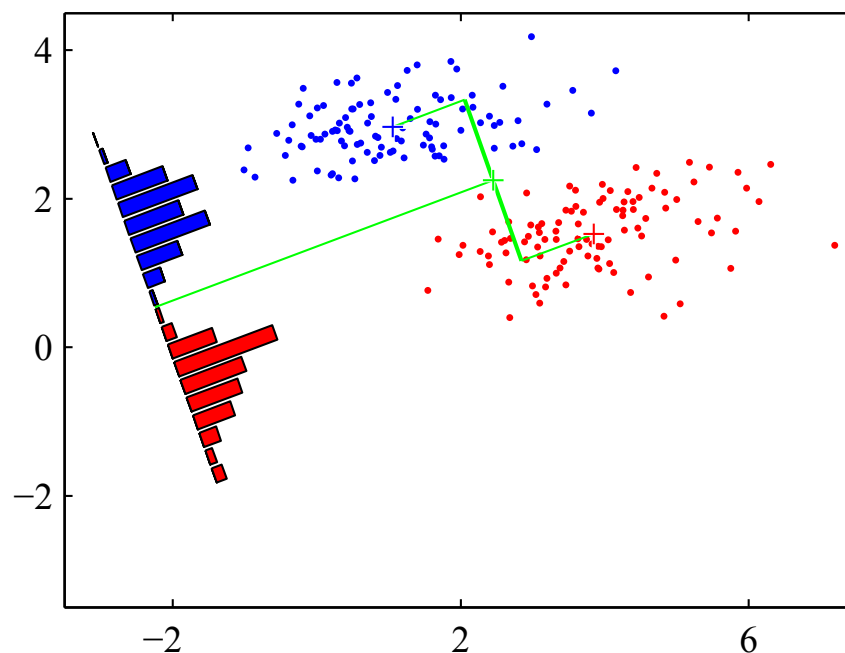
by solving

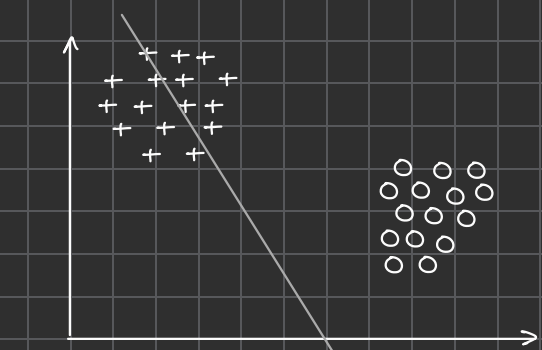
$$\frac{d}{dw} J(w) = 0$$

$$\Rightarrow w^* \propto S_W^{-1}(m_2 - m_1)$$

Fisher's linear discriminant

$$w \propto S_W^{-1}(m_2 - m_1)$$

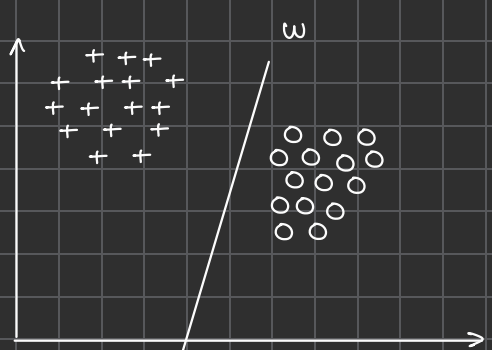




prof says "true"

if $\theta < 0$

change slow



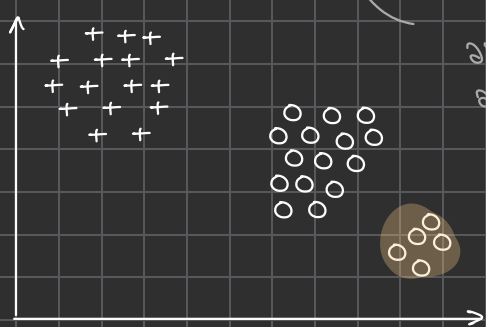
increase all terms = 0

$$\Delta w = (t_n - \text{sign}(w^T x_n))$$

if $\theta > 0$

change fast

SO PERCEPTIONS
IS ROBUST TO
OUTLIERS

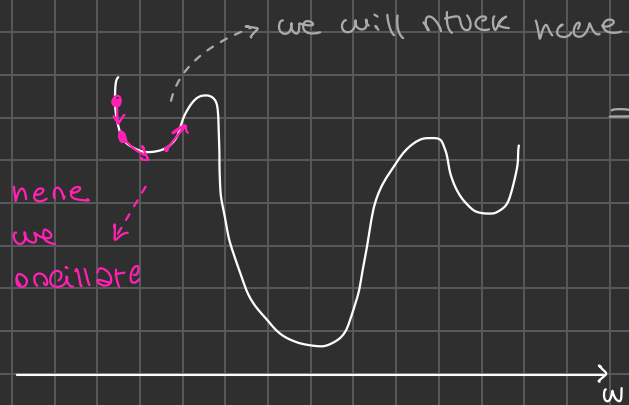


also with outlier
all this terms = 0

outlier

contribute 0 to update rule

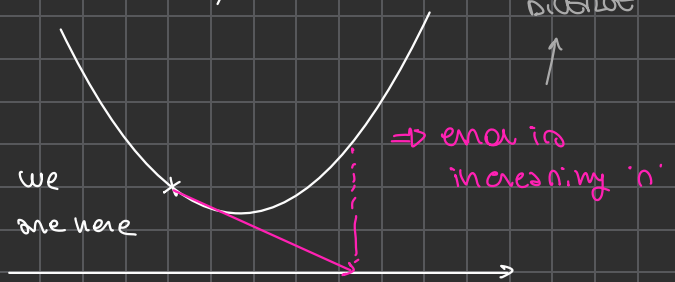
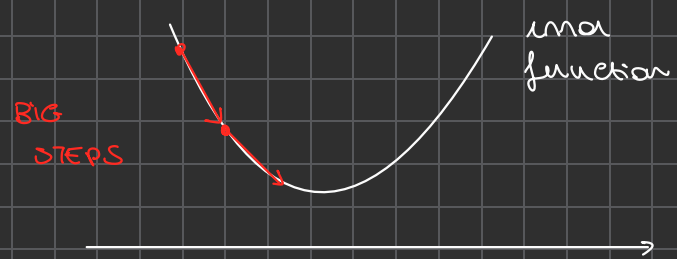
ERROR FUNCTION:

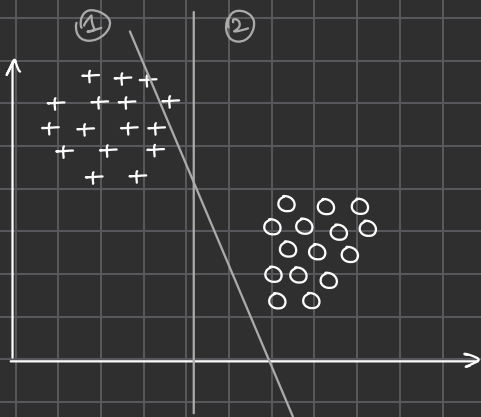


⇒ that mean
we stop in LOCAL
MINIMAL

we cannot define how
small or big must be
ERROR FUNCTION

• WHAT HAPPEN IF "learning rate" is too high

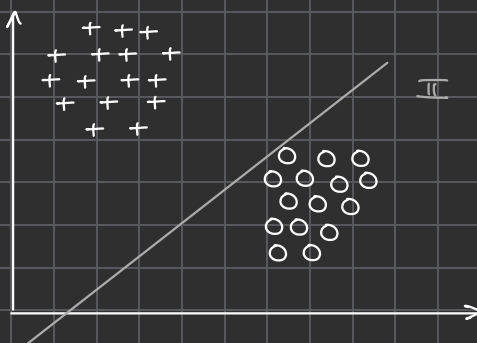
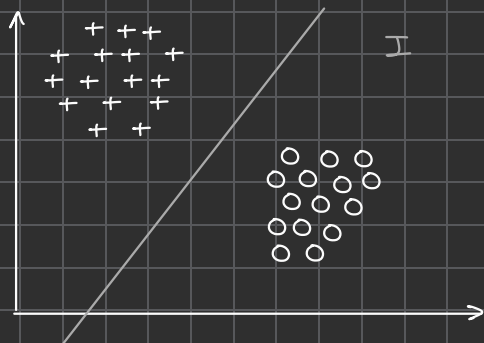




\Rightarrow this is Not Good because if we consider samples that are not in dataset this will not correct

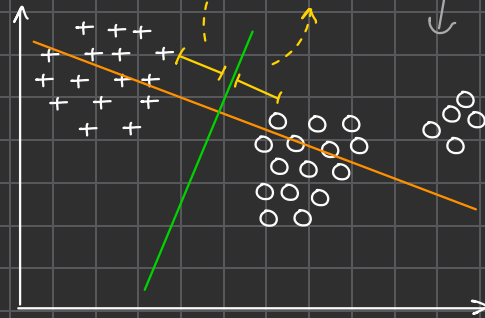
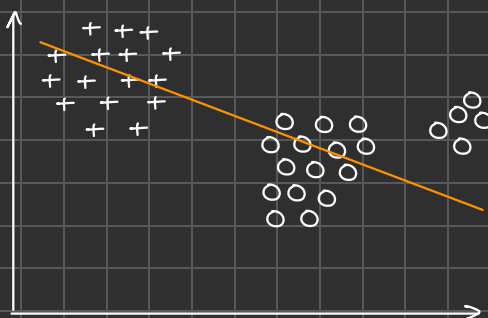
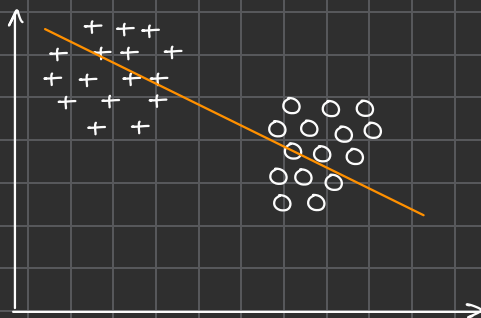
come
work is
important

\rightarrow our algorithm will stop in second iterations



if we consider a small learning rate II is more probable than I

we need some margins so will be more correct with new samples



FISHER we use data for estimation of error function no outliers we will change small slope

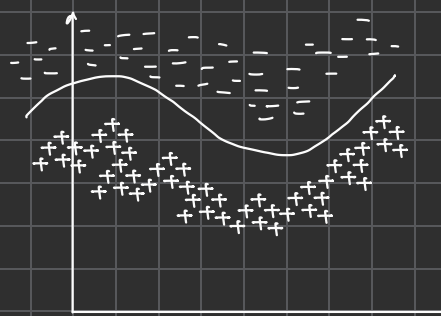
$$y(x; w) = w^T x$$

\Rightarrow linear model also x

$$y(x; w) = w^T \phi(x)$$

\Rightarrow linear model but not in x

maybe non-linear
still linear in w



EXAMPLE

NOT SEPARABLE

$$y(x^T, w) = w^T \phi(x) = w^T \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix} = w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

Fisher's linear discriminant

Summarizing, given a two classes classification problem, Fisher's linear discriminant is given by the function $y = w^T x$ and the classification of new instances is given by $y \geq -w_0$ where

$$w = S_W^{-1}(m_2 - m_1)$$

$$w_0 = w^T m$$

m is the global mean of all the data set.

Fisher's linear discriminant

Multiple classes.

$$y = W^T x$$

Maximizing

$$J(W) = \text{Tr} \left\{ (WS_W W^T)^{-1} (WS_B W^T) \right\}$$

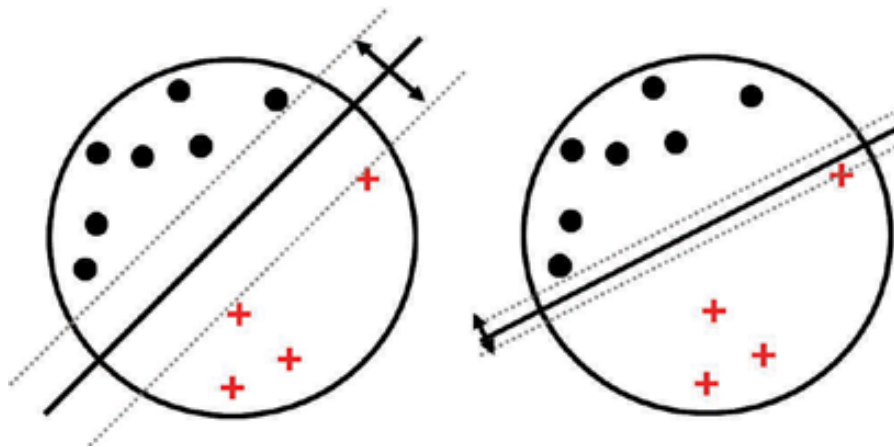
...

Support Vector Machines

→ ROBUST TO OUTLIERS THE VALUE WILL BE \Rightarrow

Support Vector Machines (SVM) for Classification aims at maximum margin providing for better accuracy.

→ similar to exam
Q: WHERE THIS TWO IMAGES COMES FROM
A: LEFT → SVM (probably)



Support Vector Machines

Let's consider binary classification $f : X \rightarrow \{+1, -1\}$ with data set $D = \{(x_n, t_n)_{n=1}^N\}$, $t_n \in \{+1, -1\}$ and a linear model

$$y(x) = w^T x + w_0$$

Assume D is linearly separable

$$\exists w, w_0 \text{ s.t. } \begin{cases} y(x_n) > 0, & \text{if } t_n = +1 \\ y(x_n) < 0, & \text{if } t_n = -1 \end{cases}$$

COMPACT FORMULA $\rightarrow t_n y(x_n) > 0 \quad \forall n = 1, \dots, N$

Support Vector Machines

Let x_k be the closest point of the data set D to the hyperplane $h : w^T x + w_0 = 0$

the *margin* (smallest distance between x_k and h) is $\frac{|y(x_k)|}{||w||}$

Given data set D and hyperplane h , the margin is computed as

$$\min_{n=1,\dots,N} \frac{|y(x_n)|}{||w||} = \dots = \frac{1}{||w||} \min_{n=1,\dots,N} [t_n(w^T x_n + w_0)]$$

using the property $|y(x_n)| = t_n y(x_n)$

Support Vector Machines

Given data set D , the hyperplane $h^* : w^{*T} x + w_0^* = 0$ with maximum margin is computed as

$$w^*, w_0^* = \operatorname{argmax}_{w, w_0} \underbrace{\frac{1}{||w||} \min_{n=1,\dots,N} [t_n(w^T x_n + w_0)]}_{\text{MAXIMUM MARGIN OPTIMIZATION PROBLEM}}$$

Support Vector Machines

Rescaling all the points does not affect the solution.

Rescale in such a way that for the closet point x_k we have

$$t_k(w^T x_k + w_0) = 1$$

Canonical representation:

$$t_n(w^T x_n + w_0) \geq 1 \quad \forall n = 1, \dots, N$$

↑
is only = 1 for the
closest point to the margin

if DATASET IS NOT SEPARABLE

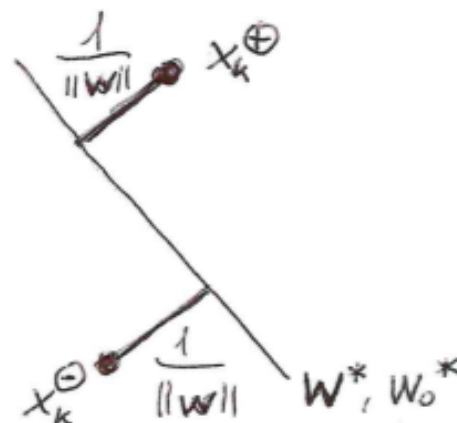
THE ALGORITHM WILL
RETURN NO SOLUTION

Support Vector Machines

When the maximum margin hyperplane w^*, w_0^* is found, there will be at least 2 closest points x_k^+ and x_k^- (one for each class).

$$w^{*T} x_k^+ + w_0^* = +1$$

$$w^{*T} x_k^- + w_0^* = -1$$



Support Vector Machines

In the canonical representation of the problem the maximum margin hyperplane can be found by solving the optimization problem

$$w^*, w_0^* = \operatorname{argmax} \frac{1}{||w||} = \operatorname{argmin} \frac{1}{2} ||w||^2$$

subject to

$$t_n(w^T x_n + w_0) \geq 1 \quad \forall n = 1, \dots, N$$

Quadratic programming problem solved with Lagrangian method.

Support Vector Machines

Solution

$$w^* = \sum_{n=1}^N a_n^* t_n x_n$$

a_i^* (Lagrange multipliers): results of the Lagrangian optimization problem

$$\tilde{L}(a) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m x_n^T x_m$$

subject to

$$a_n \geq 0 \quad \forall n = 1, \dots, N$$

$$\sum_{n=1}^N a_n t_n = 0$$

Support Vector Machines

Note:

samples x_n for which $a_n^* = 0$ do not contribute to the solution

Karush-Kuhn-Tucker (KKT) condition:

for each $x_n \in D$, either $a_n^* = 0$ or $t_n y(x_n) = 1$

thus $t_n y(x_n) > 1$ implies $a_n^* = 0$

Support vectors: x_k such that $t_k y(x_k) = 1$ and $a_k^* > 0$

$$SV \equiv \{x_k \in D \mid t_k y(x_k) = 1\}$$

Support Vector Machines

Hyperplanes expressed with support vectors

$$y(x) = \sum_{x_j \in SV} a_j^* t_j x^T x_j + w_0^* = 0$$

Note: other vectors $x_n \notin SV$ do not contribute ($a_n^* = 0$)

Support Vector Machines

To compute w_0^* :

Support vector $x_k \in SV$ satisfies $t_k y(x_k) = 1$

$$t_k \left(\sum_{x_j \in SV} a_j^* t_j x_k^T x_j + w_0^* \right) = 1$$

Multiplying by t_k and using $t_k^2 = 1$

$$w_0^* = t_k - \sum_{x_j \in SV} a_j^* t_j x_k^T x_j$$

Support Vector Machines

Instead of using one particular support vector x_k to determine w_0

$$w_0^* = t_k - \sum_{x_j \in SV} a_j^* t_j x_k^T x_j$$

a more stable solution is obtained by averaging over all the support vectors

$$w_0^* = \frac{1}{|SV|} \sum_{x_k \in SV} \left(t_k - \sum_{x_j \in S} a_j^* t_j x_k^T x_j \right)$$

Support Vector Machines

Given the maximum margin hyperplane determined by a_k^* , w_0^*

Classification of a new instance x'

$$y(x') = \text{sign} \left(\sum_{x_k \in SV} a_k^* t_k x'^T x_k + w_0^* \right)$$

Support Vector Machines

Optimization problem for determining w, w_0 (dimension $d + 1$, with $X = \mathbb{R}^d$) transformed in an optimization problem for determining a (dimension $|D|$)

Efficient when $d \ll |D|$ (most of a_i will be zero).

Very useful when d is large or infinite.

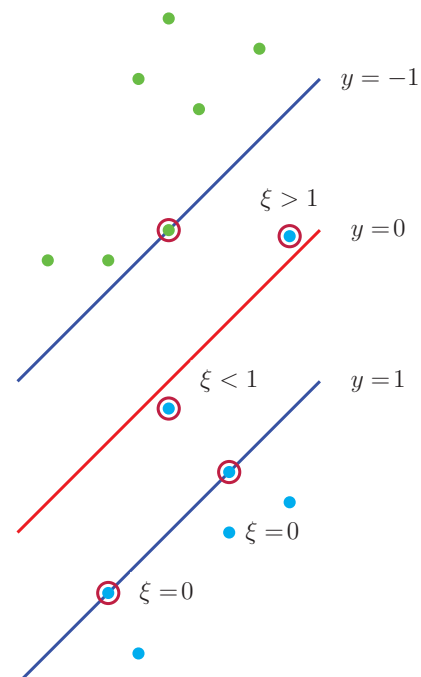
Support Vector Machines with soft margin constraints

What if data are “almost” linearly separable (e.g., a few points are on the “wrong side”)

Let us introduce *slack variables* $\xi_n \geq 0$ $n = 1, \dots, N$

Support Vector Machines with soft margin constraints

- $\xi_n = 0$ if point on or inside the correct margin boundary
- $0 < \xi_n \leq 1$ if point inside the margin but correct side
- $\xi_n > 1$ if point on wrong side of boundary



when $\xi_n = 1$, the sample lies on the decision boundary $y(x_n) = 0$
 when $\xi_n > 1$, the sample will be mis-classified

Support Vector Machines with soft margin constraints

Soft margin constraint

$$t_n y(x_n) \geq 1 - \xi_n, \quad n = 1, \dots, N$$

Optimization problem with soft margin constraints

$$w^*, w_0^* = \operatorname{argmin} \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n$$

subject to

$$t_n y(x_n) \geq 1 - \xi_n, \quad n = 1, \dots, N$$

$$\xi_n \geq 0, \quad n = 1, \dots, N$$

C is a constant (inverse of a regularization coefficient)

Support Vector Machines with soft margin constraints

Solution similar to the case of linearly separable data.

$$w^* = \sum_{n=1}^N a_n^* t_n x_n$$

$$w_0^* = \dots$$

with a_n^* computed as solution of a Lagrangian optimization problem.

Basis functions

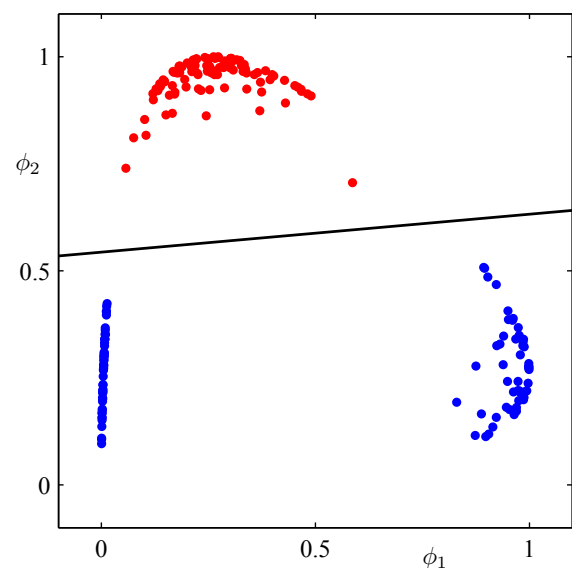
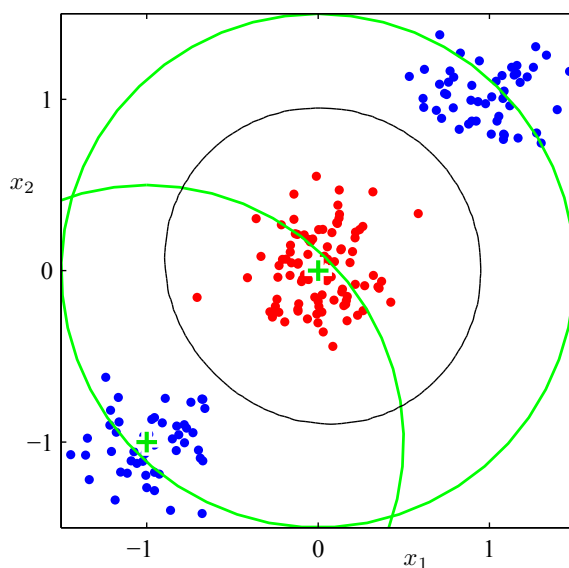
So far we considered models working directly on x .

All the results hold if we consider a non-linear transformation of the inputs $\phi(x)$ (*basis functions*).

Decision boundaries will be linear in the feature space ϕ and non-linear in the original space x

Classes that are linearly separable in the feature space ϕ may not be separable in the input space x .

Basis functions example



Basis functions examples

- Linear
- Polynomial
- Radial Basis Function (RBF)
- Sigmoid
- ...

Linear models for non-linear functions

Learning non-linear function

$$f : X \rightarrow \{C_1, \dots, C_K\}$$

from data set D non-linearly separable.

Find a non-linear transformation ϕ and learn a linear model

$$y(x) = w^T \phi(x) + w_0 \text{ (two classes)}$$

$$y_k(x) = w_k^T \phi(x) + w_{k0} \text{ (multiple classes)}$$

Summary

- Basic methods for learning linear classification functions
- Based on solution of an optimization problem
- Closed form vs. iterative solutions
- Sensitivity to outliers
- Learning non-linear functions with linear models using basis functions
- Further developed as kernel methods