

ELEC2870 - Machine learning: regression and dimensionality reduction

Linear regression

Michel Verleysen

Machine Learning Group

Université catholique de Louvain

Louvain-la-Neuve, Belgium

michel.verleysen@uclouvain.be

Outline

- Linear regression model
 - Pseudo-inverse
 - Gradient descent
 - Stochastic gradient descent
 - About the sum-of-squares criterion
- Perceptron

Linear regression

- Probably the most elementary way to perform regression
- It is a *linear* model (cannot capture nonlinear relations)

notation 1

$$y = \mathbf{w}^T \mathbf{x} + w_0$$

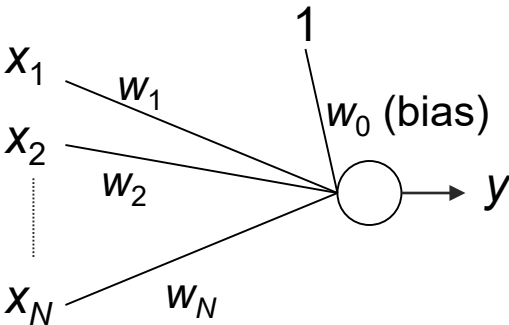


notation 2

$$y = \mathbf{w}^T \mathbf{x}$$

$$\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix} \quad \mathbf{w} = \begin{pmatrix} w_1 \\ \vdots \\ w_N \end{pmatrix}$$

$$\mathbf{x} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_N \end{pmatrix} \quad \mathbf{w} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_N \end{pmatrix}$$



Often, N is replaced by D , as it represents the Dimension of the input space

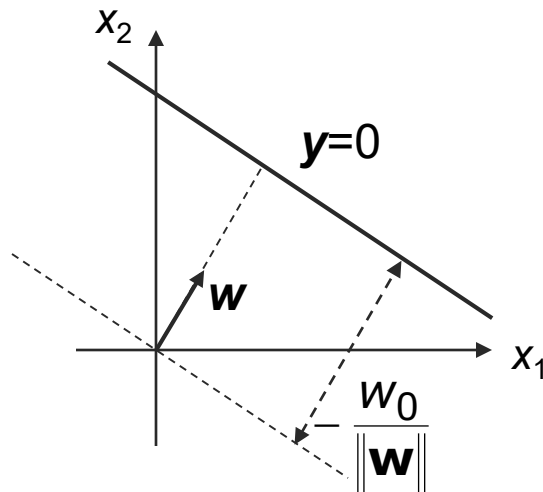
Linear discriminant function

- Notation 1

$$y = \mathbf{w}^T \mathbf{x} + w_0$$

$$\mathbf{x} = (x_1, \dots, x_N)$$

$$\mathbf{w} = (w_1, \dots, w_N)$$

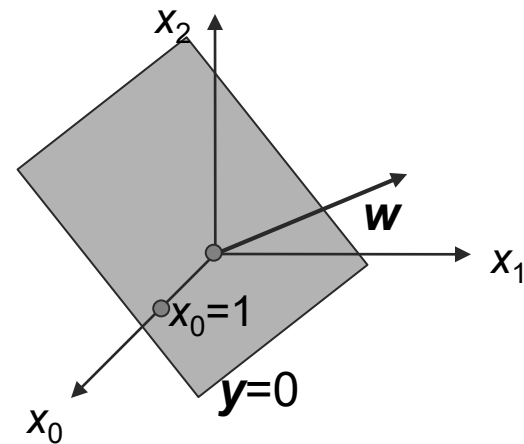


- Notation 2

$$y = \mathbf{w}^T \mathbf{x}$$

$$\mathbf{x} = (1, x_1, \dots, x_N)$$

$$\mathbf{w} = (w_0, w_1, \dots, w_N)$$



Linear regression: criterion

- Model: one linear output
- patterns (learning vectors) must follow:

$$t^p = \sum_{i=1}^D w_i x_i^p = \mathbf{w}^T \mathbf{x}^p$$

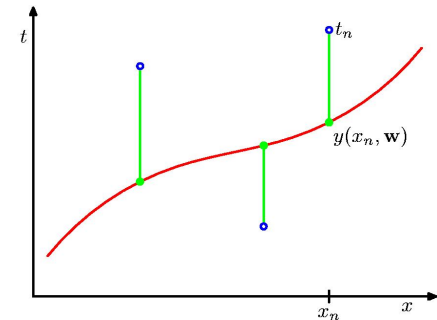
- but $P > D$
 - P patterns
 - D parameters (degrees of freedom)
- → non-ideal solution
- → optimisation (of parameters \mathbf{w}) according to a criterion:

$$E = \frac{1}{P} \sum_{p=1}^P (t^p - y^p)^2 = \frac{1}{P} \sum_{p=1}^P (t^p - \mathbf{w}^T \mathbf{x}^p)^2$$

About the sum-of-squares criterion

$$E = \frac{1}{P} \sum_{p=1}^P (t^p - y^p)^2 = \frac{1}{P} \sum_{p=1}^P (t^p - \mathbf{w}^T \mathbf{x}^p)^2$$

- The sum-of-squares criterion
 - is convenient (its derivative is linear)
 - makes the sign of the error irrelevant.
- But
 - it is not natural (the error is an interpretable distance, the square is not)
 - it gives a very large weight to large errors (the square of a large number is very large...)
 - a single or a few outlier(s) may then influence a lot the criterion, and consequently the model!
- Don't hesitate to reconsider the criterion in real settings, even at the price of a more complex model!



From: C. Bishop, Pattern Recognition and Machine Learning, Springer, 2006.

Optimizing the criterion by pseudo-inverse

- Error criterion $E = \frac{1}{P} \sum_{p=1}^P (t^p - y^p)^2 = \frac{1}{P} \sum_{p=1}^P (t^p - \mathbf{w}^T \mathbf{x}^p)^2$
- Inputs \mathbf{x}^p in a matrix and outputs t^p in a vector

$$\mathbf{X} = (\mathbf{x}^1 \mathbf{x}^2 \dots \mathbf{x}^P) = \begin{pmatrix} x_1^1 & x_1^2 & \dots & x_1^P \\ x_2^1 & x_2^2 & \dots & x_2^P \\ \vdots & \vdots & & \vdots \\ x_D^1 & x_D^2 & \dots & x_D^P \end{pmatrix}$$

$$\mathbf{t}^T = (t^1 \ t^2 \ \dots \ t^P)$$

- Error criterion

$$E = \frac{1}{P} \|\mathbf{t}^T - \mathbf{w}^T \mathbf{X}\|^2$$

Warning: sometimes (even in these lectures notes...), the definition of data matrix is the transpose of this one (columns are rows); all subsequent formulas are then transposed too. Exemple:

$$E = \frac{1}{P} \|\mathbf{X}\mathbf{w} - \mathbf{t}\|^2$$

Optimizing the criterion by pseudo-inverse

- Error criterion $E = \frac{1}{P} \|\mathbf{t}^\top - \mathbf{w}^\top \mathbf{X}\|^2$
- Gradient of error (function to minimize)
with respect to weights (free parameters)

$$\left(\frac{\partial E}{\partial \mathbf{w}} \right)^T \equiv \left(\frac{\partial E}{\partial w_1} \frac{\partial E}{\partial w_2} \dots \frac{\partial E}{\partial w_D} \right)$$

$$\begin{aligned} \frac{\partial E}{\partial w_j} &= \frac{\partial}{\partial w_j} \left(\frac{1}{P} \|\mathbf{t}^\top - \mathbf{w}^\top \mathbf{X}\|^2 \right) \\ &= \frac{\partial}{\partial w_j} \left(\frac{1}{P} (\mathbf{t}^\top - \mathbf{w}^\top \mathbf{X})(\mathbf{t} - \mathbf{X}^\top \mathbf{w}) \right) \\ &= \frac{2}{P} (\mathbf{w}^\top \mathbf{X} - \mathbf{t}^\top) \mathbf{x}_j \quad \text{where } \mathbf{x}_j = (x_j^1 \ x_j^2 \ \dots \ x_j^P) \end{aligned}$$

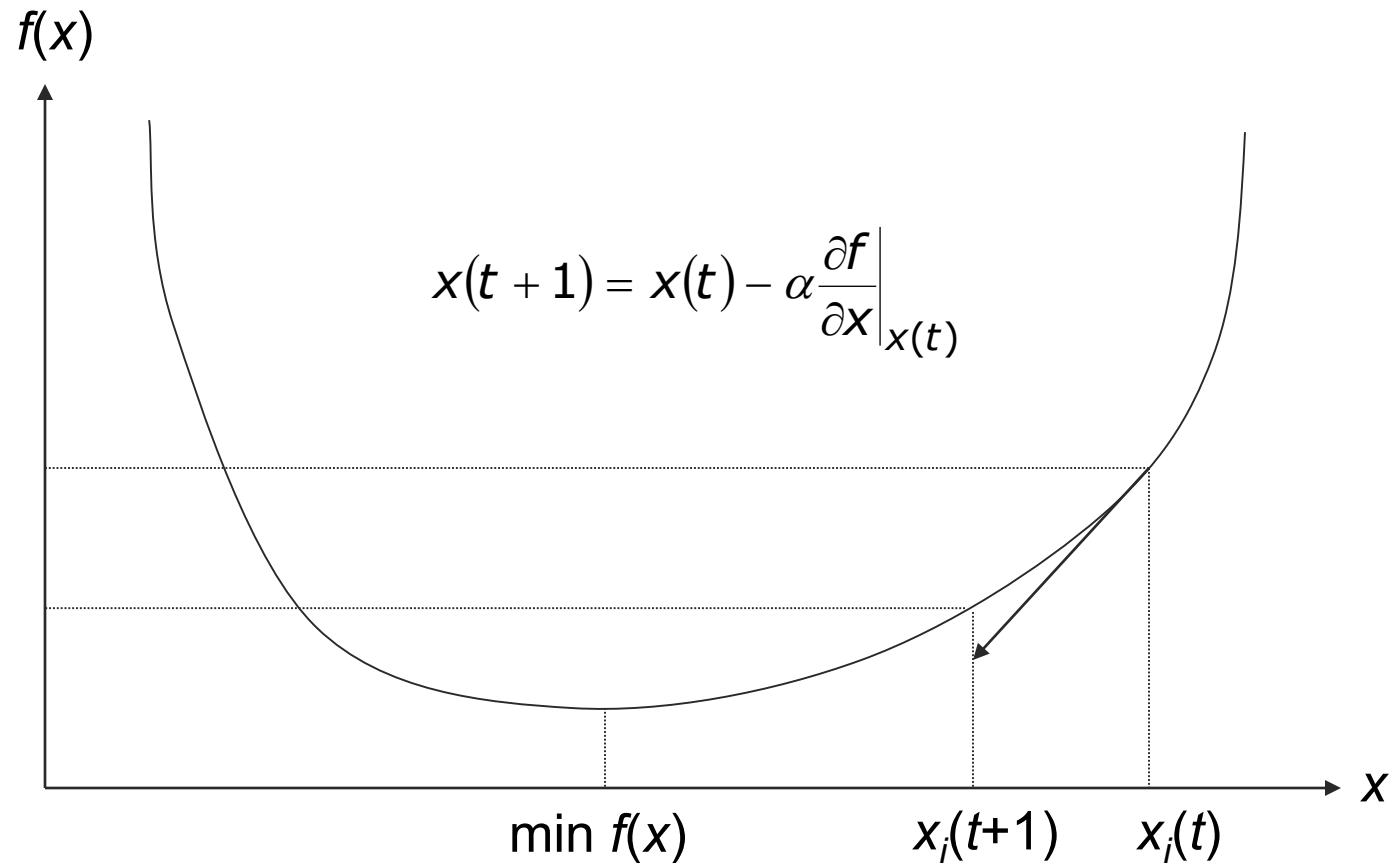
Optimizing the criterion by pseudo-inverse

- Criterion $E = \frac{1}{P} \| \mathbf{t}^\top - \mathbf{w}^\top \mathbf{X} \|^2$
- Derivative of criterion $\left(\frac{\partial E}{\partial \mathbf{w}} \right)^\top = \frac{2}{P} (\mathbf{w}^\top \mathbf{X} - \mathbf{t}^\top) \mathbf{X}^\top$
- Minimum of error $\left(\frac{\partial E}{\partial \mathbf{w}} \right)^\top = 0$

$$\longrightarrow \boxed{\mathbf{w} = (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{X}\mathbf{t}}$$

- Pseudo-inverse requires :
 - all input-output pairs (\mathbf{x}^p, t^p)
 - matrix inversion (often ill-configured)
- necessity for iterative methods without matrix inversion
- → Gradient descent !

A reminder on gradient descent



Gradient descent: elementary example

- minimum of $f(x)$

$$\begin{array}{rcl}
 x(t+1) = x(t) - \alpha \frac{\partial f}{\partial x} \Big|_{x(t)} & & f(x) = (x+2)^2 - 1 \\
 & & \frac{\partial f}{\partial x} = 2(x+2) \\
 & \swarrow & \\
 x(t+1) = x(t) - 2\alpha(x(t)+2) & & \\
 = (1-2\alpha)x(t) - 4\alpha & & \\
 = -2 + (1-2\alpha)(x(t)+2) & & \\
 & \swarrow & \\
 y(t+1) = (1-2\alpha)y(t) & & y(t) \equiv x(t)+2
 \end{array}$$

- converges to $y(t)=0$ (or $x(t)=-2$) if $0 < \alpha < 1$

Optimizing the criterion by gradient descent

- function to minimise: E
- parameters: \mathbf{w}

$$\mathbf{w}(t+1) = \mathbf{w}(t) - \alpha \left. \frac{\partial E}{\partial \mathbf{w}} \right|_{\mathbf{w}(t)}$$

$$\boxed{\mathbf{w}(t+1) = \mathbf{w}(t) + \frac{2}{P} \alpha \mathbf{X} \left(\mathbf{t} - \mathbf{X}^T \mathbf{w}(t) \right)}$$

- pseudo-inverse and gradient descent:
same error criterion -> same solution !
- Pros and cons
 - needs iterations
 - but does not require matrix inversion
 - still needs all input-output pairs (\mathbf{x}^p, t^p)

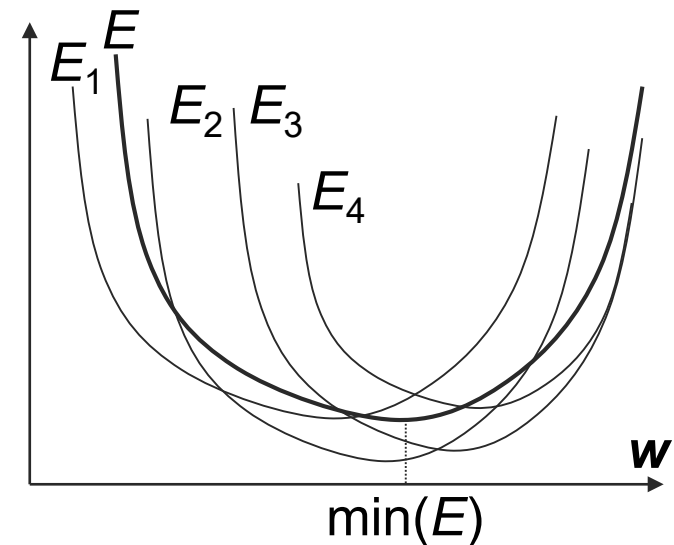
Optimizing the criterion by stochastic gradient descent

$$E = \frac{1}{P} \sum_{p=1}^P \left(t^p - \mathbf{w}^\top \mathbf{x}^p \right)^2 = \frac{1}{P} \sum_{p=1}^P E_p$$

- If data are stationery :

minimising E (or $P E$) is
equivalent to successively
minimising each E_k

$$\mathbf{w}(t+1) = \mathbf{w}(t) + 2\alpha \left(t^k - \mathbf{w}(t)^\top \mathbf{x}^k \right) \mathbf{x}^k$$



Difference between p , k and t : p and k are indices on the patterns (1... P), while t identifies iterations (may exceed P). p is the indice in the database of patterns, while k identifies the order of presentation, which may differ. The difference between p and k is not crucial here!

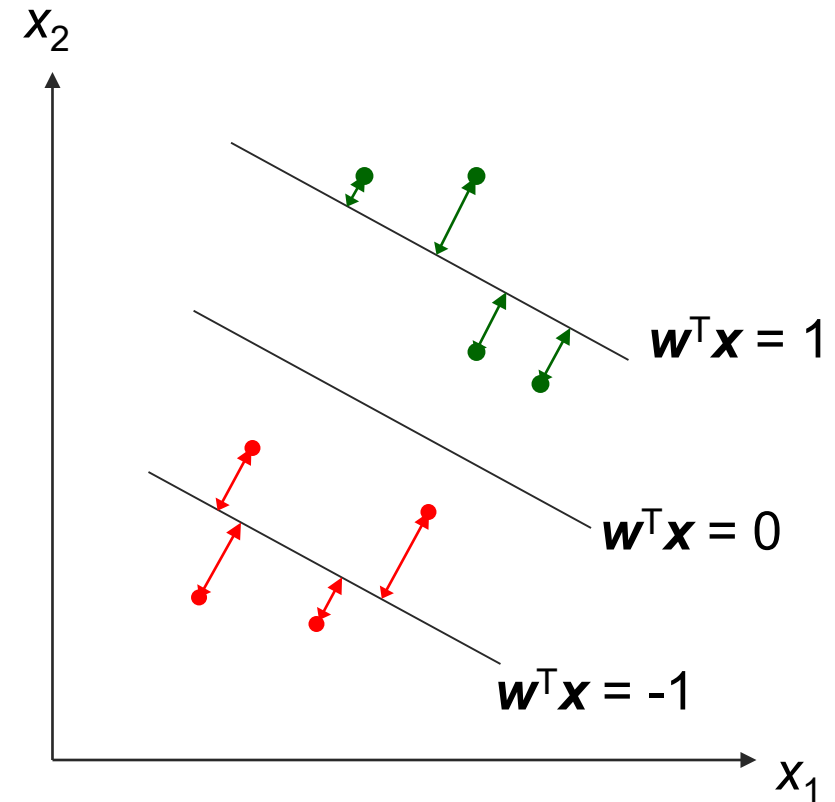
Optimizing the criterion : comparison

	Pseudo-inverse	Gradient descent	Stochastic gradient descent
Needs all (\mathbf{x}^p, t^p) at each iteration	Yes	Yes	No
# of iterations	1	Several	Many
Matrix inversion	Yes	No	No
Sensitive to order of patterns	No	No	Might be

- Note: in the (stochastic or not) gradient descent versions, the linear regression model is also called Adaline (Adaptive Linear Element)

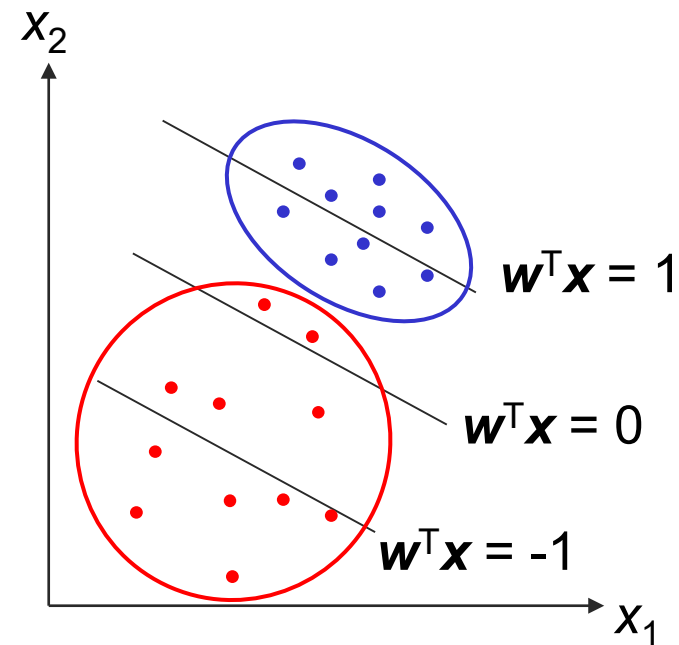
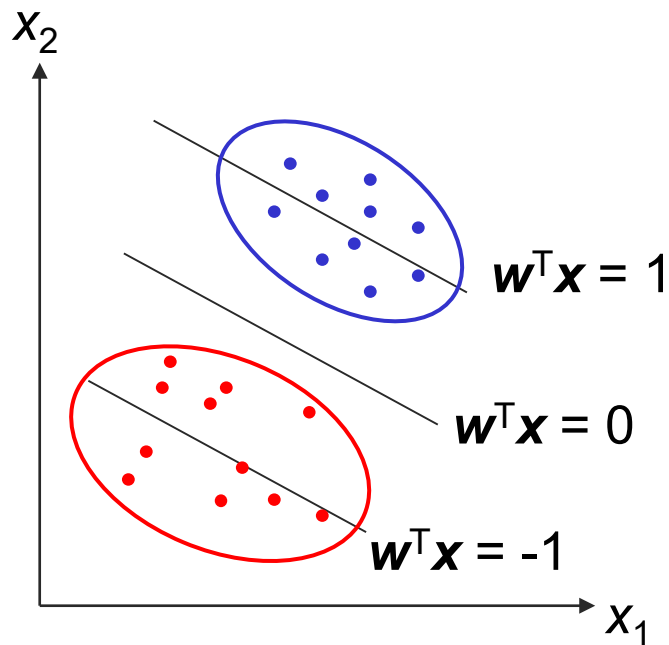
Classification with a linear regression model

- parameters \mathbf{w} are adjusted with respect to $\mathbf{w}^T \mathbf{x} = \pm 1$ (through the sum-of-squares criterion)
- separation $\mathbf{w}^T \mathbf{x} = 0$ is a consequence
- any separation $\mathbf{w}^T \mathbf{x} = A$ could be chosen



Sum-of-squares criterion in classification

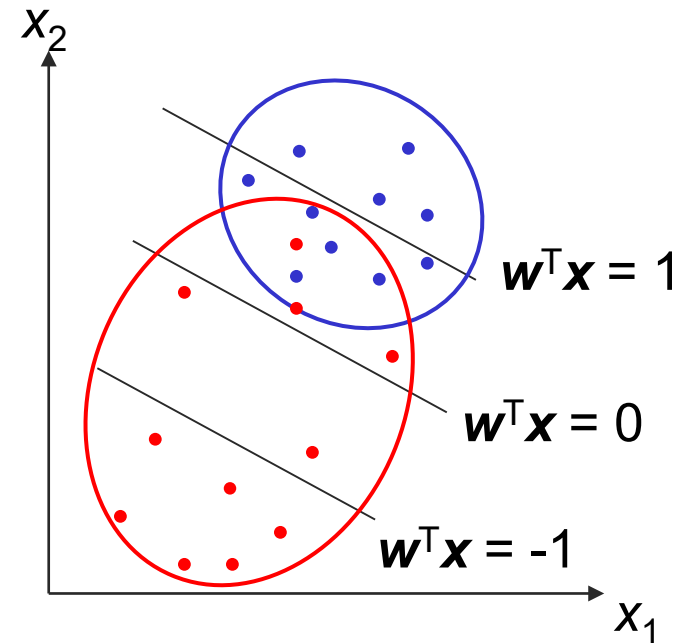
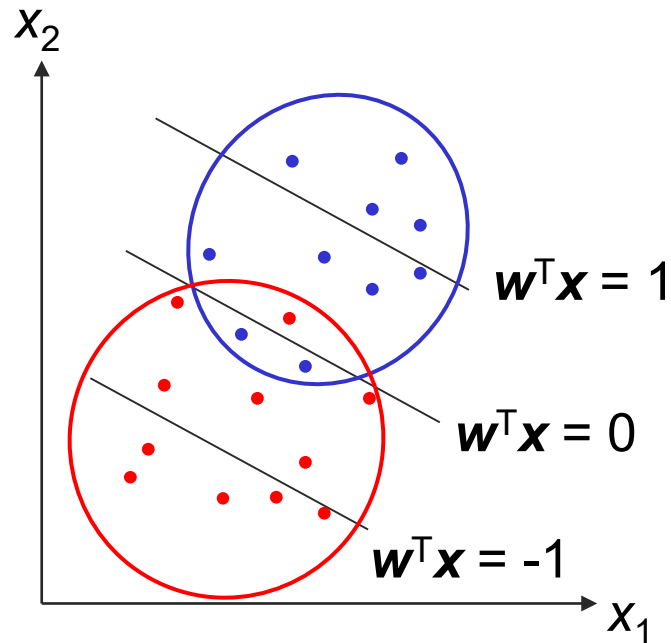
- E (sum-of-squares) is *not* equivalent to a minimum # of misclassifications



- Therefore it is *not* a good criterion for classification tasks, but still, it is widely used for its convenience!

Sum-of-squares criterion in classification

- When classes are *not* linearly separable



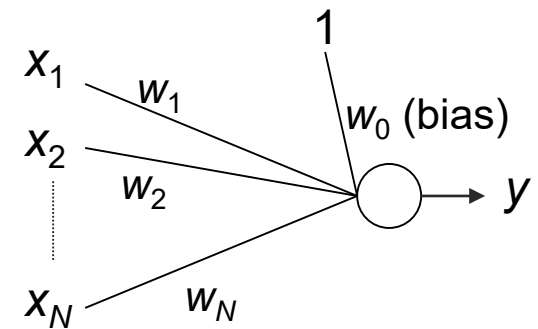
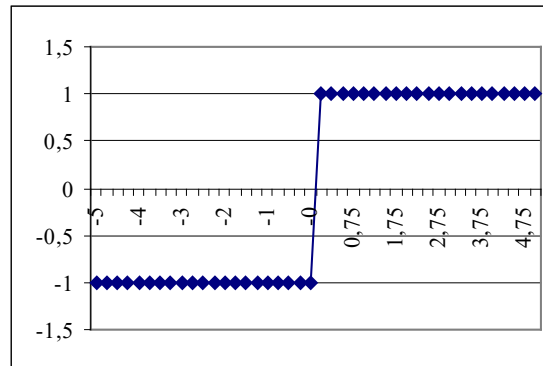
Outline

- Linear regression model
 - Pseudo-inverse
 - Gradient descent
 - Stochastic gradient descent
 - About the sum-of-squares criterion
- Perceptron

Perceptron

- The perceptron is a *classification* model
- It is introduced here as an example to emphasize the differences with respect to regression problems

- single output model with threshold (sign) as non-linear activation function
- outputs $\in \{+1, -1\}$



- error criterions:
 - Least Mean Square: cannot be used (non-continuous)
 - # of misclassifications: non-continuous too
 - \rightarrow use of perceptron criterion

Perceptron criterion

- perceptron outputs $\in \{+1, -1\}$
- class labels $t^p \in \{+1, -1\}$ for classes C^1 and C^2 respectively
- in case of correct classification

$$t = \text{sign}(\mathbf{w}^T \mathbf{x}) \begin{cases} \text{Class } C^1 & \mathbf{w}^T \mathbf{x}^k > 0 \\ \text{Class } C^2 & \mathbf{w}^T \mathbf{x}^k < 0 \end{cases} \Rightarrow \mathbf{w}^T (\mathbf{x}^k t^k) > 0$$

- An ideal criterion could be $E = - \sum_{\mathbf{w}^T \mathbf{x}^k t^k < 0} 1$
 - but this criterion is not continuous (an ε change in \mathbf{w} results in a 1 increment in E)

- perceptron criterion $E = - \sum_{\mathbf{w}^T \mathbf{x}^k t^k < 0} (\mathbf{w}^T \mathbf{x}^k t^k)$
 - continuous
 - gradient: non-continuous, piece-wise linear



Perceptron learning rule

- stochastic gradient descent on perceptron learning rule:
If \mathbf{x}^k is misclassified (and only in this case):

$$\mathbf{w}(t + 1) = \mathbf{w}(t) + \alpha \mathbf{x}^k t^k$$

- error decreases at each step

$$\begin{aligned} -\mathbf{w}(t + 1)^T \mathbf{x}^k t^k &= -\mathbf{w}(t)^T \mathbf{x}^k t^k - \alpha (\mathbf{x}^k t^k)^T \mathbf{x}^k t^k \\ &< -\mathbf{w}(t)^T \mathbf{x}^k t^k \end{aligned}$$

- perceptron convergence theorem
 - for any data set linearly separable, there is convergence to a solution in a finite number of steps

Perceptron convergence theorem

- Many proofs available
 - first one by Rosenblatt (1962)
 - here according to Bishop
- Classes are linearly separable (hypothesis) -> there exists \mathbf{w}_{sol} such that

$$\mathbf{w}_{\text{sol}}^T \mathbf{x}^p t^p > 0 \quad \forall p$$

- Hypotheses (without loss of generality):
 - $\mathbf{w}(0) = 0$
 - $\alpha = 1$

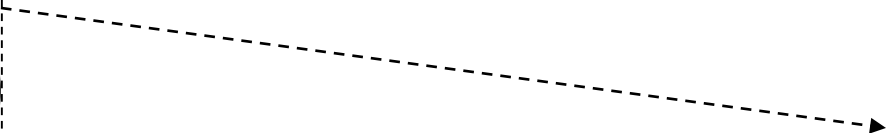
Perceptron convergence theorem

- At each step k $\mathbf{w}(t+1) = \mathbf{w}(t) + \mathbf{x}^k t^k$

- After n iterations $\mathbf{w} = \sum_p n^p \mathbf{x}^p t^p$

where n^p is the # of presentations of pattern p

- Then

$$\begin{aligned} \mathbf{w}_{\text{sol}}^T \mathbf{w} &= \sum_p n^p \mathbf{w}_{\text{sol}}^T \mathbf{x}^p t^p \\ &\geq n \min_p (\mathbf{w}_{\text{sol}}^T \mathbf{x}^p t^p) \end{aligned}$$


Perceptron convergence theorem

- Other inequality

$$\begin{aligned}\|\mathbf{w}(t+1)\|^2 &= \|\mathbf{w}(t)\|^2 + \|\mathbf{x}^p\|^2 t^{p2} + 2\mathbf{w}(t)^\top \mathbf{x}^p t^p \\ &< \|\mathbf{w}(t)\|^2 + \|\mathbf{x}^p\|^2 t^{p2}\end{aligned}$$

$$\Delta \|\mathbf{w}\|^2 = \|\mathbf{w}(t+1)\|^2 - \|\mathbf{w}(t)\|^2 \leq \max_p \|\mathbf{x}^p\|^2$$

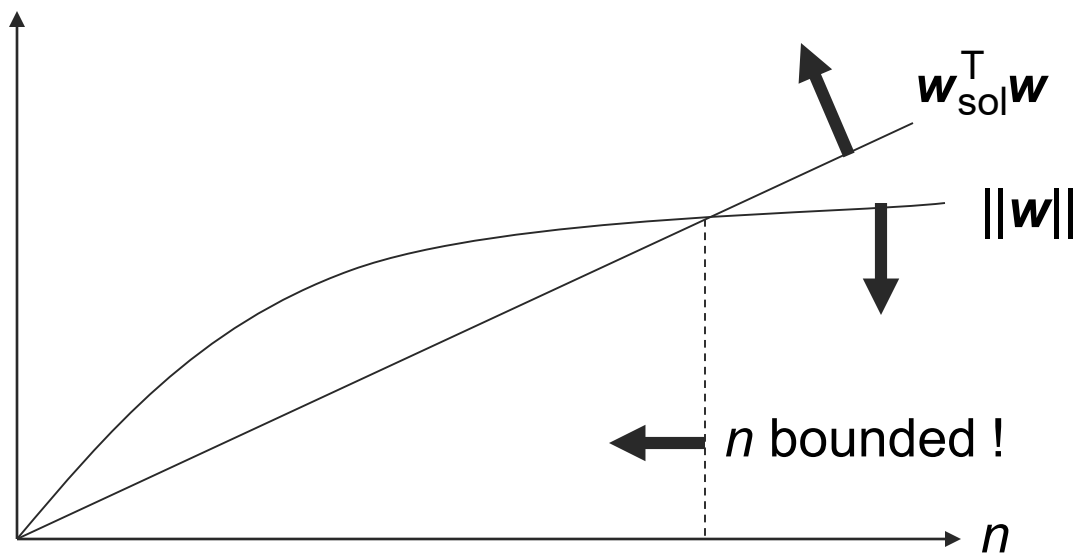
$$\|\mathbf{w}\|^2 \leq n \max_p \|\mathbf{x}^p\|^2 \longrightarrow \boxed{\|\mathbf{w}\| \leq \sqrt{n} \max_p \|\mathbf{x}^p\|}$$

- Thus

- $\|\mathbf{w}\|$ increases no faster than $n^{1/2}$
- $\mathbf{w}_{\text{sol}}^\top \mathbf{w}$ bounded below by linear function of n

Perceptron convergence theorem

- $\|\mathbf{w}\|$ increases no faster than $n^{1/2}$
- $\mathbf{w}_{\text{sol}}^T \mathbf{w}$ bounded below by linear function of n



Limitations to Perceptrons

- “Perceptrons can only solve linearly separable problems” (example: XOR not possible)
- Minsky & Papert book put a temporary end to the research in the field during many years!
- Not exactly true
 - inputs can be preprocessed
 - the problem is the model (linear), add nonlinearities!