

Perceptrons

Sven Koenig, USC

Russell and Norvig, 3rd Edition, Sections 18.7.1-18.7.4

These slides are new and can contain mistakes and typos.
Please report them to Sven (skoenig@usc.edu).

1

Perceptrons

- We now study how to acquire knowledge with machine learning.

2

Inductive Learning for Classification

- Labeled examples

Feature_1	Feature_2	Class
true	true	true
true	false	false
false	true	false

Learn $f(\text{Feature}_1, \text{Feature}_2) = \text{Class}$ from

$$f(\text{true}, \text{true}) = \text{true}$$

$$f(\text{true}, \text{false}) = \text{false}$$

$$f(\text{false}, \text{true}) = \text{false}$$

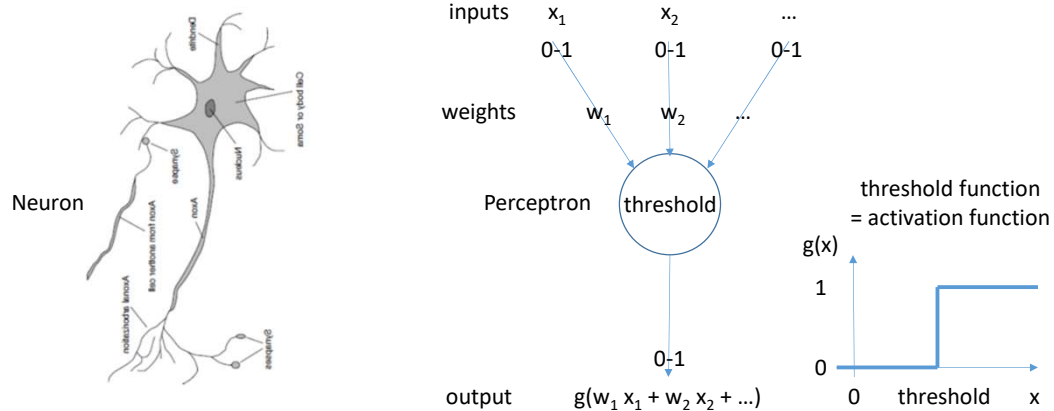
The function needs to be consistent with all labeled examples and should make the fewest mistakes on the unlabeled examples.

- Unlabeled examples

Feature_1	Feature_2	Class
false	false	?

3

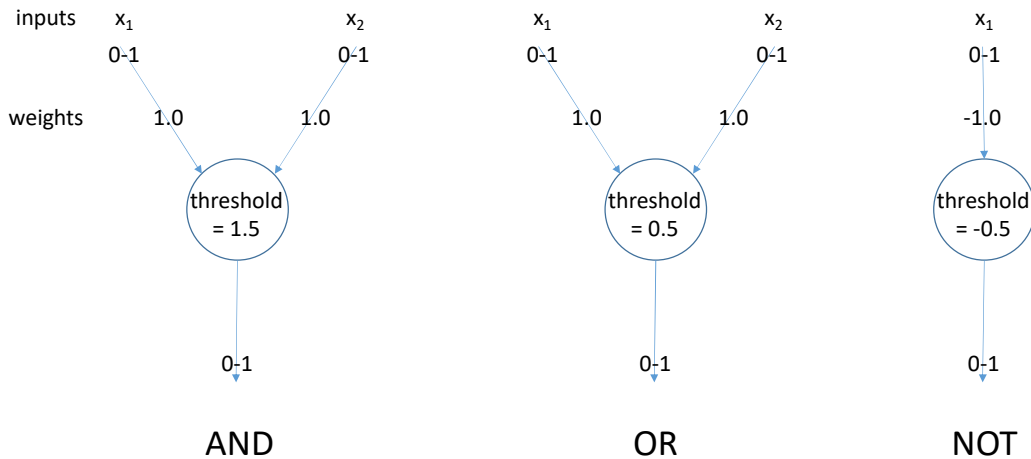
Example: Perceptron Learning



- Objective: Learn the weights for a given perceptron.
- From now on: binary (feature and class) values only (0=false, 1=true).

4

Example: Perceptron Learning

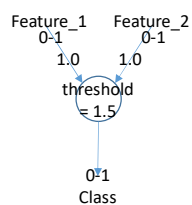


5

Example: Perceptron Learning

- Labeled examples

Feature_1	Feature_2	Class
true	true	true
true	false	false
false	true	false



- Unlabeled examples (note: classification is very fast)

Feature_1	Feature_2	Class
false	false	? (guess: false)

6

Example: Perceptron Learning

- Can perceptrons represent all Boolean functions?
 $f(\text{Feature}_1, \dots, \text{Feature}_n) \equiv \text{some propositional sentence}$

7

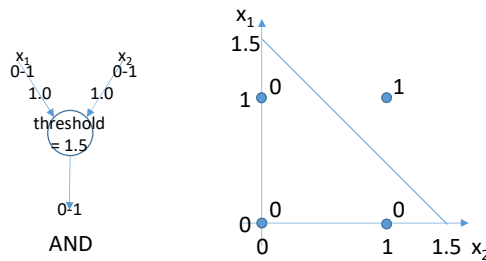
Example: Perceptron Learning

- Can perceptrons represent all Boolean functions?
 $f(\text{Feature}_1, \dots, \text{Feature}_n) \equiv \text{some propositional sentence}$
- Linear separability
 - We need to find an n -dimensional plane that separates the labeled examples with class true from the labeled examples with class false.
 - This plane determines the weights and threshold of the perceptron that can then be used to classify the unlabeled examples.

8

Example: Perceptron Learning

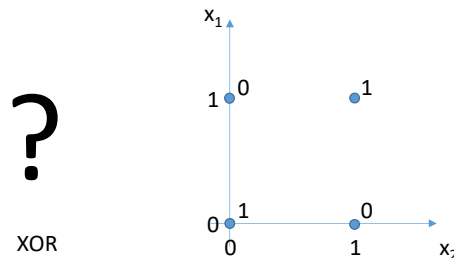
- Can perceptrons represent all Boolean functions?
 $f(\text{Feature}_1, \dots, \text{Feature}_n) \equiv$ some propositional sentence
- Linear separability
 - $w_1 x_1 + w_2 x_2 = \text{threshold}$
 - $w_1 x_1 = \text{threshold} - w_2 x_2$
 - $x_1 = (\text{threshold} / w_1) - (w_2 / w_1) x_2 = (1.5 / 1) - (1 / 1) x_2 = 1.5 - x_2$



9

Example: Perceptron Learning

- Can perceptrons represent all Boolean functions?
 $f(\text{Feature}_1, \dots, \text{Feature}_n) \equiv$ some propositional sentence
- Linear separability
 - $w_1 x_1 + w_2 x_2 = \text{threshold}$
 - $w_1 x_1 = \text{threshold} - w_2 x_2$
 - $x_1 = (\text{threshold} / w_1) - (w_2 / w_1) x_2 = (1.5 / 1) - (1 / 1) x_2 = 1.5 - x_2$



10

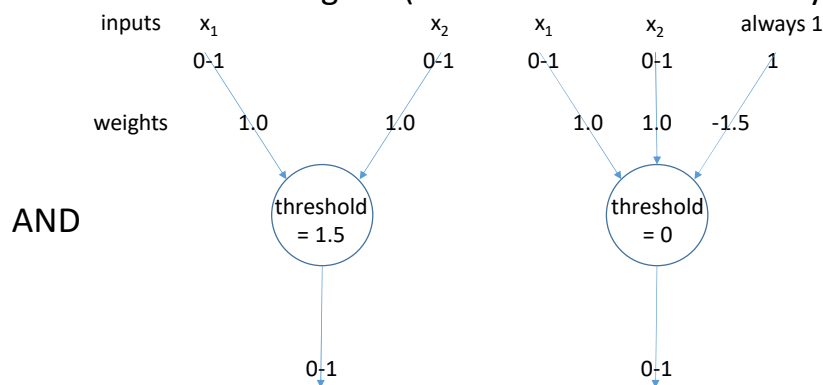
Example: Perceptron Learning

- Can perceptrons represent all Boolean functions? – no!
 $f(\text{Feature}_1, \dots, \text{Feature}_n) \equiv$ some propositional sentence
- An XOR cannot be represented with a single perceptron!
- This does not mean that single perceptrons should not be used. They will make some mistakes for some Boolean functions (that is, might not be able to classify all labeled examples correctly) but they often work well, that is, make few mistakes on the labeled and unlabeled examples. Of course, you only want to use them if they do not make too many mistakes on the labeled examples.

11

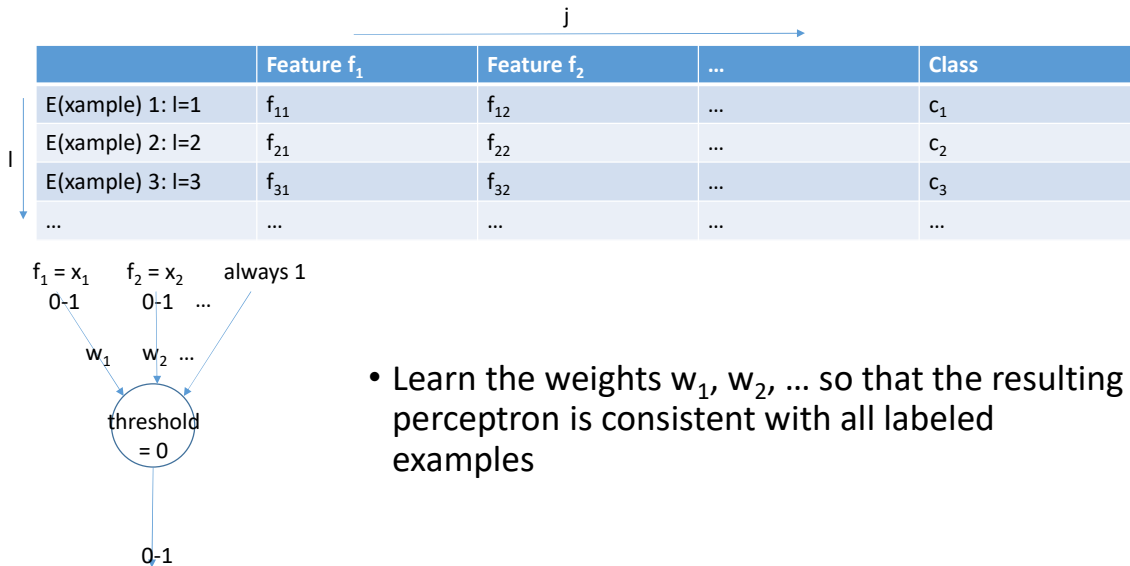
Example: Perceptron Learning

- The threshold can be expressed as a weight.
- This way, a learning algorithm only needs to learn weights instead of the threshold and the weights. (The new threshold is always zero.)



12

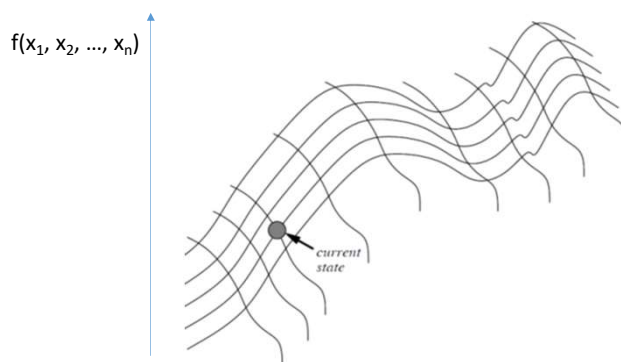
Example: Perceptron Learning



13

Gradient Descent

- Finding a **local** minimum of a **differentiable** function $f(x_1, x_2, \dots, x_n)$ with gradient descent



14

Gradient Descent

- Finding a **local** minimum of a **differentiable** function $f(x_1, x_2, \dots, x_n)$ with gradient descent
- Initialize x_1, x_2, \dots, x_n with random values
- Repeat until local minimum reached
 - Update x_1, x_2, \dots, x_n to correspond to taking a small step **against** the gradient of $f(x_1, x_2, \dots, x_n)$ at point (x_1, x_2, \dots, x_n) , where the gradient is $(d f(x_1, x_2, \dots, x_n) / d x_1, d f(x_1, x_2, \dots, x_n) / d x_2, \dots, d f(x_1, x_2, \dots, x_n) / d x_n)$.

15

Gradient Descent

- Finding a local minimum of a differentiable function $f(x_1, x_2, \dots, x_n)$ with gradient descent (for a small positive learning rate α)
- Initialize x_1, x_2, \dots, x_n with random values
- Repeat until local minimum reached
 - For all x_i in parallel
 - $x_i := x_i - \alpha d f(x_1, x_2, \dots, x_n) / d x_i$

16

Gradient Descent

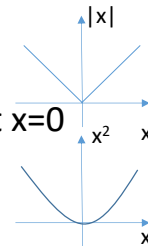
- Finding a local minimum of a differentiable function $f(x_1, x_2, \dots, x_n)$ with an **approximation** of gradient descent (for a small positive learning rate α)
- Initialize x_1, x_2, \dots, x_n with random values
- Repeat until local minimum reached
 - For all x_i
 - $x_i := x_i - \alpha \frac{d f(x_1, x_2, \dots, x_n)}{d x_i}$

17

Example: Perceptron Learning

- We use the number of misclassified labeled examples as error and learn the weights w_1, w_2, \dots with gradient descent (for a small positive learning rate α) to correspond to a (local) minimum of the error function, that is, so that the resulting perceptron is consistent with all labeled examples:

- Minimize Error := $0.5 \sum_l |o_l - c_l|$ - no: not differentiable at $x=0$
- Minimize Error := $0.5 \sum_l (o_l - c_l)^2$
- for $o_l = g(\sum_j w_j f_{ij})$, where $g()$ is the activation function.
- The 0.5 is for **beauty reasons** only (see the slide after the next one).

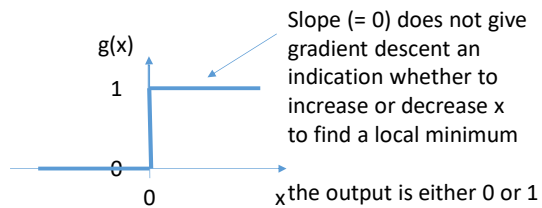


18

Example: Perceptron Learning

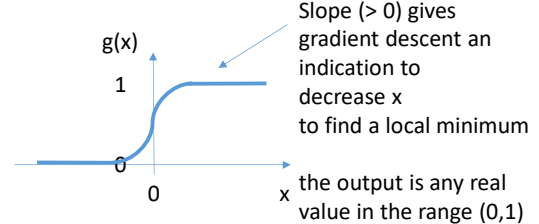
- Learn the weights w_1, w_2, \dots with gradient descent (for a small positive learning rate α) so that the resulting perceptron is consistent with all labeled examples:

Threshold function



no: not differentiable at $x=0$

Sigmoid function



$$g(x) = 1 / (1 + e^{-x})$$

$$g'(x) = e^{-x} / (1 + e^{-x})^2 = g(x) (1 - g(x))$$

19

Derivatives: Chain Rule

- Quick reminder of the chain rule (since we need it on the next slide):
 $d f(g(x)) / d x = f'(g(x)) g'(x)$
- For example, $d (2x)^2 / d x = 2(2x) 2 = 8x$
 by applying the chain rule since
 - $f(x) = x^2$ and $g(x) = 2x$
 - $f'(x) = 2x$ and $g'(x) = 2$
 - $f(g(x)) = (2x)^2$
- For example, $d (e^{2x})^2 / dx = 2(e^{2x}) e^{2x} 2 = 4 e^{4x}$
 by applying the chain rule twice in a row

20

Example: Perceptron Learning

- Learn the weights w_1, w_2, \dots with gradient descent (for a small positive learning rate α) so that the resulting perceptron is consistent with all labeled examples:
- Initialize all weights w_j with random values
- Repeat until local minimum reached
 - Let o_l be the output of the perceptron for Example l for the current weights
 - For all weights w_j in parallel
 - $w_j := w_j - \alpha \frac{d \text{Error}(w_1, w_2, \dots)}{d w_j}$ } called one epoch
- Where
 - $\frac{d \text{Error}(w_1, w_2, \dots)}{d w_j} = \frac{d}{d w_j} 0.5 \sum_l (o_l - c_l)^2 = \frac{d}{d w_j} 0.5 \sum_l (g(\sum_j w_j f_{lj}) - c_l)^2$
 $= \sum_l ((g(\sum_j w_j f_{lj}) - c_l) g'(\sum_j w_j f_{lj}) f_{lj}) = \sum_l ((o_l - c_l) g'(\sum_j w_j f_{lj}) f_{lj})$

This is the beauty reason!

21

Example: Perceptron Learning

- Learn the weights w_1, w_2, \dots with an **approximation** of gradient descent (for a small positive learning rate α) so that the resulting perceptron is consistent with all labeled examples. Each labeled example is considered individually one after the other:
- Initialize all weights w_j with random values
- Repeat until local minimum reached
 - Let o_l be the output of the perceptron for Example l for the current weights
 - For all labeled examples l
 - For all weights w_j
 - $w_j := w_j - \alpha \frac{d \text{Error}(w_1, w_2, \dots)}{d w_j}$ } called one epoch
- Where
 - $\frac{d \text{Error}(w_1, w_2, \dots)}{d w_j} = (o_l - c_l) g'(\sum_j w_j f_{lj}) f_{lj}$

22

Example: Perceptron Learning

Since the output is now any real value in the range (0,1), we consider a value less than 0.5 to be 0 and a value greater than 0.5 to be 1. So we indeed learned an AND!

- Example: Learn the weights w_1, w_2, \dots with an approximation of gradient descent (for $\alpha = 0.01$) so that the resulting perceptron is consistent with an AND [see the handout for details]

	Feature f_1	Feature f_2	Feature f_3	Class
E(xample) 1: $l=1$	0	0	1	0
E(xample) 2: $l=2$	0	1	1	0
E(xample) 3: $l=3$	1	0	1	0
E(xample) 4: $l=4$	1	1	1	1

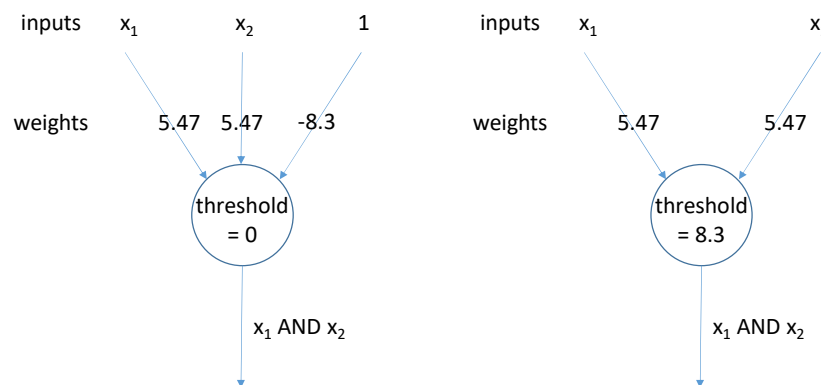
Epoch 0		Epoch 1		Epoch 2		Epoch 100		Epoch 100,000	
Weights	Outputs	Weights	Outputs	Weights	Outputs	Weights	Outputs	Weights	Outputs
$w_1 = 1.10$	$o_1 = 0.57$	1.10	0.57	1.10	0.57	1.12	0.54	5.47	0.00
$w_2 = -2.10$	$o_2 = 0.14$	-2.10	0.14	-2.10	0.14	-1.97	0.14	5.47	0.06
$w_3 = 0.30$	$o_3 = 0.80$	0.30	0.80	0.30	0.80	0.16	0.78	-8.30	0.06
	$o_4 = 0.33$		0.33		0.33		0.33		0.93

23

Example: Perceptron Learning

- Example: Learn the weights w_1, w_2, \dots with an approximation of gradient descent (for $\alpha = 0.01$) so that the resulting perceptron is consistent with an AND [see the handout for details]

- Result:



24