



Università degli Studi di Padova



Linear Predictors

Machine Learning 2022-23 UML Book Chapter 9 Slides: F. Chiariotti, 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 a hypothesis class
- Select an algorithm to find the predictor

Hypotheses Classes

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

Algorithms

- Linear Programming (for halfspaces)
- Perceptron (for halfspaces)
- Least Squares (for regression)

Affine Function Model

Class of *Affine Functions*:

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

where

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

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

- It is a linear function followed by a sum
- Two parameters: b (scalar, called bias) and w (vector)
- Dimensionality of parameters space: *d*+1

$$x \rightarrow L_d \xrightarrow{z} \phi \rightarrow y$$

Hypothesis class: $\phi \circ L_d \quad \phi \colon \mathbb{R} \to \mathcal{Y}$

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

Geometric Interpretation (2D)





- 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 $w: w \to w' = \langle b, w_1, w_2, \dots, w_d \rangle$
- Add an extra element to each vector \mathbf{x} : $\mathbf{x} \to \mathbf{x}' = \langle \mathbf{1}, x_1, x_2, \dots, x_d \rangle$

Homogeneous linear function:

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

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

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

Halfspaces Hypotheses Class

Halfspace hypothesis class

$$x \rightarrow L_d \xrightarrow{z} sign \rightarrow y$$

□ Input: X = ℝ^d (for each sample a vector of features)
 ■ Using homogeneous coordinates: x → x' = (1, x₁, x₂, ..., x_d) ∈ ℝ^{d+1}
 □ Output: Y = {-1,1} (binary classification)
 □ Loss: 0-1 loss

Halfspace Model:

 $HS_d = \operatorname{sign} \circ L_d = \{x \to sign(\langle \boldsymbol{w}, \boldsymbol{x} \rangle + b)\}, \ \boldsymbol{w} \in \mathbb{R}^d, b \in \mathbb{R}$

Using homogeneous coordinates:

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

Linear Classification: Halfspaces Hypotheses Class

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

Hypothesis class = halfspaces

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

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





Examples: Halfspaces

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

(2D space divided by a line)

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

(3D space divided by a plane)







Realizability: Linearly Separable





Linearly Separable

Not Linearly Separable



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 $Aw \ge v$

- $\pmb{w} \in \mathbb{R}^d$: vector of unknowns, $\pmb{u} \in \mathbb{R}^d$, $\mathbf{A} \in \mathbb{R}^{mxd}$, $\pmb{v} \in \mathbb{R}^m$
- Empirical Risk Minimization (ERM) for halfspaces in the *realizable* case can be expressed as a linear program (LP)
 - The ERM predictor is $L_s(h_s) = 0 \rightarrow \operatorname{sign}(\langle \boldsymbol{w}, \boldsymbol{x_i} \rangle) = y_i \ \forall i \rightarrow y_i \langle \boldsymbol{w}, \boldsymbol{x_i} \rangle > 0 \ \forall i$
 - Need some math to adapt ">" to " \geq "
- □ All solutions satisfying constraints are ok for us (\rightarrow see SVM later...)
- There exist efficient LP solvers (e.g., simplex algorithm)



Find ERM Halfspace: Perceptron





- Iterative algorithm (introduce 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 ERM Halfspace: Perceptron

ERM predictor: $L_s(h_s) = 0$ **Input:** training set $(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_m, y_m)$ \rightarrow sign $(\langle \boldsymbol{w}, \boldsymbol{x}_i \rangle) = y_i \forall i$ initialize $w^{(1)} = (0, ..., 0);$ $\rightarrow y_i \langle \boldsymbol{w}, \boldsymbol{x_i} \rangle > 0 \ \forall i$ for t = 1, 2, ... do if $\exists i \ s.t. \ y_i \langle \mathbf{w}^{(t)}, \mathbf{x}_i \rangle \leq 0$ then $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + y_i \mathbf{x}_i$; else return $\mathbf{w}^{(t)}$; At each iteration find a Interpretation of update: misclassified sample and add Note that: the sample multiplied by its label $y_i \langle \mathbf{w}^{(t+1)}, \mathbf{x}_i \rangle = y_i \langle \mathbf{w}^{(t)} + y_i \mathbf{x}_i, \mathbf{x}_i \rangle$ **y**_t·**x**_t $= y_i \langle \mathbf{w}^{(t)}, \mathbf{x}_i \rangle + ||\mathbf{x}_i||^2$ \Rightarrow update guides w to be "more" W1+1 correct" on (\mathbf{x}_i, y_i) . $\|\boldsymbol{x}_i\|^2 > 0$ and target is $y_i(\boldsymbol{w}^{(t)}, \boldsymbol{x}_i) > 0$: from $x_t, y_t=1$ step t to t+1 "more correct" on i-th sample Termination? Depends on the realizability assumption!



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
- $B = \min\{||\mathbf{w}||: y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \ge 1 \quad \forall i = 1, ..., m\}$
- $R = max \|\boldsymbol{x}_i\|$

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



Theorem: proof

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

$$\frac{\sqrt{T}}{RB} = \frac{T}{\sqrt{T}RB} \le \frac{\left\langle \boldsymbol{w}^*, \boldsymbol{w}^{(T+1)} \right\rangle}{\|\boldsymbol{w}^*\| \|\boldsymbol{w}^{(T+1)}\|} \le 1 \Rightarrow T < (RB)^2$$

- 4. Proceed in 2 parts:
 - a) Numerator: demonstrate that $\langle w^*, w^{(T+1)} \rangle \geq T$
 - b) Denominator: demonstrate that $\|\boldsymbol{w}^*\| \|\boldsymbol{w}^{(T+1)}\| \leq \sqrt{T} RB$

Theorem: proof (Numerator)

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

- First iteration : $w^{(1)} = (0, ..., 0) \rightarrow \langle w^*, w^{(1)} \rangle = 0$
- At each step $\langle w^*, w^{(t+1)} \rangle \langle w^*, w^{(t)} \rangle \ge 1$ (using perceptron update rule and recalling definition of w^* , see *)

• After T iterations: $\langle w^*, w^{(T+1)} \rangle \ge T$ (see *)





Theorem: proof (Denominator)

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

a)
$$\|w^{(T+1)}\|^2 \leq TR^2 \rightarrow \|w^{(T+1)}\| \leq \sqrt{T}R$$
 (**)
b) $\|w^*\| = B$ (by definition)

$$(**) \|w^{(T+1)}\|^{2} = \sum_{t=1}^{T} \left(\|w^{(t+1)}\|^{2} - \|w^{(t)}\|^{2} \right)$$

$$= \sum_{t=1}^{T} \left(\|w^{(t)} + y_{i}x_{i}\|^{2} - \|w^{(t)}\|^{2} \right)$$

$$= \sum_{t=1}^{T} \left(2y_{i}\langle w^{t}, x_{i}\rangle + \|x_{i}\|^{2} \right) \leq TR^{2}$$

$$\leq 0 \text{ by algorithm}$$

$$(\text{perceptron update condition:} \quad (\text{definition of R})$$

$$R = max \|x_{i}\|$$

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

On non separable data:

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





















VC Dimension of Halfspaces





VC Dimension of Halfspaces: Demonstration (1)

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

Demonstration (homogenous case):

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

VC Dimension of Halfspaces: Demonstration (2)

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

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

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

- 3. Define $I = \{i: a_i > 0\}$, $J = \{j: a_j < 0\}$: either *I* or *J* are non-empty
- 4. Assume both non-empty: $\sum_{i \in I} a_i x_i = \sum_{j \in J} |a_j| x_j$
- 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 :

 $0 < \sum_{i \in \mathbf{I}} a_i \langle \mathbf{x}_i, \mathbf{w} \rangle = \langle \sum_{i \in \mathbf{I}} a_i \mathbf{x}_i, \mathbf{w} \rangle = \langle \sum_{j \in J} |a_j| \langle \mathbf{x}_j, \mathbf{w} \rangle = \sum_{j \in J} |a_j| \langle \mathbf{x}_j, \mathbf{w} \rangle < 0$

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

Linear Regression

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

 \Box Domain set : $\mathcal{X} \in \mathbb{R}^d$, label set : $\mathcal{Y} = \mathbb{R}$

□ Find $h \in \mathcal{H}_{reg}$: $\mathbb{R}^d \to \mathbb{R}$ that best approximates the relation between input and output

□ Hypothesis class (linear regression):

 $\mathcal{H}_{reg} = L_d = \{ \boldsymbol{x} \to < \boldsymbol{w}, \boldsymbol{x} > +b : \boldsymbol{w} \in \mathbb{R}^d, b \in \mathbb{R} \}$

□ Loss function: squared loss (L2) is commonly used but other functions are possible (e.g., mean absolute error) $\ell(h, (x, y)) \stackrel{\text{def}}{=} (h(x) - y)^2$

D Empirical Risk function: *Mean Squared Error* on training set $L_s(h) = \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)^2$

Linear Regression

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



Example: Linear Regression (1)

DIPARTIMENTO DI INGEGNERIA DELL'INFORMAZIONE

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	



Example: Linear Regression (2)

12

14

DIPARTIMENTO DI INGEGNERIA DELL'INFORMAZIONE

Car Age (years)	Price (€)	
4	6300	
4	5800	
5	5700	
5	4500	7000
7	4500	6000
7	4200	5000 -
8	4100	A ▲
9	3100	Car
10	2100	3000
11	2500	2000
12	2200	
L		

Example: Linear Regression (3)

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



Prediction quality

- 1. Linear regression gets us a good prediction
- 2. We have a good prediction of the average, but a high error variance
- E[Y|X]=E[Y] (linear regression does not tell us much...)
- 4. Linear regression is misleading

We can intuitively see the good candidates for (simple) linear regression





Least Squares

$$\underset{w}{\arg\min L_s(h_w)} = \underset{w}{\arg\min \frac{1}{m}} \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2$$

- 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 w and set to 0

Least Squares: Solution

- Compute gradient w.r.t \boldsymbol{w} and set to 0 $\frac{\partial L_s}{\partial \boldsymbol{w}} = \frac{2}{m} \sum_{i=1}^m (\langle \boldsymbol{w}, \boldsymbol{x}_i \rangle - y_i) \boldsymbol{x}_i = 0 \rightarrow \sum_{i=1}^m \langle \boldsymbol{w}, \boldsymbol{x}_i \rangle \boldsymbol{x}_i = \sum_{i=1}^m y_i \boldsymbol{x}_i$
- Set

 $A = \left(\sum_{i=1}^{m} \mathbf{x}_{i} \mathbf{x}_{i}^{T}\right) = \begin{bmatrix} \vdots \\ \mathbf{x}_{1} \\ \vdots \end{bmatrix} \dots \begin{bmatrix} \vdots \\ \mathbf{x}_{m} \\ \vdots \end{bmatrix} \begin{bmatrix} \cdots & \mathbf{x}_{1} & \cdots \\ \vdots \\ \cdots & \mathbf{x}_{m} & \cdots \end{bmatrix} \qquad \mathbf{b} = \sum_{i=1}^{m} y_{i} \mathbf{x}_{i} = \begin{bmatrix} \vdots & \cdots & \vdots \\ \mathbf{x}_{1} & \cdots & \mathbf{x}_{m} \\ \vdots & \cdots & \vdots \end{bmatrix} \begin{bmatrix} y_{1} \\ \vdots \\ y_{m} \end{bmatrix}$

• The solution is:

$$\sum_{i=1}^{m} \langle \boldsymbol{w}, \boldsymbol{x}_i \rangle \boldsymbol{x}_i = \sum_{i=1}^{m} y_i \boldsymbol{x}_i \to A \boldsymbol{w} = \boldsymbol{b} \to \boldsymbol{w} = A^{-1} \boldsymbol{b}$$

- The unknown is *w*, A: dxd matrix, b and w: d-dimensional vectors
- The case in which A is not invertible requires a special handling (*not part of the course*)

Polynomial Regression

- Polynomial regression: find the one dimensional polynomial of degree n that better predicts the data
 - $\circ \quad p(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_n x^n$
 - Need to estimate the coefficient vector *a*
 - 1D polynomial pred. deg. $n: \mathcal{H}_{poly}^n = \{x \to p(x)\}, \mathcal{X} = \mathbb{R}, \mathcal{Y} = \mathbb{R}$
- Reduce the problem to a *n*-dimensional linear regression using the mapping:

 $\psi \colon \mathbb{R} \to \mathbb{R}^{n+1} \quad \psi(x) = (1, x, x^2, \dots, x^n)$

We obtain:

 $< a, \psi(x) >= a_0 + a_1 x + a_2 x^2 + \dots + a_n x^n$

- 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



In pseudo-3D space





Logistic Regression

Probability of passing exam versus hours of studying

- 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 $\mathcal{H}: \phi_{sig} \circ L_d$ where $\phi: \mathbb{R} \to [0,1]$ is the sigmoid function and L_d a linear function



Sigmoid Function



Bigger than ½ for positive values and smaller for negative ones

- Tends to 1 for large positive values and to -1 for negative ones
- Can be viewed as a scaled and shifted "soft" sign function

Loss for Logistic Regression

□ Instead of hard choice \rightarrow use probability of correct label being 0 or 1

□
$$H_{sig} = \phi_{sig} \circ L_d = \{ \mathbf{x} \to \phi_{sig} (\langle \mathbf{w}, \mathbf{x} \rangle) : \mathbf{w} \in \mathbb{R}^d \}$$

□ Hypothesis class : $h_{\mathbf{w}}(\mathbf{x}) = \frac{1}{1 + e^{-\langle \mathbf{w}, \mathbf{x} \rangle}}$

□ Loss function: $\ell(h_w, (x, y)) = \log(1 + e^{-y < w, x >})$

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

Logistic Loss Function

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

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

Case $y=1$: need $h_w(x) \rightarrow 1$
 $h_w(x) = \frac{1}{1+e^{-\langle w,x \rangle}} = \frac{1}{1+e^{-y\langle w,x \rangle}}$

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

Case $y=-1$: need $h_w(x) \rightarrow 0 \Rightarrow 1 - h_w(x) \rightarrow 1$
 $1 - h_w(x) = 1 - \frac{1}{1+e^{-\langle w,x \rangle}} = \frac{1+e^{-\langle w,x \rangle}-1}{1+e^{-\langle w,x \rangle}} = \frac{1}{e^{\langle w,x \rangle}+1} = \frac{1}{1+e^{-y\langle w,x \rangle}}$
Same as before :

▶ If denominator small $1 - h_w(x) \rightarrow 1$ good case

▶ if denominator large $1 - h_w(x) \rightarrow 0$ error

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

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

Not part of the course

MLE approach:

- 1. Given training set $S = ((x_1, y_1), ..., (x_m, y_m))$, assume each (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: $L(S; \theta) = \log(P[S|\theta])$
 - \circ log: monotonic \rightarrow same maximum, but simpler to differentiate
- 4. Maximum Likelihood Estimator (MLE): $\hat{\theta} = argmax_{\theta}L(S;\theta)$

MLE and Logistic Regression

Not part of the course

MLE solution is equivalent to ERM solution for logistic regression

Logistic Regression:

- 1. Assume training set $S = ((x_1, y_1), \dots, (x_m, y_m))$
- 2. $P[y_i = 1] = h_w(x_i) = \frac{1}{1 + e^{-\langle w, x_i \rangle}} = \frac{1}{1 + e^{-y_i \langle w, x_i \rangle}}$ (since $y_i = 1$)
- 3. $P[y_i = -1] = 1 h_w(x_i) = \frac{1}{1 + e^{\langle w, x_i \rangle}} = \frac{1}{1 + e^{-y_i \langle w, x_i \rangle}}$ (first equality recall logistic loss, 2nd since $y_i = -1$)
- 4. Likelihood of training set (joint probability P[S|w]): $\prod_{i=1}^{m} \left(\frac{1}{1+e^{-y_i(w,x_i)}}\right)$
- 5. Log likelihood : $\log(P[S|w]) = \log \prod_{i=1}^{m} \left(\frac{1}{1+e^{-y_i \langle w, x_i \rangle}} \right) = -\sum_{i=1}^{m} \log(1+e^{-y_i \langle w, x_i \rangle})$ $\rightarrow corresponds to (-1)^* logistic loss$
- **Maximum Likelihood Estimator:**

$$argmax_{w}L(S; w) = argmax_{w} \log(P[S|w]) = argmin_{w} \sum_{i=1}^{m} log(1 + e^{-y_{i}\langle w, x_{i} \rangle})$$

Recall: argmax(-x)=argmin(x) They have the same target !!!



Example: MLE for Gaussian PDF (1)



Assume that the data is produced by a Gaussian distribution

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



f1 ~ N (10, 2.25) f2 ~ N (10, 9), f3 ~ N (10, 0.25) f4 ~ N (8, 2.25)

find μ, σ maximizing the joint probability of data

Example: MLE for Gaussian PDF (2)



