

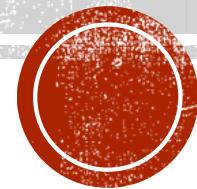
THE PERCEPTRON

CS 412 Introduction to Machine Learning

Prof. Zheleva

February 1, 2018

Reading Assignment: CML: 4



LAST LECTURE: NEAREST NEIGHBORS

- From real-world problem to data to model to solution
- Using geometry for classification
 - Curse of dimensionality
- kNN: Very simple non-linear classification technique
 - Works well with large amounts of data
 - Related to very sophisticated machine learning technique: kernel methods!
- Decision boundaries

FROM DATA TO ML SOLUTION

Why one of Africa's worst conflicts is getting worse

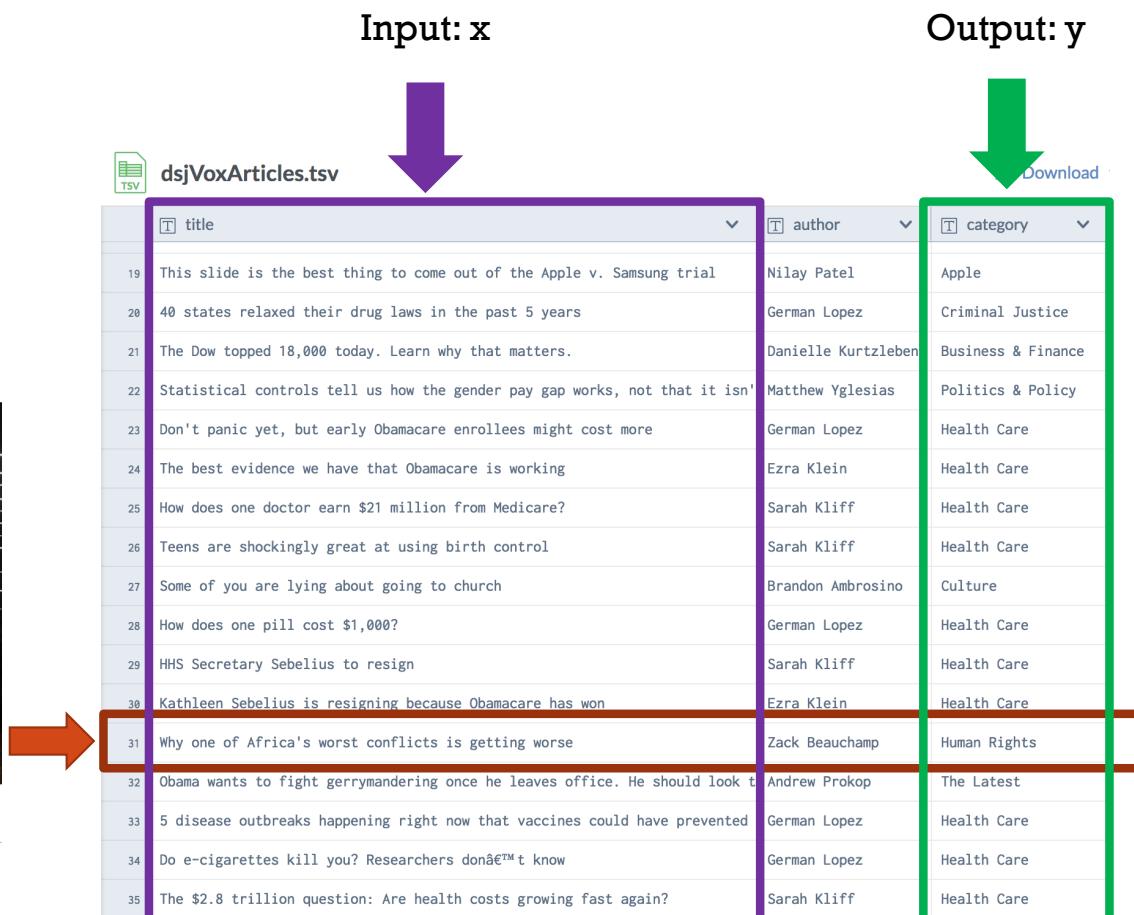
By Zack Beauchamp | @zackbeauchamp | zack@vox.com | Apr 12, 2014, 11:00am EDT

[SHARE](#) [MORE](#)



CAR civilians besieged by the anti-balaka militias take shelter under French military protection.
| Miguel Medina/AFP/Getty Images

The conflict in the Central African Republic (CAR) is a huge international issue — thousands of people have been killed since it began in December 2012, and its worsening rapidly enough that the UN has just greenlit a



Source: <https://data.world/elenadata/vox-articles>

GEOMETRIC VIEW OF DATA

- Each feature is a dimension
- The feature values of each data instance correspond to a point in n-dimensional space where n is the number of features

Features: bag of words



Feature values: presence of a word

aardvark	africa	beach	conflict	...	worse	zebra
0	1	0	1	...	1	0

Feature values: count of a word

aardvark	africa	beach	conflict	...	worse	zebra
0	5	0	2	...	1	0



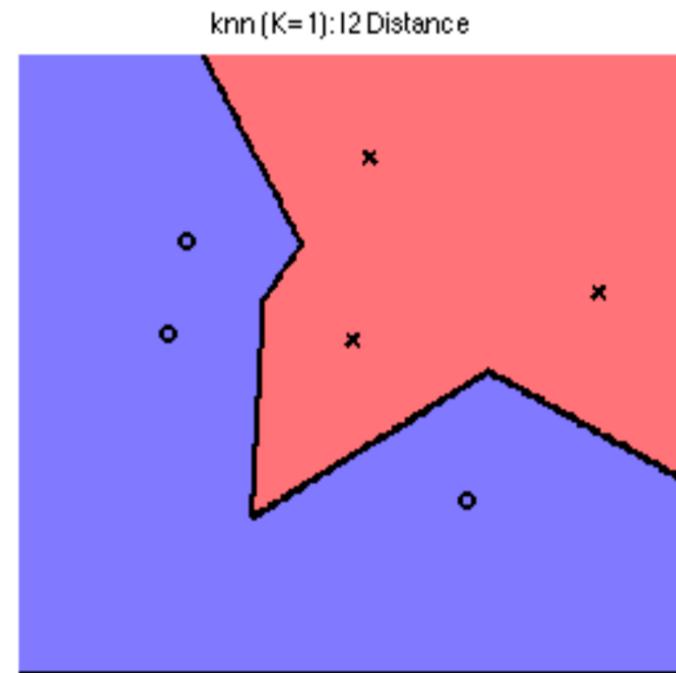
COMPONENTS OF A K-NN CLASSIFIER

- Distance metric
 - How do we measure distance between instances
 - Determines the layout of the example space
- The k hyperparameter
 - How large a neighborhood should we consider?
 - Determines the complexity of the hypothesis space



DECISION BOUNDARY OF A CLASSIFIER

- The line that separates positive and negative regions in the feature space
- Why is it useful?
 - It helps us visualize how examples will be classified for the entire feature space
 - It helps us visualize the complexity of the learned model
- Decision boundary for 1NN



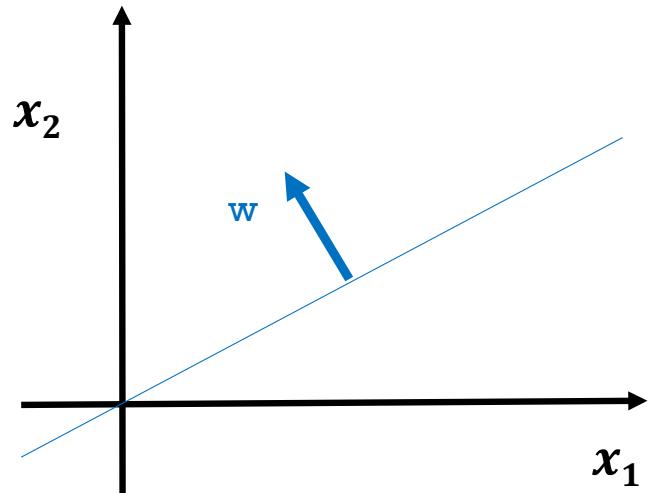
TODAY: THE PERCEPTRON

- The perceptron: a new model and algorithm
 - Its variants: voted, averaged
 - Convergence proof
- Fundamental machine learning concepts
 - Online vs. batch learning
 - Error-driven learning
 - Linear separability and margin of a dataset



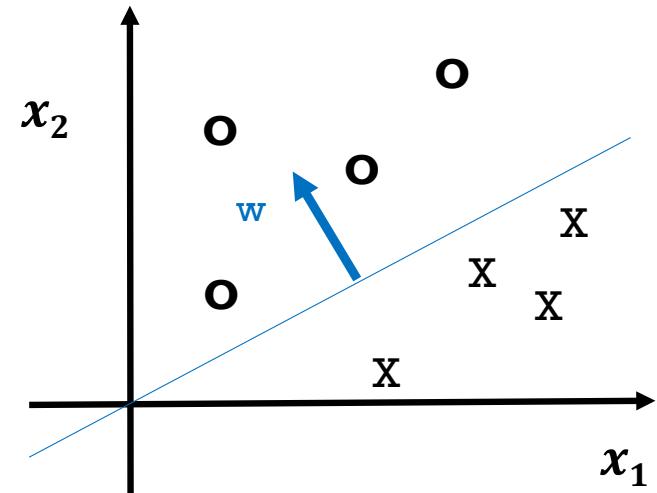
GEOMETRY CONCEPT: HYPERPLANE

- 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

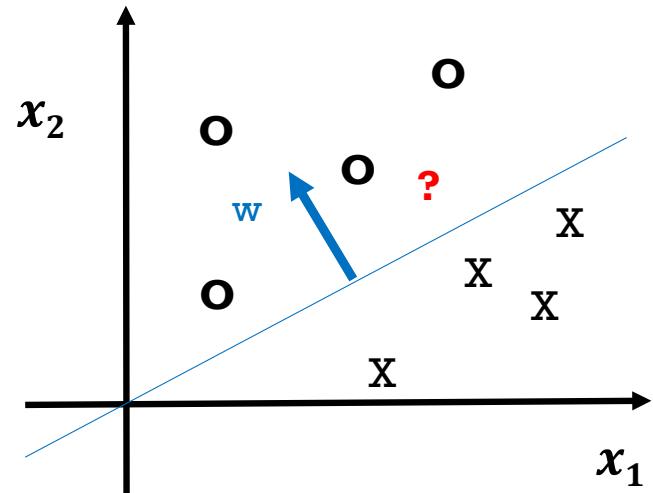
- Let's assume that the decision boundary is a hyperplane
- Then, training consists in finding a hyperplane w that separates positive from negative examples



BINARY CLASSIFICATION VIA HYPERPLANES

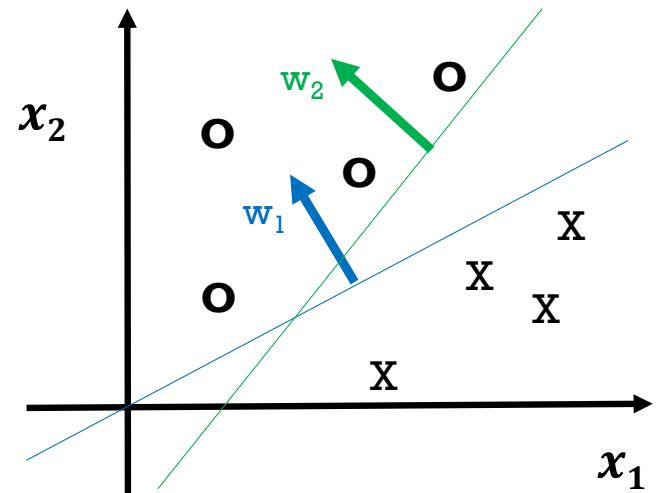
- At test time, we check on what side of the hyperplane examples fall
- Classifier = hyperplane that separates positive from negative examples

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



FUNCTION APPROXIMATION WITH PERCEPTRON

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



PERCEPTRON: PREDICTION ALGORITHM

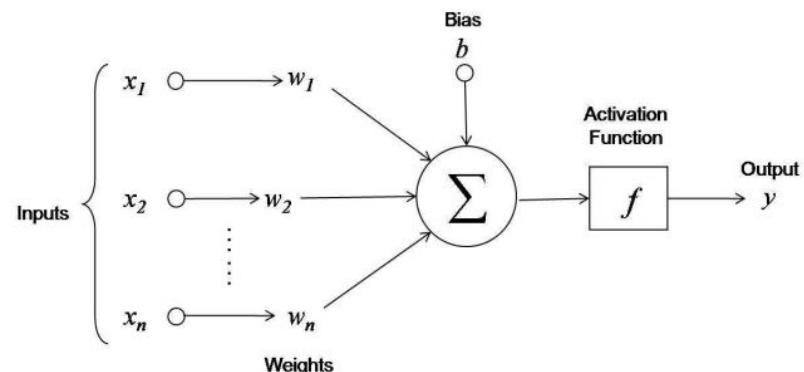
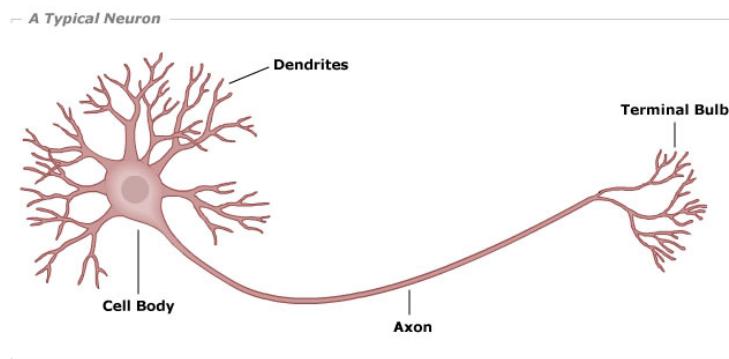
Algorithm 6 PERCEPTRONTEST($w_0, w_1, \dots, w_D, b, \hat{x}$)

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

- The perceptron as a neuron
 - Can receive both stimulatory and inhibitory signals
 - If adequate stimulus, it produces a nerve impulse



$$a \leftarrow \sum_{d=1}^D w_d x_d + b$$



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  $w_0, w_1, \dots, w_D, b$ 
```

epoch

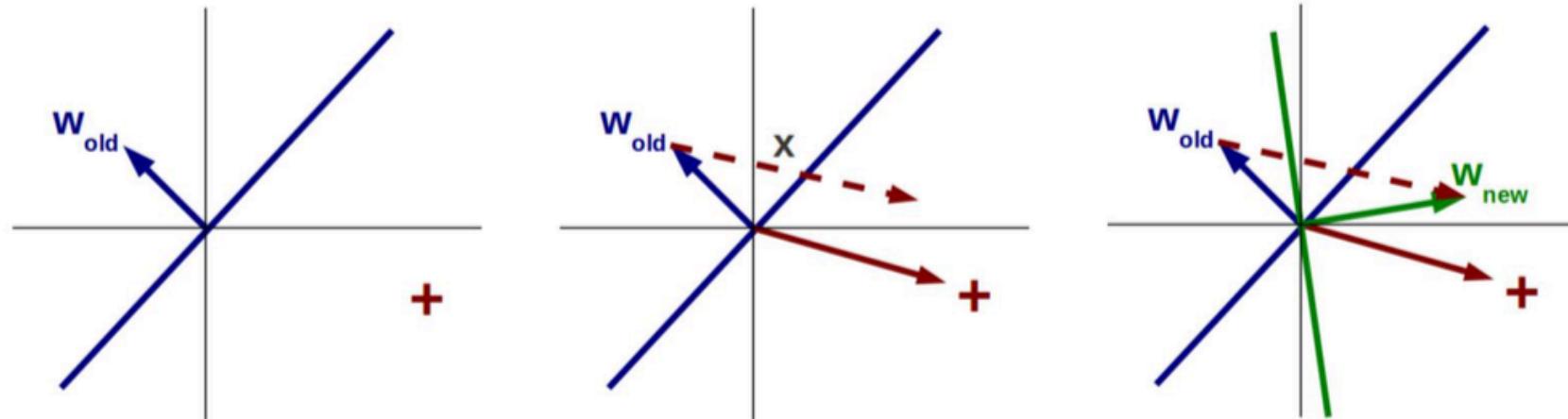
Non-obvious why this update works!



PERCEPTRON UPDATE: GEOMETRIC INTERPRETATION

Update for a misclassified positive example (assume $b=0$):

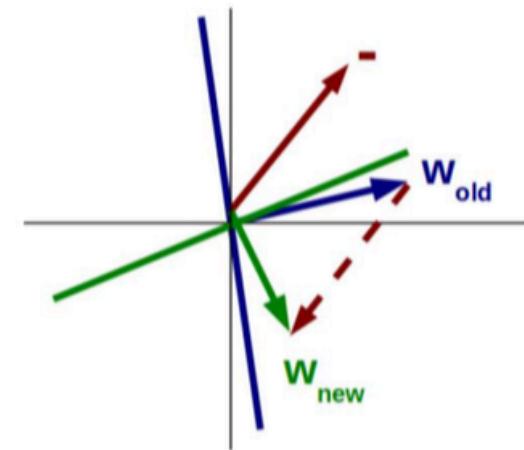
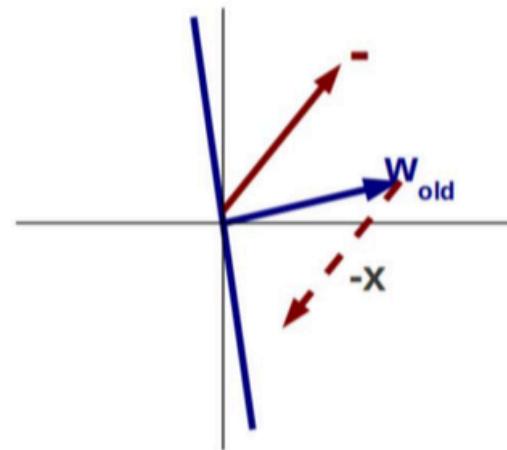
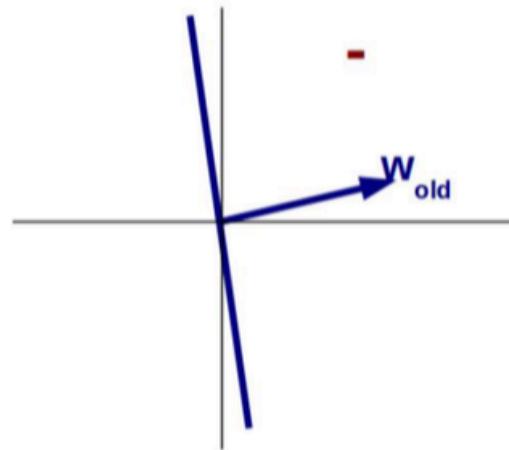
$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} + \mathbf{x}$$



PERCEPTRON UPDATES

Update for a misclassified negative example:

$$\mathbf{w}_{new} = \mathbf{w}_{old} - \mathbf{x}$$



PREDICTION VARIANTS

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( $a$ )
```

- The voted perceptron
 - Better than vanilla, can be slow

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

- The averaged perceptron

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

- Requires keeping track of “survival time” of weight vectors $\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(K)}$



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 TRAINING

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

can be re-written as

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



AVERAGED PERCEPTRON TRAINING

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

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



PERCEPTRON TRAINING ALGORITHM

- **Online algorithm**
 - We look at one example at a time, and update the model when 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!
 - Randomized is better
 - Permuting the order in every iteration makes learning much faster
- Early stopping
 - Good strategy to avoid overfitting
- Simple modifications dramatically improve performance
 - Voting or averaging



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



CONVERGENCE OF TRAINING

- The perceptron has converged if it can classify every training example correctly
 - i.e. if it has found a hyperplane that correctly separates positive and negative examples
- Under which conditions does the perceptron converge and how long does it take?



CONVERGENCE OF PERCEPTRON

Theorem (Block & Novikoff, 1962):

If the training data $D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$ is linearly separable with margin γ by a unit norm hyperplane w_* , assuming $\|x^{(i)}\| \leq R$ for all i and $b=0$, then perceptron converges after $\frac{R^2}{\gamma^2}$ errors during training.



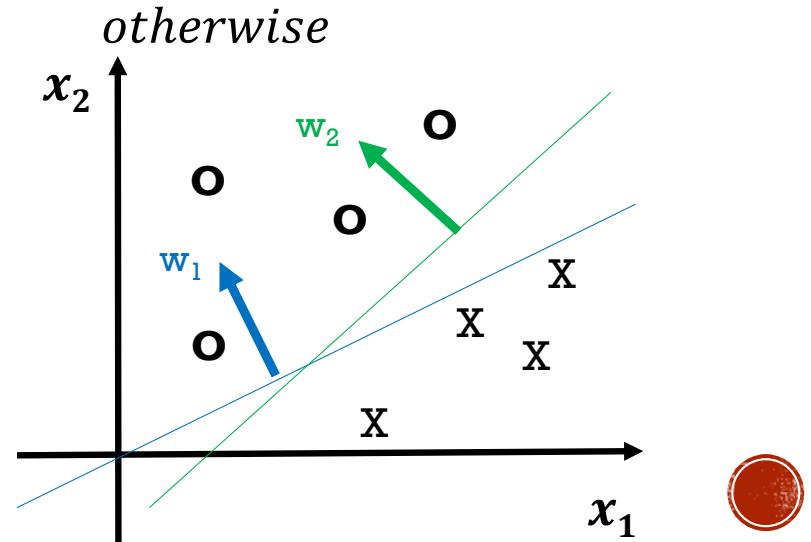
DATASET MARGIN

- Hyperplane margin on D : Distance between the hyperplane (w, b) and the nearest point in D

$$\text{margin}(D, w, b) = \begin{cases} \min_{(x,y) \in D} y(w \cdot x + b) & \text{if } w \text{ separates } D \\ -\infty & \text{otherwise} \end{cases}$$

- Dataset margin: Largest attainable margin on D

$$\text{margin}(D) = \sup_{w,b} \text{margin}(D, w, b)$$



CONVERGENCE OF PERCEPTRON

Theorem (Block & Novikoff, 1962):

If the training data $D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$ is linearly separable with margin γ by a unit norm hyperplane w_* , assuming $\|x^{(i)}\| \leq R$ for all i and $b=0$, then perceptron converges after $\frac{R^2}{\gamma^2}$ errors during training.



CONVERGENCE OF PERCEPTRON

Theorem (Block & Novikoff, 1962):

If the training data $D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$ is linearly separable with margin γ by a unit norm hyperplane w_* , assuming $\|x^{(i)}\| \leq R$ for all i and $b=0$, then perceptron converges after $\frac{R^2}{\gamma^2}$ errors during training.

What does it mean?

- Perceptron converges quickly when margin is large, slowly when it is small
- Bound does not depend on number of training examples N , nor on number of features
- Proof guarantees that perceptron converges but not necessarily to the max margin separator



DECISION TREES VS. KNN VS. PERCEPTRON

Properties of classification problem	Decision trees	k-NN	Perceptron
Binary features			
Numeric features			
Categorical features			
Robust to noisy training examples			
Fast classification is crucial			
Many irrelevant features			
Relevant features have very different scale			



SUMMARY

- Perceptron concepts
 - training/prediction algorithms (standard, voting, averaged)
 - convergence theorem and what practical guarantees it gives us
 - how to draw/describe the decision boundary of a perceptron classifier
- Fundamental ML concepts
 - Determine whether a data set is linearly separable and define its margin
 - Error driven algorithms, online vs. batch algorithms

ANNOUNCEMENTS

- HW 2 is out
 - due February 20, 11:59pm



ACKNOWLEDGEMENTS

- These slides use materials by Marine Carpuat, Piyush Rai and Hal Daume III

