Outline



- Properties of the perceptron
- Input-output pairs
- Perceptron learning method
- Perceptron learning example
- Proof of convergence
- Good material:

http://hagan.okstate.edu/4 Perceptron.pdf

The Perceptron

- Receives input through its synapsis (x_i)
- Synapsis are weighted (*w_i*)
- A b value biases the sum to enable asymmetric behavior
- A weighted sum is calculated
- Activation function applied





x_i : input vector

 w_{ki} : weight coefficient vector of neuron k b_k : bias value of neuron k o_k : output value of neuron k

9/17/2019.

P-ITEEA-0011

Lecture 2

Neural Networks



Perceptron is an Input \rightarrow Output device



As opposed to **Traditional Computers** where

- the math of the functionality is known
- the known math should be programmed

At Neural Networks

the math behind the functionality is unknown the functionality is "illustrated" with examples

Function illustrated by examples

• Given a set of input-output pairs

 $x_j \rightarrow d_j$ (x_j: input vector; d_j : desired output)

- Number of input vectors
 - Finite/limited set (e.g. AND function)
 - Equivalent with a look-up-table (LUT), math known
 - Mathematically it is correct to define a function by listing all the IO pairs
 - Goal: generate a simpler than LUT decision making device through learning
 - Infinite/open set (customers of a bank asking for a loan)
 - Math behind is unknown, cannot be coded directly
 - Goal: generate the function through learning
 - It should predict well the output of a previously unknown/untested input (<u>GENERALIZATION</u>)

	X				d
#	age	gender	debt	salary	
1	25	M(1)	25	100	Y(1)
2	22	F(2)	18	80	Y(1)
3	65	M(1)	3000	200	N(0)
					•





Linear separability



- Today, we assume that the IO sets are linearly separable
- The decision boundary is a hyperplane defined:

$$\mathbf{w}^T \mathbf{x} = \mathbf{0}$$

- Positive side of the hyperplane is classified: +1 (yes)
- Negative side of the hyperplane is classified : 0 (no).













X2

Lecture 2

P-ITFFA-0011

X



Maximum Margin:

Define the margin

of a linear classifier as the width that the boundary could be increased by before hitting a data point.





Maximum Margin:

Define the margin

of a linear classifier as the width that the boundary could be increased by before hitting a data point.

What does learning mean?

- Given an annotated dataset $\mathbf{x}_i \rightarrow \mathbf{d}_i$
- Given the parametric equation of the perceptron
 - $y = sign(\mathbf{w}^T \mathbf{x})$
- Goal: find the optimal \mathbf{w}_{opt} weights (parameters), where for each j $d_i = sign(\mathbf{w}_{opt}^T \mathbf{x}_i)$





The learning algorithm: Datasets

- Training set
 - Set of input desired output pairs
 - Will be used for training

 $X^+ = \{ \mathbf{x} : d = +1 \}$ $X^- = \{ \mathbf{x} : d = 0 \}$

- Test set
 - Used, when we have large set of input vectors (not used today)
 - Set of input desired output pairs
 - Will be used for testing and scoring the result
- We assumed that X⁺ and X⁻ must be linearly separable
- We are looking for an optimal parameter set:

$$X^{+} = \left\{ \mathbf{x} : \mathbf{w}_{\text{opt}}^{\text{T}} \mathbf{x} > 0 \right\},$$
$$X^{-} = \left\{ \mathbf{x} : \mathbf{w}_{\text{opt}}^{\text{T}} \mathbf{x} < 0 \right\}.$$

The learning algorithm: Recursive algorithm



- We have to develop a recursive algorithm called learning, which can learn the weight step by step, based on observing
 - the (i) input,
 - the (ii) weight vector,
 - the (iii) desired output, and
 - the (iv) actual output of the system.
- This can be described formally as follows:

$$\mathbf{w}(k+1) = \Psi(\mathbf{x}(k), \mathbf{w}(k), d(k), y(k)) \rightarrow \mathbf{w}_{opt}$$

The learning algorithm: Perceptron Learning Algorithm

- In a more ambitious way it can be called intelligent, because
 - perceptron can learn through examples (adapt),
 - even the function parameters are fully hidden.
- Perceptron learning was introduced by Frank Rosenblatt 1958
 - Built a 20x20 image sensor
 - With analog perceptron
 - 400 weights controlled by electromotors



The learning algorithm: Recursive steps

- 1. Initialization.
 - Set **w**(0)=**0** or **w**(0)=**rand**
- 2. Activation.

Select a $\mathbf{x}_k \rightarrow \mathbf{d}_k$ pair

- 3. Computation of actual response $y(k) = sign(w^T(k)x(k))$
- 4. Adaptation of the weight vector $\mathbf{w}(k+1) = \Psi(\mathbf{x}(k), \mathbf{w}(k), d(k), y(k))$
- 5. Continuation

Until all responses of the perceptron are OK

Lecture 2

P-ITFFA-0011



- Given a 3 input vector example
- Assume that bias is zero • (decision boundary will cross the origo)
- Random initialization

to the decision boundary!!!

Decision boundary: $x_1 - 0.8x_2 = 0$

Its orthogonal vector is: (1, -0.8) **P-ITEEA-0011** Lecture 2

Remember: the weight vector is orthogonal

$$\mathbf{x}_{1} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, d_{1} = 1;$$

$$\mathbf{x}_{2} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}, d_{2} = 0;$$

$$\mathbf{x}_{3} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, d_{3} = 0;$$

$$\mathbf{x}_{3} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, d_{3} = 0;$$

 $\mathbf{w}^{T}(1) = \begin{bmatrix} 1 & -0.8 \end{bmatrix};$

$$2 \circ 1$$

3



17





Weight update: very simple example

• Test with the first input vector

$$\mathbf{x}_{1} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, d_{1} = 1; \qquad d_{j} \cdot y_{j} > 0$$
$$\mathbf{w}^{T}(1) = \begin{bmatrix} 1 & -0.8 \end{bmatrix}; \qquad \mathbf{w}^{T}(1) = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \mathbf{w}^{T}(1$$

$$y_1(1) = sign(w^T(1)x_1) = sign(\begin{bmatrix} 1 & -0.8 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix}) = sign(1 - 1.6) = 0$$

The result is not OK! Positive misclassification: Instead of 1, the result is 0!! (The normal vector points to the positive side of the decision boundary.)



Idea: <u>add</u> the vector pointing to the positively misclassified point to the orthogonal vector of the decision boundary, to rotate it <u>towards</u> the point! $w(k+1)=w(k)+x_1$



 $\mathbf{w}^{T}(2) = [1 + 1 -0.8 + 2] = [2 1.2];$



Weight update: very simple example

• Test with the second input vector $\mathbf{x}_2 = \begin{bmatrix} -1 \\ 2 \end{bmatrix}, d_1 = 0;$ $\mathbf{w}^T(2) = \begin{bmatrix} 2 & 1.2 \end{bmatrix};$ $y_2(2) = sign(w^T(2)x_2) = sign(\begin{bmatrix} 2 & 1.2 \end{bmatrix}, \begin{bmatrix} -1 \\ 2 \end{bmatrix}) = sign(-2 + 2.4) = 1$

The result is not OK! Negative misclassification: Instead of 0, the result is 1!!



Idea: **<u>subtract</u>** the vector pointing to the negatively misclassified point to the orthogonal vector of the decision boundary, to rotate it <u>away</u> the point! $w(k+2)=w(k+1)-x_2$

$$\mathbf{w}^{T}(3) = [2 - (-1) \quad 1.2 - 2] = [3 \quad -0.8]$$



Weight update: very simple example

Test with the third input vector

Test with the third input vector
$$\mathbf{x}_3 = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, d_3 = 0;$$

 $\mathbf{w}^T(3) = \begin{bmatrix} 3 & -0.8 \end{bmatrix};$
 $y_3(3) = sign(w^T(3)x_3) = sign(\begin{bmatrix} 3 & -0.8 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix}) = sign(0 + 0.8) = 1$

The result is not OK! Negative misclassification: Instead of 0, the result is 1!!



Again: **subtract** the vector pointing to the negatively misclassified point to the orthogonal vector of the decision boundary, to rotate it away the point! $w(k+3)=w(k+2)-x_3$



$$\mathbf{w}^{T}(4) = \begin{bmatrix} 3 - 0 & -0.8 - (-1) \end{bmatrix} = \begin{bmatrix} 3 & 0.2 \end{bmatrix};$$

ality!!! 2 C

Test with the again with the second vector

The result is OK!

- Do not modify!!!
- Test with the again with the third vector

The result is OK!

- Do not modify!!!
- Since all input vectors are correctly classified: we are ready

Weight update: very simple example

- Start again:
 - Test with the again with the first vector

The result is OK!

Do not modify!!!









- Positive misclassification : ADD $\varepsilon = d_j \cdot y_j = 1$ $w(k+1)=w(k)+x_j$
- Negative misclassification : SUBTRACT $\varepsilon = d_j \cdot y_j = -1$ w(k+1)=w(k)-x_j
- Correct classification : DO NOTHING $\varepsilon = d_j \cdot y_j = 0$ w(k+1)=w(k)
- In general:

w(k+1)=w(k)+
$$\varepsilon$$
 x_j





The learning algorithm: Adaptation

We were looking for a recursive function:

 $\mathbf{w}(k+1) = \Psi(\mathbf{x}(k), \mathbf{w}(k), d(k), y(k))$

In general:

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \varepsilon \eta \mathbf{x}_j$$

where is the error function

$$\varepsilon(k) = d(k) - y(k)$$

$$d(k) = \begin{cases} 0 & \text{if } \mathbf{x}(k) \text{ belongs to class } X^+ \\ -1 & \text{if } \mathbf{x}(k) \text{ belongs to class } X^-, \end{cases}$$

and

 η is the learning rate (η controls the learning speed and should be positive)





Weight update strategy



- Apply all the input vectors in one after the others, selecting them randomly
- Instance update
 - Update the weights after each input
- Batch update
 - Add up the modifications
 - Update the weights with the sum of the modifications, after all the inputs were applied
- Mini batch
 - Select a smaller batch of input vectors, and do with that as in the batch mode

9/17/2019.

Perceptron Convergence theorem (1)

Assumptions:

- **w**(0)=0
- the input space is linearly separable, therefore w_o (stands for $w_{optimal}$) exists:

$$x \in X^+$$
: $w_o^T x > 0: d = 1$
 $x \in X^-$: $w_o^T x < 0: d = 0$

- Let us denote
$$\tilde{x} = -x$$

 $\tilde{x} \in \tilde{X}^-$: $w_{o}^T \tilde{x} > 0$: $d = 1$

For the proof, see also: Simon Haykins: Neural Networks and Learning Machines, Section 1.3: <u>http://dai.fmph.uniba.sk/courses/NN/haykin.neural-networks.3ed.2009.pdf</u>



Perceptron Convergence theorem (2)



- Idea:
 - During the training, the network will be activated with those input vectors (one after the other), where the decision is wrong, hence non zero adaptation is needed:

$$x(j) \in X^+$$
: $w^T(j)x(j) < 0$, $y = 0$, $d = 1$
 $x(j) \in \tilde{X}^-$: $w^T(j)x(j) < 0$, $y = 0$, $d = 1$

Note: The error function is always positive ($\mathcal{E} = 1$)



Perceptron Convergence theorem (3)

- According to the learning method:
- $w(n+1)=w(0)+\eta x(0)+\eta x(1)+\eta x(2)+\eta x(3)+...+\eta x(n)$
 - where

$$x(j) \in X^+$$
: $w^T(j)x(j) < 0, y = -1, d = 1$

or

$$x(j) \in \widetilde{X}^{-}$$
: $w^{T}(j)x(j) < 0, y = -1, d = 1$

- The decision boundary will be:

 $\eta w^T x = 0$

which means that η is a scaling factor, therefore it can be choosen for any positive number.

Let us use $\eta = 1$, therefore $\eta \varepsilon = 1$

Perceptron Convergence theorem (4)



• We will calculate $||w(n+1)||^2$ in two ways, and give an upper and a lower boundary, and it will turn out that an n_{max} exists, and beyond that the lower boundary is higher than the upper boundary (squeeze theorem, sandwitch lemma (*közrefogási elv, rendőr elv*))

Perceptron Convergence theorem (5) lower limit (1)



According to the learning method, the presented input vectors are added up:

 $w(n+1) = w(0) + x(0) + x(1) + \dots + x(n)$ w(0)=0

Multiply it with w_o^T from the left:

$$w_o^T w(n+1) = w_o^T x(0) + w_o^T x(1) + \dots + w_o^T x(n)$$

$$0 < \alpha \le w_o^T x(j)$$
 Because each input vector (or its opposite) were selected that way.

$$0 < \alpha = \min_{x(n) \in \{X^+, \tilde{X}^-\}} w_o^T x(n)$$

$$w_o^T w(n+1) \ge n\alpha$$

Perceptron Convergence theorem (6) lower limit (2)



 $w_o^T w(n+1) \ge n\alpha$ We apply Cauchy Schwarty inequality $||a||^2 ||b||^2 \ge ||a^T b||^2$

$$\|w_0^T\|^2 \|w(n+1)\|^2 \ge \|w_o^T w(n+1)\|^2 \ge n^2 \alpha^2$$

Lower limit:

$$\|w(n+1)\|^2 \ge \frac{n^2 \alpha^2}{\|w_0^T\|^2}$$

Lower limit proportional with n²

Perceptron Convergence theorem (7) upper limit (1)



Let us have a different synthetization approach of w(n+1):

w(k+1) = w(k) + x(k)

for k= 0 ... n

Squared Euclidian norm:

$$\|w(k+1)\|^{2} = \|w(k)\|^{2} + \|x(k)\|^{2} + 2w(k)^{T}x(k)$$

Because each input vector (or its opposite) were selected that way.

$$\|w(k+1)\|^{2} \le \|w(k)\|^{2} + \|x(k)\|^{2}$$
$$\|w(k+1)\|^{2} - \|w(k)\|^{2} \le \|x(k)\|^{2}$$

 $w(k)^T x(k) < 0$

for k= 0 ... n

Perceptron Convergence theorem (8) upper limit (2)



$$\|w(k+1)\|^2 - \|w(k)\|^2 \le \|x(k)\|^2$$

Summing up the upper term:

$$\sum_{k=0}^{n} \left(\left\| w(k+1) \right\|^{2} - \left\| w(k) \right\|^{2} \right) \le \sum_{k=0}^{n} \left\| x(k) \right\|^{2}$$

Telescoping sum: $\sum_{i=1}^{n} (a_{i+1} - a_i) = a_{n+1} - a_1$ Example: $\sum_{i=1}^{4} (a_{i+1} - a_i) = a_2 - a_1 + a_3 - a_2 + a_4 - a_3 + a_5 - a_4 = a_5 - a_1$

Note that there is a telescoping sum in the left hand side.



Linear upper limit and squared lower limit cannot grow unlimitedly

n_{max} should exist

$$n_{\max} = \frac{\beta \|w_0\|^2}{\alpha^2}$$