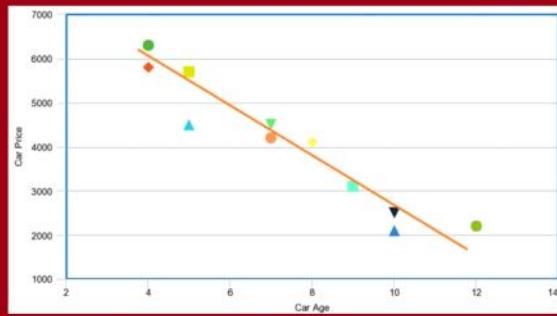




DIPARTIMENTO
DI INGEGNERIA
DELL'INFORMAZIONE



UNIVERSITÀ
DEGLI STUDI
DI PADOVA



Linear Predictors

Machine Learning 2021
UML Book Chapter 9
Slides P. Zanuttigh

Some material from F. Vandin, N. Ailon, S. Shwartz, J. Janecek



Linear Predictors

The design of a ML strategy requires 2 main steps:

- Select an hypothesis class
- Select an algorithm to find the predictor

→ algo that perform ERM and find the optimal predictor \hat{a}

Hypotheses Classes

- Halfspaces (classification)
- Linear Regression (regression)
- Logistic Regression (classification modeled as a regression problem)

Algorithms

- Linear Programming (for halfspaces) } } less used
- Perceptron (for halfspaces)
- Least Squares (for regression)

Fundamental
to consider
the previous
knowledge of
the system
to choose
the hyp class



Affine Function Model

Class of Affine Functions:

$$L_d = \{h_{w,b}, w \in \mathbb{R}^d, b \in \mathbb{R}\}$$

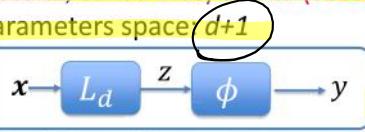
where

$$h_{w,b} = \langle w, x \rangle + b = \left(\sum_{i=1}^d w_i x_i \right) + b$$

Each member of L_d is a function $x \rightarrow \langle w, x \rangle + b$, $w \in \mathbb{R}^d, b \in \mathbb{R}$

Two parameters: b (scalar, called bias) and w (vector)

Dimensionality of parameters space: $d+1$



Hypothesis class: $\phi \circ L_d : \mathbb{R} \rightarrow \mathcal{Y}$

- Binary classification $\mathcal{Y} = \{-1, 1\} \rightarrow \phi(z) = \text{sign}(z)$
- Regression $\mathcal{Y} = \mathbb{R} \rightarrow \phi(z) = z$

→ 2 params

- $w \in \mathbb{R}^d$ vector of-dimensional
- $b \in \mathbb{R}$ scalar

↓
TO LEARN

\vec{x} INPUT VALUES
vector

y OUTPUT
scalar

$$= \vec{w} \cdot \vec{x} + b$$

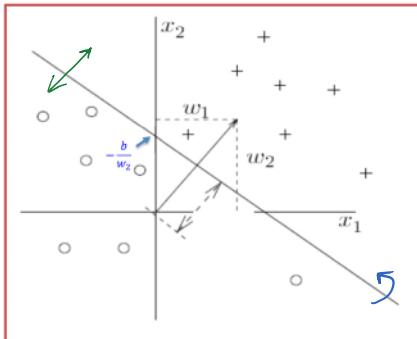
↓

OUTPUT → sign(y)

$$\mathcal{H} \rightarrow \underbrace{\text{sign}}_{\text{func}} \circ \underbrace{L_d}_{\text{affine function}}$$



Geometric Interpretation (2D)



$$x_2 = -\frac{w_1}{w_2}x_1 - \frac{b}{w_2}$$

$d = 2$
params $\rightarrow w_1, w_2, b$

shift
from
the
origin

- The bias is proportional to the offset of the line from the origin
- The weights determine the slope of the line
- The weight vector is perpendicular to the line



Homogeneous Linear Functions

Homogeneous coordinates:

- Idea: incorporate b into w as an extra dimension/coordinate
- Add an extra dimension to w : $w \rightarrow w' = \langle b, w_1, w_2, \dots, w_d \rangle \rightarrow w' \in \mathbb{R}^{d+1}$
- Add an extra element to each vector x : $x \rightarrow x' = \underbrace{\begin{pmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix}}_{x'} \rightarrow x' \in \mathbb{R}^{d+1}$

Homogeneous linear function:

Rewrite affine functions: $L_d = \{h_{w,b}, w \in \mathbb{R}^d, b \in \mathbb{R}\}$ using homogeneous coord.

$$h_{w,b} = \langle w, x \rangle + b = \left(\sum_{i=1}^d w_i x_i \right) + b = b + w_1 x_1 + \dots + w_d x_d$$
$$\langle w', x' \rangle = \left(\sum_{i=1}^{d+1} w'_i x'_i \right) = b * 1 + w_1 x_1 + \dots + w_d x_d$$

equal

- $\langle w, x \rangle + b = \langle w', x' \rangle$, rewrite affine function as a linear model
- Get rid of bias (incorporated in the weights vector)
- The affine function becomes a linear function !

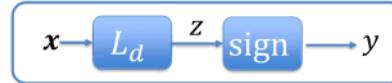
Linear Model

→ How to
rewrite the
affine model
as a linear
model
↓
incorporate the
 b (as into w)
 b



Halfspaces Hypotheses Class

Halfspace hypothesis class



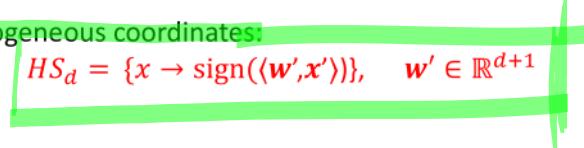
- Input: $\mathcal{X} = \mathbb{R}^d$ (for each sample a vector of features)
 - Using homogeneous coordinates: $x \rightarrow x' = (1, x_1, x_2, \dots, x_d) \in \mathbb{R}^{d+1}$
- Output: $\mathcal{Y} = \{-1, 1\}$ (binary classification)
- Loss: 0-1 loss
 - $\begin{cases} 1 & \rightarrow \text{correct output} \\ 0 & \rightarrow \text{otherwise} \end{cases}$

Halfspace Model:

$$HS_d = \text{sign} \circ L_d = \{x \rightarrow \text{sign}(\langle w, x \rangle + b)\}, \quad w \in \mathbb{R}^d, b \in \mathbb{R}$$

Using homogeneous coordinates:

$$HS_d = \{x \rightarrow \text{sign}(\langle w', x' \rangle)\}, \quad w' \in \mathbb{R}^{d+1}$$





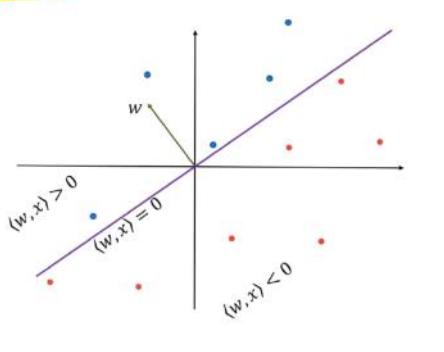
Linear Classification: Halfspaces Hypotheses Class

$\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \{-1, 1\}$, 0-1 loss

Hypothesis class = halfspaces

$$HS_d = \text{sign} \circ L_d = \{\mathbf{x} \rightarrow \text{sign}(h_{w,b}(\mathbf{x})) : h_{w,b} \in L_d\} \quad (2)$$

Example: $\mathcal{X} = \mathbb{R}^2$



$$d=2$$

$$w_1, w_2, b$$

$$\vec{x} = [x_1, x_2]$$

$$\vec{w} = [w_1, w_2]$$

$$b$$

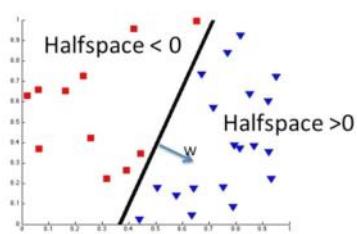
$$\vec{x}' = [1, x_1, x_2] \quad \left. \begin{array}{l} \text{AFFINE} \\ \downarrow \\ \vec{w}' = [b, x_1, x_2] \end{array} \right\} \text{HOMOGEN}$$



Examples: Halfspaces

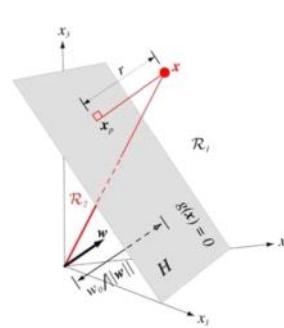
$$\mathcal{X} = \mathbb{R}^2 \text{ } (d=2)$$

(2D space divided by a line)



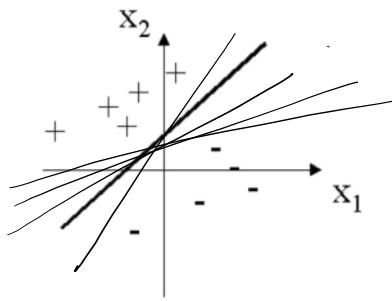
$$\mathcal{X} = \mathbb{R}^3 \text{ } (d=3)$$

(3D space divided by a plane)





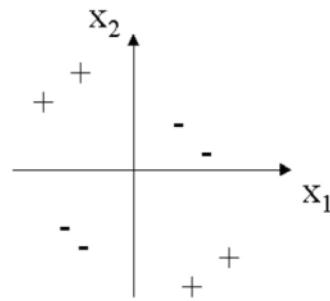
Realizability: Linearly Separable



Linearly Separable

Can separate
the data
using a line

$$\mathcal{L}_S(h) = 0$$



Not Linearly Separable

CANNOT separate
the data
with a line

can't solve
the problem
with linear
model

$$\rightarrow \mathcal{L}_S(h) \neq 0$$



Linear Programming (not part of the course)

Linear Program (LP): maximize a linear function subject to linear inequalities

Target: find $\max_{w \in \mathbb{R}^d} \langle u, w \rangle$ subject to $Aw \geq v$

$w \in \mathbb{R}^d$: vector of unknowns, $u \in \mathbb{R}^d$, $A \in \mathbb{R}^{m \times d}$, $v \in \mathbb{R}^m$

- ❑ Empirical Risk Minimization (ERM) for halfspaces in the *realizable* case can be expressed as a linear program (LP)
- ❑ All solutions satisfying constraints are ok for us (\rightarrow see SVM later...)
- ❑ There exist efficient LP solvers (e.g., simplex algorithm)

All INPUT
CORRECTLY
classified

$$\mathcal{L}_S(h_S) = 0$$

$$\text{sign}(\langle w, x_i \rangle) = y_i \quad \forall i$$

LABELS

$$y_i \langle w, x_i \rangle > 0 \quad \forall i$$

$\oplus \rightarrow$ CLASSIFY OK

$\ominus \rightarrow$ CLASSIFY X

| OUT | LABEL |
|------------------|-------|
| $\oplus \cdot +$ | OK |
| $- \cdot -$ | OK |
| $- \cdot +$ | X |
| $+ \cdot -$ | X |



Find ERM Halfspace: Perceptron



- ❑ Iterative algorithm (introduced by Rosenblatt in 1958)
- ❑ Target: find separating hyperplane
- ❑ Find vector w representing separating hyperplane (in homogeneous coordinates)
- ❑ At each step focus on a misclassified sample and guide the algorithm to be "more correct" on it
- ❑ In the realizable case always converge to a (ERM) solution correctly classifying all points
- ❑ Simple and fast in most cases (but there exists critical situations)

FIND
 \vec{w}'

change \vec{x}'
to have a more
correct solution

$L_s(a_s) = 0$

find the correct
solutions

the algo
CONVERGES

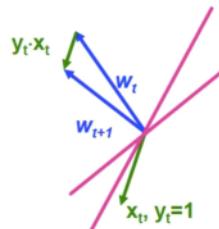


Find ERM Halfspace: Perceptron

$w^{(t)}$
↓
weight
vector
at step
 t

Input: training set $(x_1, \dots, x_m, y_1, \dots, y_m)$
initialize $w^{(1)} = (0, \dots, 0)$
for $t = 1, 2, \dots$ do
 if $\exists i$ s.t. $y_i \langle w^{(t)}, x_i \rangle \leq 0$ then
 else return $w^{(t)}$

Interpretation of update:



Termination? Depends on the realizability assumption!

→ ERM predictor $L_S(h_S) = 0$



→ OUTPUT $_i = \text{LABEL}_i + i$

→ $y_i : \langle w, x_i \rangle > 0 \quad \forall i$

⊕ → OK

⊖ → X

$$y_i = \begin{cases} +1 \\ -1 \end{cases}$$

Note that...
At each iteration find a misclassified sample and add the sample multiplied by its label

$$\begin{aligned} y_i \langle w^{(t+1)}, x_i \rangle &= y_i \langle w^{(t)} + y_i x_i, x_i \rangle \\ &= y_i \langle w^{(t)}, x_i \rangle + ||x_i||^2 \\ &\stackrel{\text{sol. step}(t)}{>} 0 \end{aligned}$$

update guides w to be "more correct" on (x_i, y_i)

|| x_i || $^2 > 0$ and target is $y_i \langle w^{(t)}, x_i \rangle > 0$: from step t to $t+1$ "more correct" on i -th sample



Perceptron on Linearly Separable Data

Halfspaces: *realizability assumption* corresponds to *linearly separable data*

Theorem:

- $(x_1, y_1), \dots, (x_m, y_m)$ is linearly separable
- $B = \min\{\|w\|: y_i \langle w, x_i \rangle \geq 1 \ \forall i = 1, \dots, m\}$
- $R = \max\|x_i\|$

TRAIN set
is L. separabile

- The perceptron algorithm stops after at most $(RB)^2$ iterations
- When it stops, it holds that $\forall i \in \{1, \dots, m\}: y_i \langle w^{(t)}, x_i \rangle > 0$
- Notice: by design, the algorithm stops when there are no more wrongly classified samples



Theorem: Demonstration

1. Define:

- w^* : vector achieving the *min* in definition of B
- T : number of iterations before stopping

$t = 1, \dots, T \rightarrow$

2. Consider: $\frac{\langle w^*, w^{(T+1)} \rangle}{\|w^*\| \|w^{(T+1)}\|} \rightarrow$ it is smaller than 1 (cosine of angle)

3. We need to demonstrate that:

$$\frac{\sqrt{T}}{RB} = \frac{T}{\sqrt{TRB}} \leq \frac{\langle w^*, w^{(T+1)} \rangle}{\|w^*\| \|w^{(T+1)}\|} \leq 1 \Rightarrow T < (RB)^2$$

4. Proceed in 2 parts:

- a) Numerator: demonstrate that $\langle w^*, w^{(T+1)} \rangle \geq T$
- b) Denominator: demonstrate that $\|w^*\| \|w^{(T+1)}\| \leq \sqrt{T} RB$

$$\left| \frac{\sqrt{T}}{RB} < 1 \right|$$

total # of steps

$w^{(T+1)}$ final weight, product of the algo

Numerator bigger

Denominator smaller



Theorem: Demonstration (Numerator)

Numerator: demonstrate that $\langle w^*, w^{(T+1)} \rangle \geq T$

- First iteration: $w^{(1)} = (0, \dots, 0) \rightarrow \langle w^*, w^{(1)} \rangle = 0$
- At each step $\langle w^*, w^{(t+1)} \rangle - \langle w^*, w^{(t)} \rangle \geq 1$ (using perceptron update rule, see *)
- After T iterations: $\langle w^*, w^{(T+1)} \rangle \geq T$ (see *)

$$\begin{aligned} (*) \langle w^*, w^{(T+1)} \rangle &= \sum_{t=1}^T (\langle w^*, w^{(t+1)} \rangle - \langle w^*, w^{(t)} \rangle) \\ &= \sum_{t=1}^T \langle w^*, w^{(t+1)} - w^{(t)} \rangle = \sum_{t=1}^T \langle w^*, y_i x_i \rangle \geq T \end{aligned}$$

(perceptron update rule) (definition of B and w^*)

$$\begin{aligned} w^{(t+1)} &= w^{(t)} + y_i x_i \\ w^{(t+1)} - w^{(t)} &= y_i x_i \end{aligned}$$

Def of B and w^*
 $\|w^*\| : y_i \langle w^*, x_i \rangle \geq 1$
: T steps
 $\geq T$



Theorem: Demonstration (Denominator)

Denominator: demonstrate that $\|w^*\| \|w^{(T+1)}\| \leq \sqrt{T} RB$

a) $\|w^{(T+1)}\|^2 \leq TR^2 \rightarrow \|w^{(T+1)}\| \leq \sqrt{TR}$ (**)

b) $\|w^*\| = B$ (by definition)

$$\begin{aligned} (\text{**}) \|w^{(T+1)}\|^2 &= \sum_{t=1}^T (\|w^{(t+1)}\|^2 - \|w^{(t)}\|^2) \\ &= \sum_{t=1}^T (\|w^{(t)} + y_t x_i\|^2 - \|w^{(t)}\|^2) \quad \text{the two } \|w^{(t)}\|^2 \text{ cancel out} \\ &= \sum_{t=1}^T (2y_t \langle w^{(t)}, x_i \rangle + \|x_i\|^2) \leq TR^2 \end{aligned}$$

Algorithm
(perceptron update condition;
select missclassified sample)

≤ 0 by algorithm
 $\leq R^2$ by algorithm
(definition of R)

$$\dots T \text{ steps} \leq TR^2$$



Perceptron: Notes

It is simple to implement.

On separable data:

- Convergence is guaranteed
- Convergence depends on B , which can be exponential in d
 - If the input vectors are not normalized and arranged in some unfavorable ways the running time can be very long
 - A Linear Programming (LP) approach may be better to find ERM solution in some cases
- Potentially multiple solutions, which one is picked depends on starting values

→ better have
normalized
input
vectors

On non separable data:

- Run for some time and keep best solution found up to that point (pocket algorithm)

Learning Example

Initial Values:

$$\eta = 0.2$$

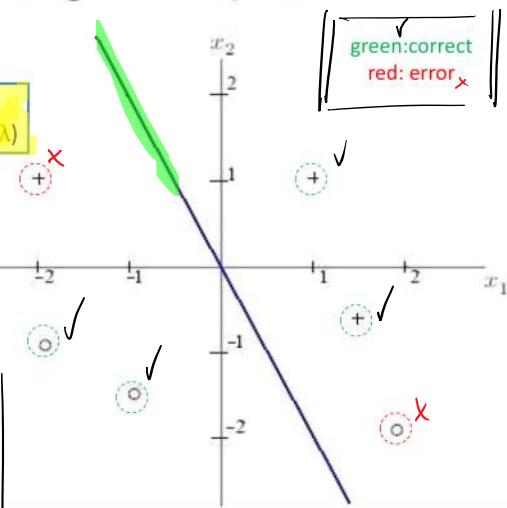
Learning rate
(also denoted with λ)

$$w = \begin{pmatrix} 0 \\ 1 \\ 0.5 \end{pmatrix}$$

bias = 6
Initial line position

$$\left| \begin{array}{l} 0 = w_0 + w_1 x_1 + w_2 x_2 \\ = 0 + x_1 + 0.5x_2 \\ \Rightarrow x_2 = -2x_1 \end{array} \right|$$

Perceptron with learning rate: $w^{(t+1)} = w^{(t)} + \lambda y_i x_i$



η or $\lambda \rightarrow$ to slow and stabilize the algd

Learning Example

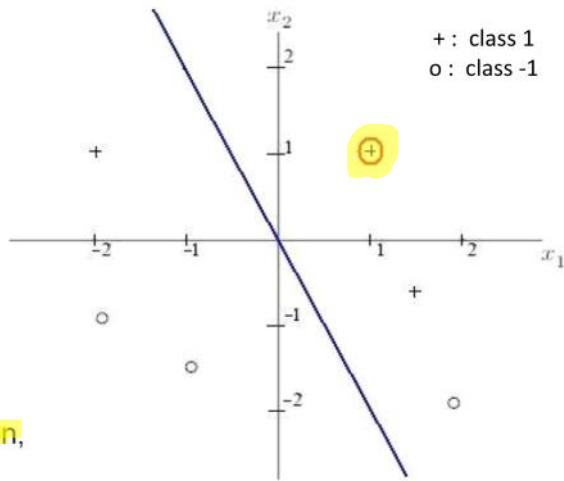
$$\eta = 0.2$$

$$w = \begin{pmatrix} 0 \\ 1 \\ 0.5 \end{pmatrix}$$

$$x_1 = 1, x_2 = 1$$

$$w^T x > 0$$

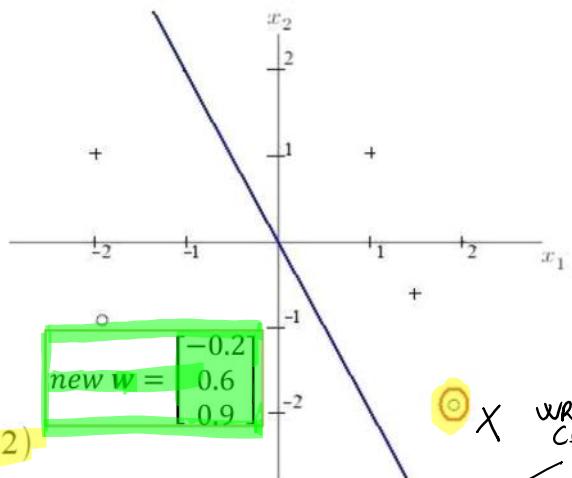
Correct classification,
no action.



Learning Example

$$\eta = 0.2$$

$$w = \begin{pmatrix} 0 \\ 1 \\ 0.5 \end{pmatrix}$$



$$x_1 = 2, x_2 = -2$$

$$w_0 = w_0 - 0.2 * 1$$

$$w_1 = w_1 - 0.2 * 2$$

$$w_2 = w_2 - 0.2 * (-2)$$

$$\text{new } w = \begin{pmatrix} -0.2 \\ 0.6 \\ 0.9 \end{pmatrix}$$

need to update weight w

WRONG Classify

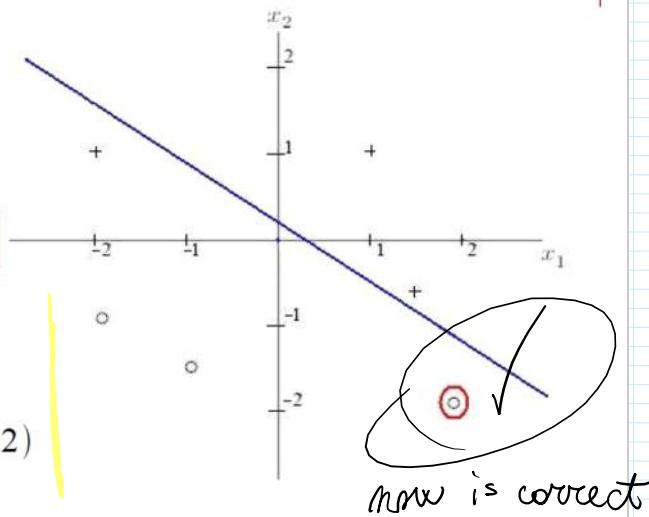
Learning Example

$$\eta = 0.2$$

$$w = \begin{pmatrix} -0.2 \\ 0.6 \\ 0.9 \end{pmatrix}$$

w = $\begin{pmatrix} 0 \\ 1 \\ 0.5 \end{pmatrix}$

$$\begin{aligned} x_1 &= 2, x_2 = -2 \\ w_0 &= w_0 - 0.2 * 1 \\ w_1 &= w_1 - 0.2 * 2 \\ w_2 &= w_2 - 0.2 * (-2) \end{aligned}$$



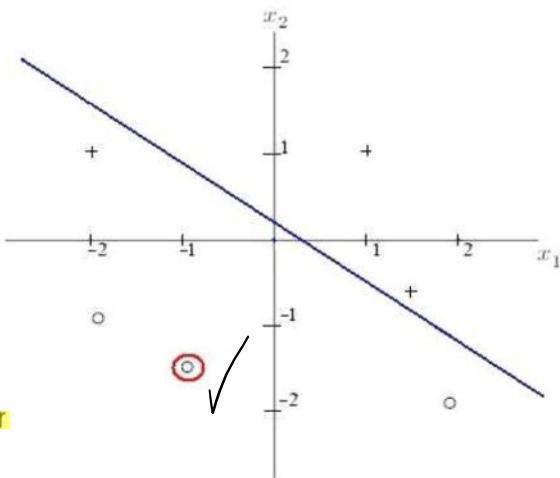
Learning Example

$$\eta = 0.2$$

$$w = \begin{pmatrix} -0.2 \\ 0.6 \\ 0.9 \end{pmatrix}$$

$$x_1 = -1, x_2 = -1.5$$
$$w^T x < 0$$

Correct classification
no action



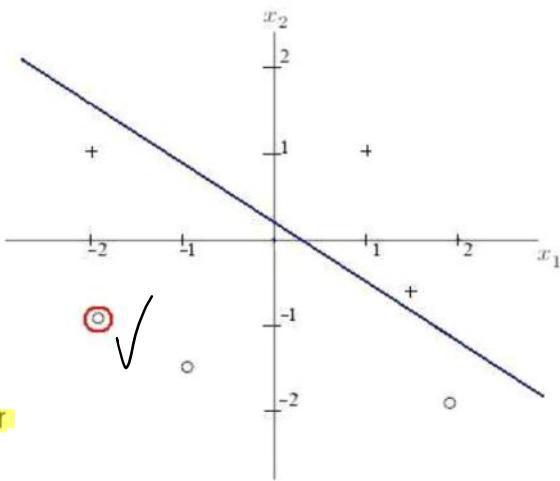
Learning Example

$$\eta = 0.2$$

$$w = \begin{pmatrix} -0.2 \\ 0.6 \\ 0.9 \end{pmatrix}$$

$$x_1 = -2, x_2 = -1$$
$$w^T x < 0$$

Correct classification
no action



Learning Example

$$\eta = 0.2$$

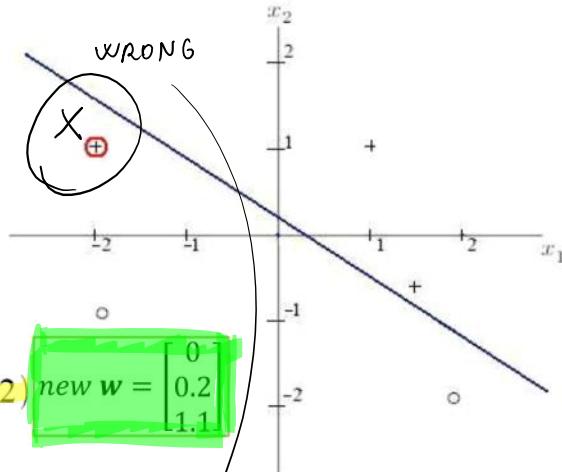
$$w = \begin{pmatrix} -0.2 \\ 0.6 \\ 0.9 \end{pmatrix}$$

$$x_1 = -2, x_2 = 1$$

$$w_0 = w_0 + 0.2 * 1$$

$$w_1 = w_1 + 0.2 * (-2)$$

$$w_2 = w_2 + 0.2 * 1$$



need to update
 w

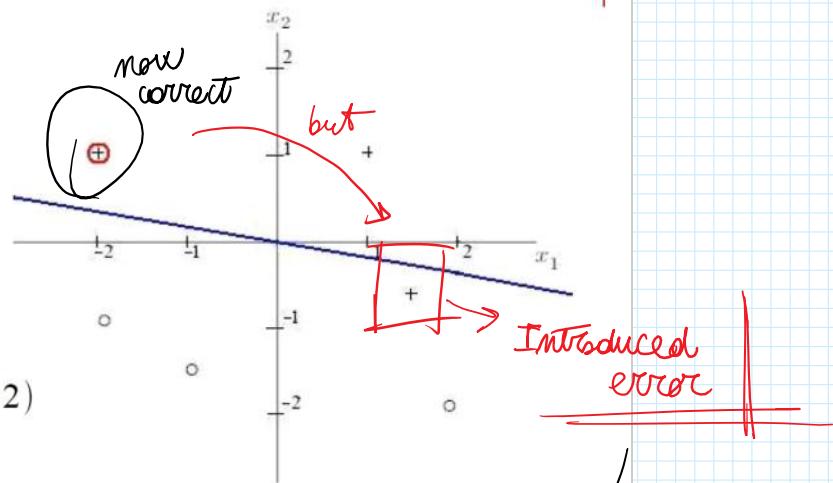
Learning Example

$$\eta = 0.2$$

$$w = \begin{pmatrix} 0 \\ 0.2 \\ 1.1 \end{pmatrix}$$

previous w
 $w = \begin{pmatrix} -0.2 \\ 0.6 \\ 0.9 \end{pmatrix}$

$$\begin{aligned} x_1 &= -2, x_2 = 1 \\ w_0 &= w_0 + 0.2 * 1 \\ w_1 &= w_1 + 0.2 * (-2) \\ w_2 &= w_2 + 0.2 * 1 \end{aligned}$$

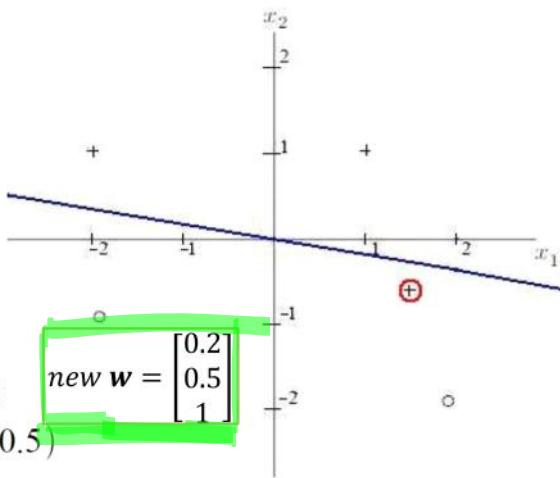


Learning Example

$$\eta = 0.2$$

$$w = \begin{pmatrix} 0 \\ 0.2 \\ 1.1 \end{pmatrix}$$

$$\begin{aligned} x_1 &= 1.5, x_2 = -0.5 \\ w_0 &= w_0 + 0.2 * 1 \\ w_1 &= w_1 + 0.2 * 1.5 \\ w_2 &= w_2 + 0.2 * (-0.5) \end{aligned}$$



Learning Example

$$\eta = 0.2$$

$$w = \begin{pmatrix} 0.2 \\ 0.5 \\ 1 \end{pmatrix}$$

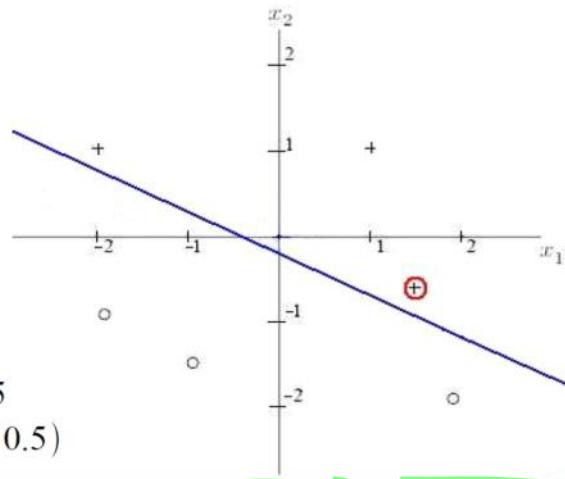
w previous w
 $\boxed{\begin{pmatrix} 0 \\ 1.1 \end{pmatrix}}$

$$x_1 = 1.5, x_2 = -0.5$$

$$w_0 = w_0 + 0.2 * 1$$

$$w_1 = w_1 + 0.2 * 1.5$$

$$w_2 = w_2 + 0.2 * (-0.5)$$



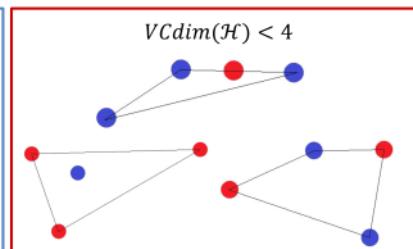
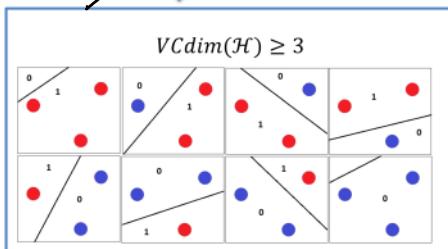


VC Dimension of Halfspaces

VC-Dim of linear functions (hyperplane)

- In the plane (hyperplane = line):
 $\rightarrow 2D \Rightarrow d=2$
- VC(Hyperplanes) is at least 3
 - VC(Hyperplanes) < 4 since there is no set of 4 points, which can be shattered by a line.
- $\Rightarrow VC(H)=3$. In general, for a k-dimension space $VC(H)=k+1$
- NB: It is useless selecting a set of linearly independent points

- For the halfspace hypothesis class:
- The VC dimension of the class of homogenous halfspaces in \mathbb{R}^d is d
 - The VC dimension of the class of non-homogenous halfspaces in \mathbb{R}^d is $d+1$



Have to demonstrate that exists a set of size d that can be shattered by the hyp class

$$VC \geq d$$

Every set of size $d+1$ can't be shattered

$$VC < d+1 \rightarrow \text{contradiction}$$

$$\underline{\underline{VC = d}}$$



VC Dimension of Halfspaces: Demonstration (1)

The VC dimension of the class of homogenous halfspaces in \mathbb{R}^d is d

Demonstration (homogenous case):

set of size d that can be
shattered

a) $VCdim(HS_d) \geq d$

1. Consider the set e_1, \dots, e_d where $\forall i: e_i = (0, \dots, 0, 1, 0, \dots, 0)$
 - i.e., all "0" except "1" in the i -th coordinate
2. The set is shattered by the homogeneous halfspaces: to obtain labeling e_1, \dots, e_d set $w = (y_1, \dots, y_d) \Rightarrow \langle w, e_i \rangle \geq y_i \quad \forall i$
 - for each vector only the multiplication with the corresponding label is $\neq 0$



VC Dimension of Halfspaces: Demonstration (2)

b) $VCdim(HS_d) < d + 1$

- x_1, \dots, x_{d+1} generic set of $d+1$ vectors in \mathbb{R}^d
- They must be linearly dependent:

$$\exists a_1, \dots, a_{d+1} \in \mathbb{R} \text{ (not all zero): } \sum_{i=1}^{d+1} a_i x_i = 0 \quad (*)$$

3. Define $I = \{i: a_i > 0\}$, $J = \{j: a_j < 0\}$: either I or J are non-empty

4. Assume both non-empty: $\sum_{i \in I} a_i x_i = \sum_{j \in J} |a_j| x_j \quad (*_2)$

5. By contradiction: assume that the set is shattered: \exists a vector w such that $\langle w, x_i \rangle > 0 \forall i \in I$ and $\langle w, x_j \rangle < 0 \forall j \in J$

6. It follows a contradiction :

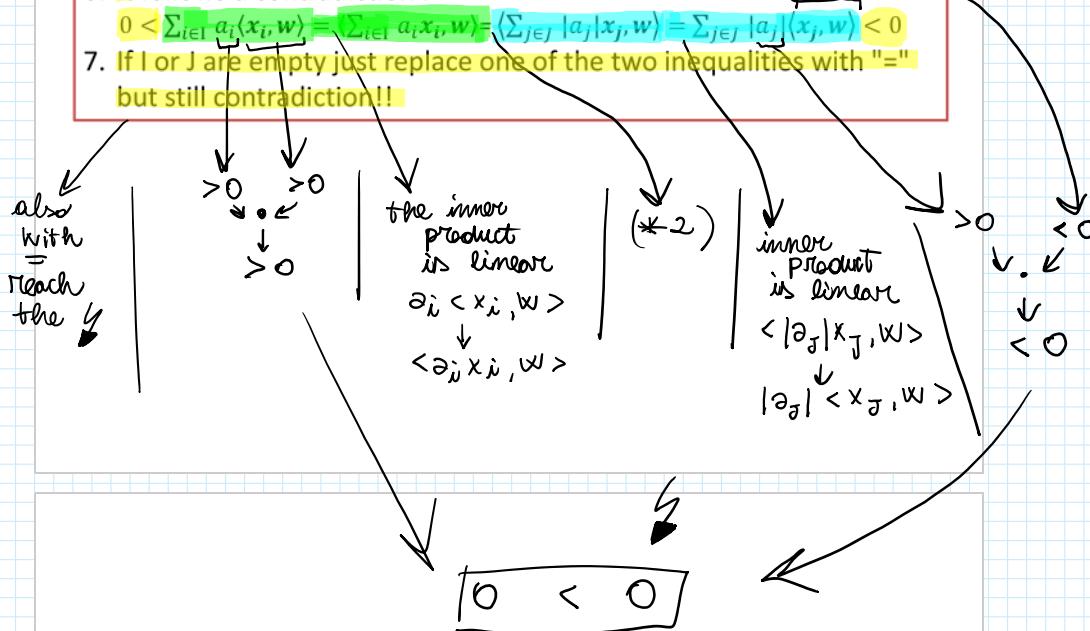
$$0 < \sum_{i \in I} a_i \langle x_i, w \rangle = \langle \sum_{i \in I} a_i x_i, w \rangle = \langle \sum_{j \in J} |a_j| x_j, w \rangle = \sum_{j \in J} |a_j| \langle x_j, w \rangle < 0$$

7. If I or J are empty just replace one of the two inequalities with "=" but still contradiction!!

$d+1$ vectors in \mathbb{R}^d

→ linearly dep

↓
(based formed
by d elements)



Linear Regression

Linear regression: estimate the relation between some explanatory variables and some real valued outcome

- Domain set: $x \in \mathbb{R}^d$, label set: $y = \mathbb{R}$
- Find $h \in \mathcal{H}_{reg}: \mathbb{R}^d \rightarrow \mathbb{R}$ that best approximates the relation between input and output

$$(b, \vec{w}) \quad (1, \vec{x})$$

- Find $h \in \mathcal{H}_{reg}$: $\mathbb{R}^d \rightarrow \mathbb{R}$ that best approximates the relation between input and output.
- Hypothesis class:

$$\mathcal{H}_{reg} = L_d = \{x - \langle w, x \rangle + b : w \in \mathbb{R}^d, b \in \mathbb{R}\}$$

- Loss function: squared loss (L2) is commonly used but other functions are possible (e.g., mean absolute error)

$$\ell(h, (x, y)) \stackrel{\text{def}}{=} (h(x) - y)^2$$

- Empirical Risk function: Mean Squared Error on training set

$$L_s(h) = \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)^2$$

$$(b, \vec{w}) \rightarrow (1, \vec{x})$$

the output (correct/
incorrect)
is not ± 1 anymore

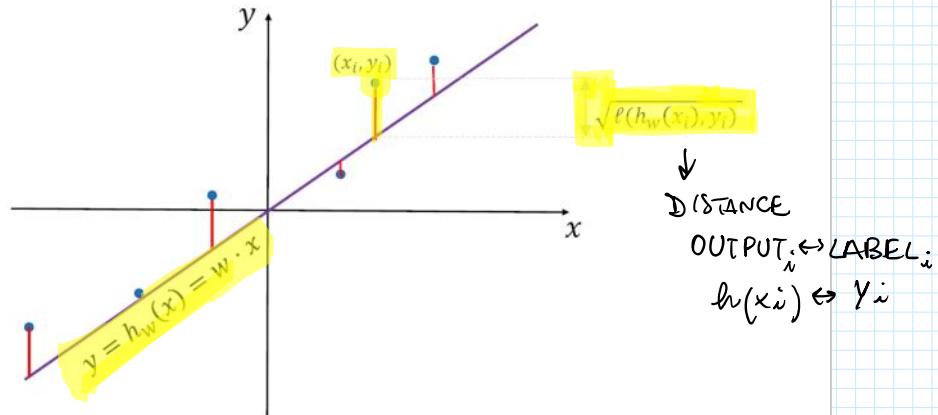
can be more correct
or less correct

another
loss and
error



Linear Regression

$$\mathcal{X} = \mathbb{R}^1 \quad \mathcal{Y} = \mathbb{R}$$

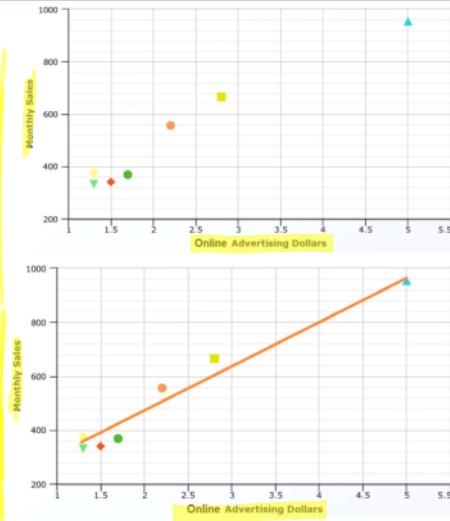




DIPARTIMENTO
DI INGEGNERIA
DELL'INFORMAZIONE

Example: Linear Regression (1)

| Online Store | Monthly Sales (in 1000 \$) | Online Advertising Dollars (1000 \$) |
|--------------|----------------------------|--------------------------------------|
| 1 | 368 | 1.7 |
| 2 | 340 | 1.5 |
| 3 | 665 | 2.8 |
| 4 | 954 | 5.0 |
| 5 | 331 | 1.3 |
| 6 | 556 | 2.2 |
| 7 | 376 | 1.3 |



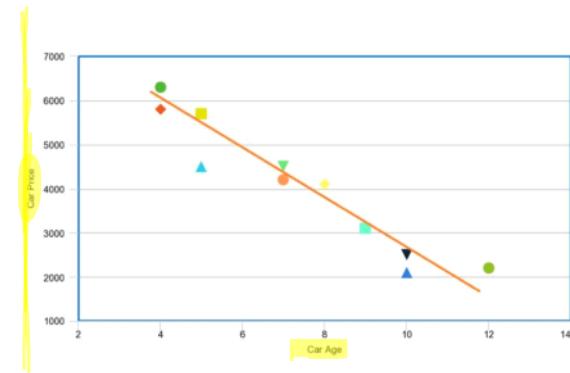


DIPARTIMENTO
DI INGEGNERIA
DELL'INFORMAZIONE

Example:

Linear Regression (2)

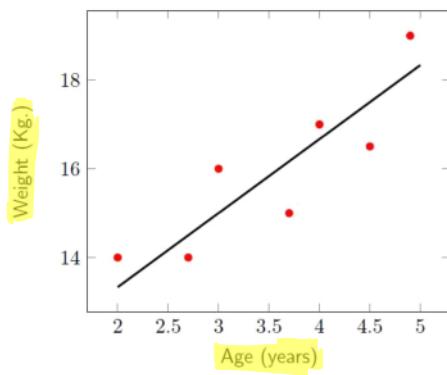
| Car Age (years) | Price (€) |
|--------------------|--------------|
| 4 | 6300 |
| 4 | 5800 |
| 5 | 5700 |
| 5 | 4500 |
| 7 | 4500 |
| 7 | 4200 |
| 8 | 4100 |
| 9 | 3100 |
| 10 | 2100 |
| 11 | 2500 |
| 12 | 2200 |





Example: Linear Regression (3)

- $\mathcal{X} \subset \mathbb{R}^d, \mathcal{Y} \subset \mathbb{R}, \mathcal{H} = \{x \mapsto \langle w, x \rangle : w \in \mathbb{R}^d\}$
- Example: $d = 1$, predict weight of a child based on his age.





Least Squares

$$\arg \min_{\mathbf{w}} L_S(h_{\mathbf{w}}) = \arg \min_{\mathbf{w}} \frac{1}{m} \sum_{i=1}^m (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)^2$$

- The *least squares* algorithm solves the ERM problem for linear regression predictors with the squared loss
- Find the parameters vector that minimize the MSE between the estimated and training values
- To solve the problem: calculate gradient w.r.t vector \mathbf{w} and set to 0

want to find
the \vec{w}
that minimize



Least Squares: Solution

- Compute gradient w.r.t w and set to 0

$$\frac{\partial L}{\partial w} = \frac{2}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i) x_i = 0 \rightarrow \sum_{i=1}^m \langle w, x_i \rangle x_i = \sum_{i=1}^m y_i x_i$$

$$\arg \min_w L_S(w) = \arg \min_w \frac{1}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2$$

MATRIX FORM

$$A = \left(\sum_{i=1}^m x_i x_i^T \right) \quad \begin{matrix} \rightarrow m \\ \downarrow \\ d \end{matrix} \quad \begin{matrix} \rightarrow d \\ \downarrow \\ m \end{matrix} \quad \begin{matrix} \cdots \\ x_1 \\ \vdots \\ x_m \\ \cdots \\ x_m \end{matrix} \quad \begin{matrix} m \\ \downarrow \\ m \end{matrix}$$

$$b = \sum_{i=1}^m y_i x_i \quad \begin{matrix} \downarrow \\ d \times m \\ \uparrow \\ m \end{matrix} \quad \begin{matrix} m \\ \downarrow \\ \cdots \\ x_1 \\ \vdots \\ x_m \\ \cdots \\ y_m \end{matrix} \quad \begin{matrix} y_1 \\ \vdots \\ y_m \end{matrix} \quad \begin{matrix} m \\ \downarrow \\ m \end{matrix}$$

- The solution is:

$$Aw = b \rightarrow w = A^{-1}b$$

- Inversion of A is the most expensive operation

- The case in which A is not invertible requires a special handling (not part of the course)

$$\sum_{i=1}^m \langle w, x_i \rangle x_i = \sum_{i=1}^m y_i x_i$$

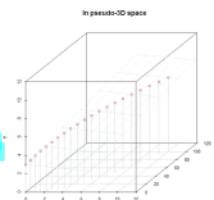
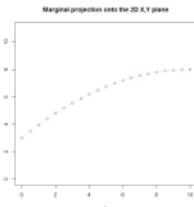
$$\hat{A} \cdot w = b \Rightarrow w = \hat{A}^{-1} \cdot b$$

$$\hat{A} \rightarrow d \times d \text{ matrix} \quad \parallel \quad \vec{b}, \vec{w} \rightarrow 1 \times d$$



Polynomial Regression

- ❑ Polynomial regression: find the one dimensional polynomial of degree n that better predicts the data
 - $p(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$
 - Need to estimate the coefficient vector a
 - 1D polynomial pred. deg. n : $\mathcal{H}_{poly}^n = \{x \rightarrow p(x)\}, X = \mathbb{R}, Y = \mathbb{R}$
- ❑ Reduce the problem to a n -dimensional linear regression using the mapping:
- ❑ We obtain:
$$\langle a, \psi(x) \rangle = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$$
- ❑ Find the vector of coefficients a using the Least Square algorithm
- ❑ Non-linear relation becomes linear in the higher dimensional space
- ❑ Notice that the variables are not independent
- ❑ The optimization can become unstable for large n



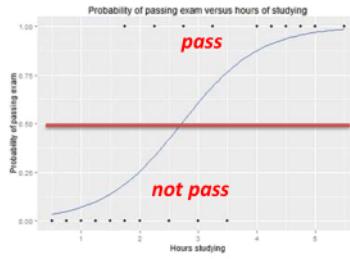
$$x = 3$$
$$\psi(x) = (1, 3, 9, 27, \dots)$$
$$x^0, x^1, x^2, x^3 \dots$$

each dimension is correlated a different power

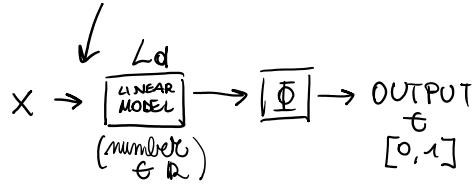


Logistic Regression

classification method



- ❑ Reframe a classification problem as a regression one
- ❑ Target as in regression:
 - learn a function $h: \mathbb{R}^d \rightarrow [0; 1]$
 - the output of h is a real number
- ❑ Used for classification:
 - interpret the output of h as the probability that the label is 1
 - regression-like output for classification!
- ❑ We'll deal with binary classification but the approach can be extended to the multi-class setting
- ❑ Hypothesis class $\mathcal{H}: \phi_{\text{sig}} \circ L_d$ where $\phi: \mathbb{R} \rightarrow [0, 1]$ is the sigmoid function and L_d an affine function



$$\left| \begin{array}{l} P(1) \rightarrow \\ > 0,5 \rightarrow \text{OUT} \rightarrow 1 \\ < 0,5 \rightarrow \text{OUT} \rightarrow -1 \end{array} \right.$$

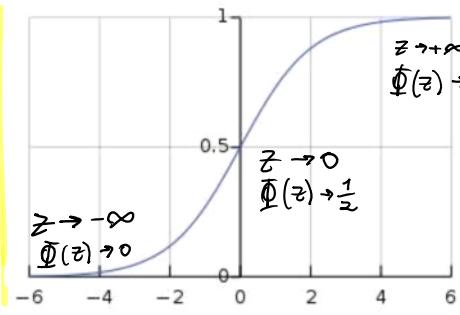
$$h: \mathbb{R}^d \rightarrow [0, 1]$$

↓
a number between 0 and 1
es. ↓

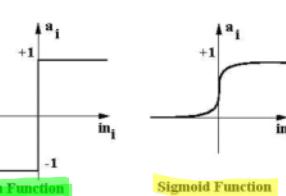
$$x_i \in \mathbb{R}^d \rightarrow P(y_i = 1)$$



Sigmoid Function



$$\phi_{sig}(z) = \frac{1}{1 + e^{-z}}$$



- Bigger than $\frac{1}{2}$ for positive values and smaller for negative ones
- Tends to 1 for large positive values and to 0 for negative ones
- Can be viewed as a scaled and shifted "soft" sign function

Can set sort of thresholds to the outputs
outcomes

≈

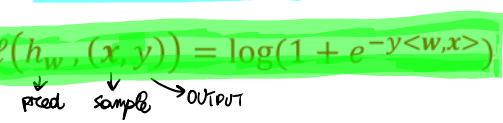
$P(x) \geq 0,7 \rightarrow \text{OUT } 1 \text{ (MALE)}$

$0,3 < P(x) < 0,7 \rightarrow \text{UNSURE}$

$P(x) \leq 0,3 \rightarrow \text{OUT } -1 \text{ (FEMALE)}$



Loss for Logistic Regression

- ❑ Instead of hard choice → use probability of correct label being 0 or 1
- ❑ $H_{sig} = \phi_{sig} \circ L_d = \{x \rightarrow \phi_{sig}(\langle w, x \rangle) : w \in \mathbb{R}^d\}$
- ❑ Hypothesis class: $h_w(x) = \frac{1}{1+e^{-\langle w, x \rangle}}$
- ❑ Loss function: $\ell(h_w, (x, y)) = \log(1 + e^{-y\langle w, x \rangle})$

- ❑ ERM Problem: $\operatorname{argmin}_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y_i \langle w, x_i \rangle})$



Logistic Loss Function

- Loss function: $\ell(h_w(x), y) = \log(1 + e^{-y \langle w, x \rangle})$, why?

- Consider $h_w(x) = \frac{1}{1+e^{-(w,x)}} \leftrightarrow y \in \{+1, -1\}$

- Case $y=1$: need $h_w(x) \rightarrow 1$

$$h_w(x) = \frac{1}{1+e^{-(w,x)}} = \frac{1}{1+e^{-y(w,x)}}$$

- If denominator small $h_w(x) \rightarrow 1$ good case
- If denominator large $h_w(x) \rightarrow 0$ error

- Case $y=-1$: need $h_w(x) \rightarrow 0 \Rightarrow 1 - h_w(x) \rightarrow 1$

$$1 - h_w(x) = 1 - \frac{1}{1+e^{-(w,x)}} = \frac{1+e^{-(w,x)} - 1}{1+e^{-(w,x)}} = \frac{e^{-(w,x)}}{e^{-(w,x)} + 1} = \frac{1}{1+e^{-y(w,x)}}$$

- Same as before!
- If denominator small $1 - h_w(x) \rightarrow 1$ good case, if large $1 - h_w(x) \rightarrow 0$ error

Loss need to increase with $1 + e^{-y(w,x)}$ and log is monotonic

$$1 + e^{-y \langle w, x \rangle} = \begin{cases} 0 & \text{good} \\ \infty & \text{error} \end{cases}$$

log is monotonic

(also simple to differentiate)



Maximum Likelihood Estimation (MLE)

Maximum Likelihood Estimation (MLE) is a statistical approach for finding the parameters that maximize the joint probability of a given dataset assuming a specific parametric probability function

- ❑ MLE essentially assumes a generative model for the data
- ❑ MLE solution is equivalent to ERM solution for logistic regression

MLE approach:

1. Given training set $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m))$, assume each (\mathbf{x}_i, y_i) is i.i.d. from some probability distribution (that is characterized by some parameters)
2. Consider $P[S|\theta]$ (likelihood of data given parameters)
3. log likelihood: $L(S; \theta) = \log(P[S|\theta])$
 - \log : monotonic \rightarrow same maximum, but simpler to differentiate
4. **Maximum Likelihood Estimator (MLE):** $\hat{\theta} = \operatorname{argmax}_{\theta} L(S; \theta)$

(NOT PART
of the
course)



MLE and Logistic Regression

MLE solution is equivalent to ERM solution for logistic regression

Logistic Regression:

1. Assume training set $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m))$
2. $P[y_i = 1] = h_{\mathbf{w}}(\mathbf{x}_i) = \frac{1}{1+e^{-\langle \mathbf{w}, \mathbf{x}_i \rangle}} = \frac{1}{1+e^{-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle}}$ (since $y_i = 1$)
3. $P[y_i = -1] = 1 - h_{\mathbf{w}}(\mathbf{x}_i) = \frac{1}{1+e^{\langle \mathbf{w}, \mathbf{x}_i \rangle}} = \frac{1}{1+e^{-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle}}$ (first equality recall logistic loss, 2nd since $y_i = -1$)
4. Likelihood of training set (joint probability $P[S|\mathbf{w}]$): $\prod_{i=1}^m \left(\frac{1}{1+e^{-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle}} \right)$
5. Log likelihood: $\log(P[S|\mathbf{w}]) = \log \prod_{i=1}^m \left(\frac{1}{1+e^{-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle}} \right) = -\sum_{i=1}^m \log(1 + e^{-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle})$
 \rightarrow corresponds to (-1)*logistic loss

Maximum Likelihood Estimator:

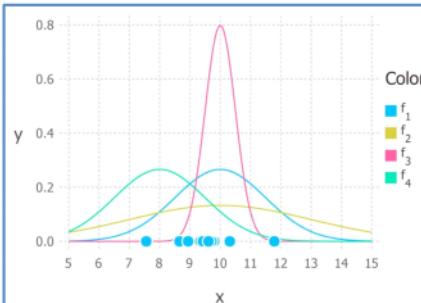
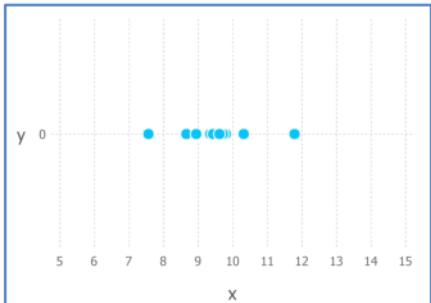
$$\operatorname{argmax}_{\mathbf{w}} L(S; \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \log(P[S|\mathbf{w}]) = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^m \log(1 + e^{-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle})$$

Recall: $\operatorname{argmax}(-x) = \operatorname{argmin}(x)$ They have the same target !!!



DIPARTIMENTO
DI INGEGNERIA
DELL'INFORMAZIONE

Example: MLE for Gaussian PDF (1)



Assume that the data is produced
by a Gaussian distribution

$$P(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

$f_1 \sim N(10, 2.25)$ $f_2 \sim N(10, 9)$,
 $f_3 \sim N(10, 0.25)$ $f_4 \sim N(8, 2.25)$

find μ, σ maximizing the joint
probability of data



Example: MLE for Gaussian PDF (2)

Joint probability

$$P(9, 9.5, 11; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(9-\mu)^2}{2\sigma^2}\right) \times \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(9.5-\mu)^2}{2\sigma^2}\right) \times \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(11-\mu)^2}{2\sigma^2}\right)$$

Assume
3 samples:
9, 9.5, 11

Log likelihood

$$\ln(P(x; \mu, \sigma)) = \ln\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \frac{(9-\mu)^2}{2\sigma^2} + \ln\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \frac{(9.5-\mu)^2}{2\sigma^2} + \ln\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \frac{(11-\mu)^2}{2\sigma^2}$$

Recall
 $\log(ab) = \log(a) + \log(b)$

Get optimal mean

$$\frac{\partial \ln(P(x; \mu, \sigma))}{\partial \mu} = \frac{1}{\sigma^2} [9 + 9.5 + 11 - 3\mu] = 0$$

The same can be done for σ

Differentiate w.r.t μ and set to 0

$$\mu = \frac{9 + 9.5 + 11}{3} = 9.833$$

set derivative to 0