

Slides adapted from Prof. Carpuat



CMSC 422 Introduction to Machine Learning

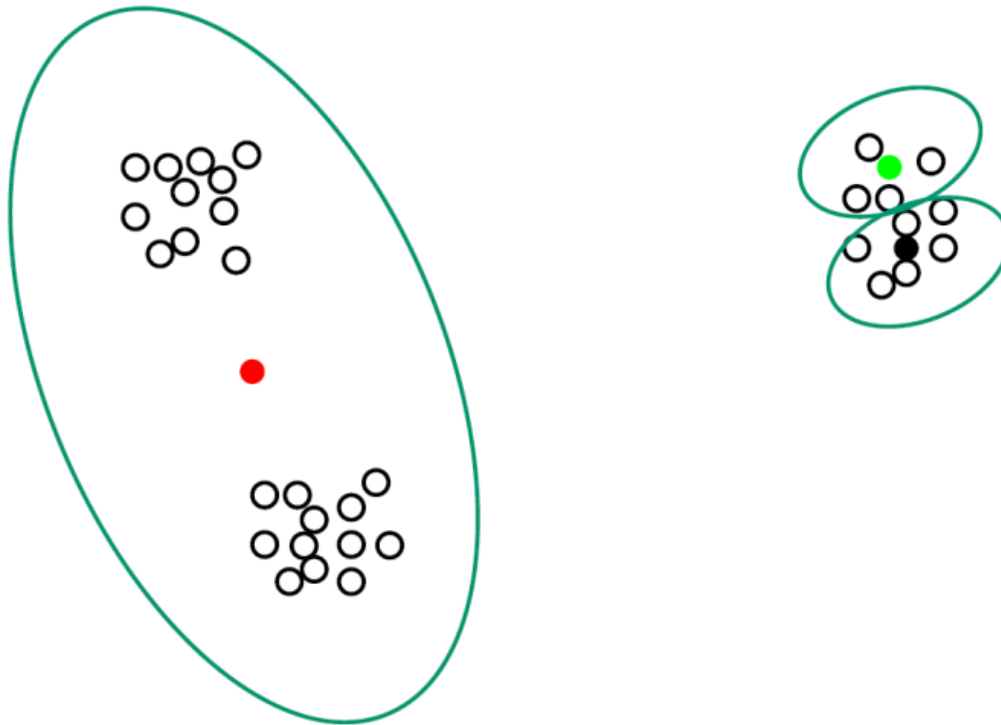
Lecture 6 The Perceptron

Furong Huang / furongh@cs.umd.edu



UNIVERSITY OF
MARYLAND

Impact of initialization



Optimization View of K-means

Given a set of observations (x_1, x_2, \dots, x_n) where each observation is a d -dimensional real vector

k-means clustering aims to partition the n observations into $k(\leq x)$ sets $S = \{S_1, S_2, \dots, S_k\}$ so as to minimize the within-cluster sum of squares.

Formally, the objective is to find:

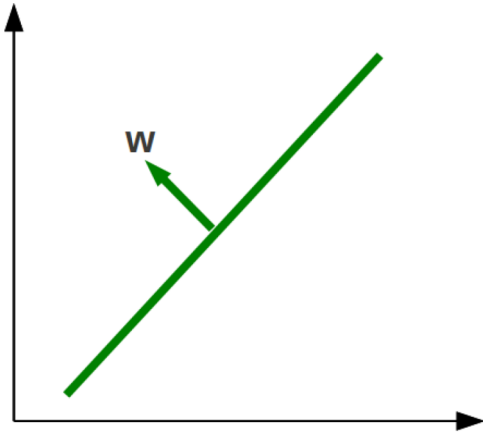
$$\arg \min_S \sum_{i=1}^k \sum_{x \in S_i} \|x - \mu_i\|^2 = \arg \min_S \sum_i^k |S_i| \text{Var}(S_i)$$

where μ_i is the mean of points in S_i .

This week

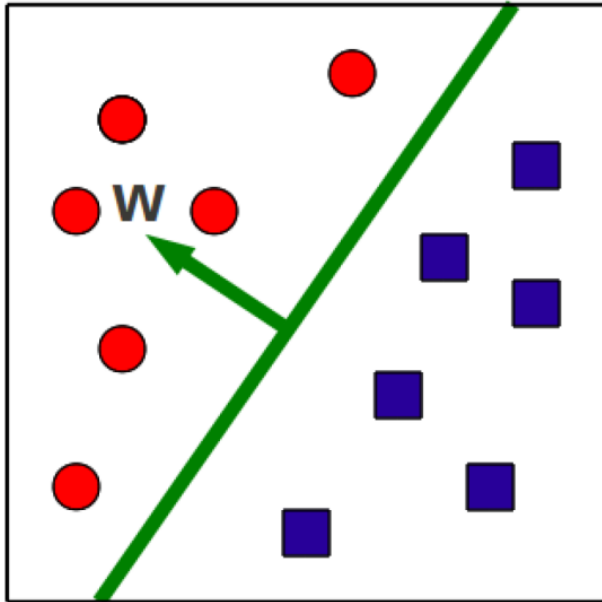
- A new model/algorithm
 - the perceptron
 - and its variants: voted, averaged
- Fundamental Machine Learning Concepts
 - Online vs. batch learning
 - Error-driven learning
- Project 1 coming soon!

Geometry concept: Hyperplane



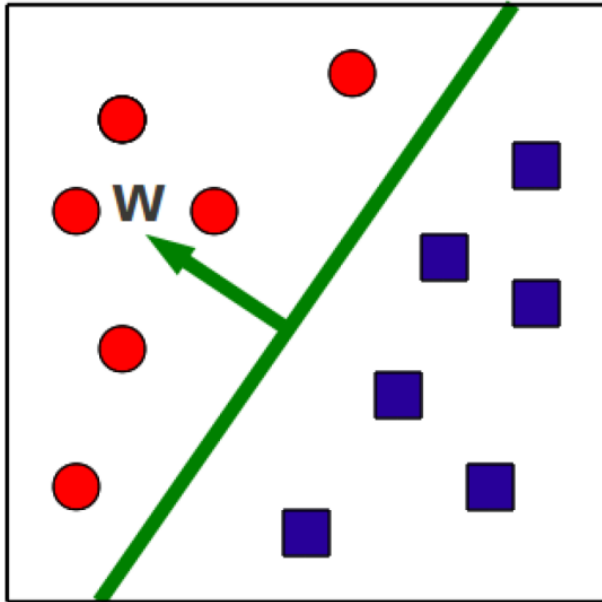
- ❑ Separates a D-dimensional space into two half-spaces
- ❑ Defined by an outward pointing normal vector $w \in \mathbb{R}^D$
 - ❑ w is **orthogonal** to any vector lying on the hyperplane
- ❑ Hyperplane passes through the origin, unless we also define a **bias** term b

Binary classification via hyperplanes



- ❑ Let's assume that the decision boundary is a hyperplane
- ❑ Then, training consists in finding a hyperplane w that separates positive from negative examples

Binary classification via hyperplanes



At test time, we check on
what side of the hyperplane
examples fall

$$\hat{y} = \text{sign}(w^T x + b)$$

Function Approximation with Perceptron

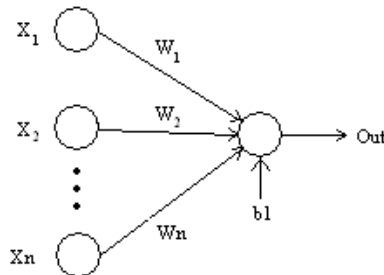
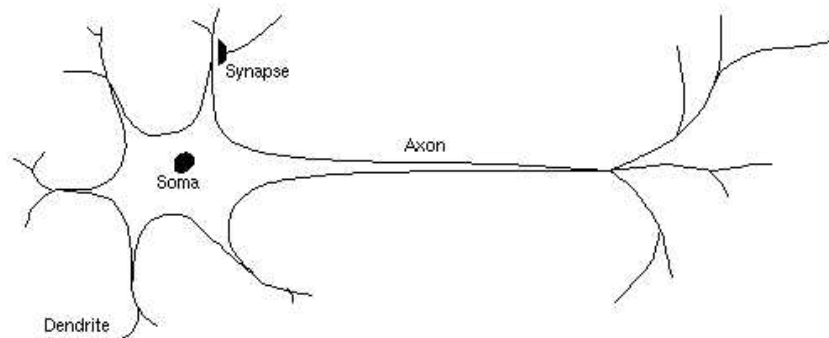
- Problem setting
- Set of possible instances X
 - Each instance $x \in X$ is a feature vector $x = [x_1, \dots, x_D]$
- Unknown target function $f: X \rightarrow Y$
 - Y is binary valued $\{-1; +1\}$
- Set of function hypotheses $H = \{h \mid h: X \rightarrow Y\}$
 - Each hypothesis h is a hyperplane in D -dimensional space
- Input
- Training examples $\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$ of unknown target function f
- Output
- Hypothesis $h \in H$ that best approximates target function f

Perception: Prediction Algorithm

Algorithm 6 PERCEPTRONTEST($w_0, w_1, \dots, w_D, b, \hat{x}$)

1: $a \leftarrow \sum_{d=1}^D w_d \hat{x}_d + b$ // compute activation for the test example
2: **return** SIGN(a)

Aside: biological inspiration



Analogy: the
perceptron
as a neuron

Perceptron Training Algorithm

Algorithm 5 PERCEPTRONTRAIN(\mathbf{D} , $MaxIter$)

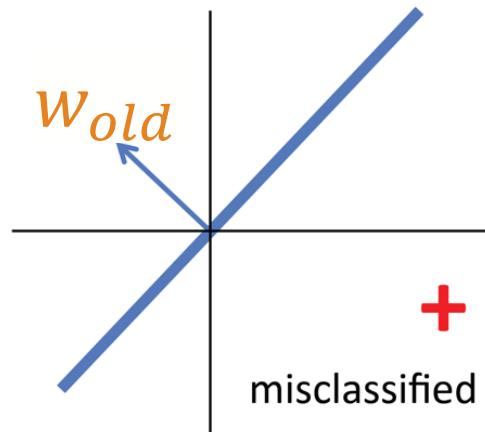
```
1:  $w_d \leftarrow 0$ , for all  $d = 1 \dots D$  // initialize weights
2:  $b \leftarrow 0$  // initialize bias
3: for  $iter = 1 \dots MaxIter$  do
4:   for all  $(x, y) \in \mathbf{D}$  do
5:      $a \leftarrow \sum_{d=1}^D w_d x_d + b$  // compute activation for this example
6:     if  $ya \leq 0$  then
7:        $w_d \leftarrow w_d + yx_d$ , for all  $d = 1 \dots D$  // update weights
8:        $b \leftarrow b + y$  // update bias
9:     end if
10:   end for
11: end for
12: return  $w_0, w_1, \dots, w_D, b$ 
```

Perceptron update: geometric interpretation

A training example (\mathbf{x}, y) is misclassified, i.e.,

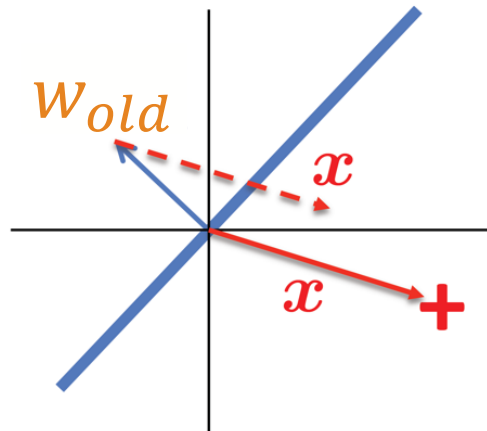
$$\text{sign}(\mathbf{w}_{old}^T \mathbf{x} + b) \neq y$$

Let's say $y = +1$



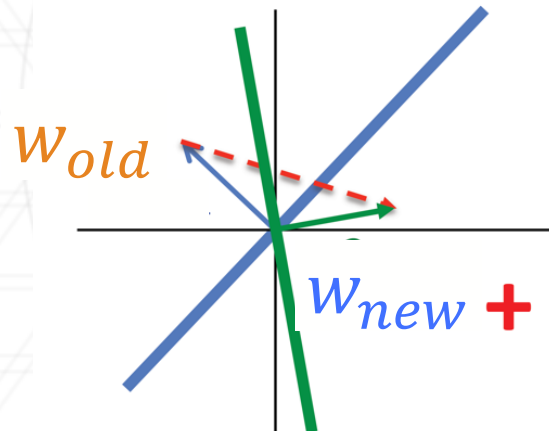
Perceptron update: geometric interpretation

Update: $\mathbf{w}_{new} \leftarrow \mathbf{w}_{old} + y\mathbf{x}$, i.e., $\mathbf{w}_{new} \leftarrow \mathbf{w}_{old} + \mathbf{x}$



Perceptron update: geometric interpretation

Update: $\mathbf{w}_{new} \leftarrow \mathbf{w}_{old} + y\mathbf{x}$, i.e., $\mathbf{w}_{new} \leftarrow \mathbf{w}_{old} + \mathbf{x}$



Properties of the Perceptron training algorithm

❑ Online

- ❑ We look at one example at a time, and update the model as soon as we make an error
- ❑ **As opposed to batch** algorithms that update parameters after seeing the entire training set

❑ Error-driven

- ❑ We only update parameters/model if we make an error

Practical considerations

- ❑ The order of training examples matters!
 - ❑ Random is better
- ❑ Early stopping
 - ❑ Good strategy to avoid overfitting
- ❑ Simple modifications dramatically improve performance
 - ❑ voting or averaging

Predicting with

- The voted perceptron

$$\hat{y} = \text{sign} \left(\sum_{k=1}^K c^{(k)} \text{sign} \left(\mathbf{w}^{(k)} \cdot \hat{\mathbf{x}} + b^{(k)} \right) \right)$$

- The averaged perceptron

$$\hat{y} = \text{sign} \left(\sum_{k=1}^K c^{(k)} \left(\mathbf{w}^{(k)} \cdot \hat{\mathbf{x}} + b^{(k)} \right) \right)$$

Require keeping track of “survival time” of weight vectors

$$c^{(1)}, \dots, c^{(K)}$$

How would you modify this algorithm for voted perceptron?

Algorithm 5 PERCEPTRONTRAIN(\mathbf{D} , $MaxIter$)

```
1:  $w_d \leftarrow 0$ , for all  $d = 1 \dots D$  // initialize weights
2:  $b \leftarrow 0$  // initialize bias
3: for  $iter = 1 \dots MaxIter$  do
4:   for all  $(x, y) \in \mathbf{D}$  do
5:      $a \leftarrow \sum_{d=1}^D w_d x_d + b$  // compute activation for this example
6:     if  $ya \leq 0$  then
7:        $w_d \leftarrow w_d + yx_d$ , for all  $d = 1 \dots D$  // update weights
8:        $b \leftarrow b + y$  // update bias
9:     end if
10:   end for
11: end for
12: return  $w_0, w_1, \dots, w_D, b$ 
```

How would you modify this algorithm for averaged perceptron?

Algorithm 5 PERCEPTRONTRAIN(\mathbf{D} , $MaxIter$)

```
1:  $w_d \leftarrow 0$ , for all  $d = 1 \dots D$  // initialize weights
2:  $b \leftarrow 0$  // initialize bias
3: for  $iter = 1 \dots MaxIter$  do
4:   for all  $(x, y) \in \mathbf{D}$  do
5:      $a \leftarrow \sum_{d=1}^D w_d x_d + b$  // compute activation for this example
6:     if  $ya \leq 0$  then
7:        $w_d \leftarrow w_d + yx_d$ , for all  $d = 1 \dots D$  // update weights
8:        $b \leftarrow b + y$  // update bias
9:     end if
10:   end for
11: end for
12: return  $w_0, w_1, \dots, w_D, b$ 
```

Averaged perceptron decision rule

$$\hat{y} = \text{sign} \left(\sum_{k=1}^K c^{(k)} \left(\boldsymbol{w}^{(k)} \cdot \hat{\boldsymbol{x}} + b^{(k)} \right) \right)$$

can be rewritten as

$$\hat{y} = \text{sign} \left(\left(\sum_{k=1}^K c^{(k)} \boldsymbol{w}^{(k)} \right) \cdot \hat{\boldsymbol{x}} + \sum_{k=1}^K c^{(k)} b^{(k)} \right)$$

Averaged Perceptron Training

Algorithm 7 AVERAGEDPERCEPTRONTRAIN(\mathbf{D} , $MaxIter$)

```
1:  $\mathbf{w} \leftarrow \langle 0, 0, \dots, 0 \rangle$  ,  $b \leftarrow 0$  // initialize weights and bias
2:  $\mathbf{u} \leftarrow \langle 0, 0, \dots, 0 \rangle$  ,  $\beta \leftarrow 0$  // initialize cached weights and bias
3:  $c \leftarrow 1$  // initialize example counter to one
4: for  $iter = 1 \dots MaxIter$  do
5:   for all  $(\mathbf{x}, y) \in \mathbf{D}$  do
6:     if  $y(\mathbf{w} \cdot \mathbf{x} + b) \leq 0$  then
7:        $\mathbf{w} \leftarrow \mathbf{w} + y \mathbf{x}$  // update weights
8:        $b \leftarrow b + y$  // update bias
9:        $\mathbf{u} \leftarrow \mathbf{u} + y c \mathbf{x}$  // update cached weights
10:       $\beta \leftarrow \beta + y c$  // update cached bias
11:     end if
12:      $c \leftarrow c + 1$  // increment counter regardless of update
13:   end for
14: end for
15: return  $\mathbf{w} - \frac{1}{c} \mathbf{u}$ ,  $b - \frac{1}{c} \beta$  // return averaged weights and bias
```

Can the perceptron always find a hyperplane to separate positive from negative examples?

This week

A new model/algorithm

the perceptron

and its variants: voted, averaged

Fundamental Machine Learning Concepts

Online vs. batch learning

Error-driven learning

Project 1 coming soon!



Furong Huang

3251 A.V. Williams, College Park, MD 20740

301.405.8010 / furongh@cs.umd.edu