

Perceptron

Sriram Sankararaman

The instructor gratefully acknowledges Fei Sha, Ameet Talwalkar, Eric Eaton, and Jessica Wu whose slides are heavily used, and the many others who made their course material freely available online.

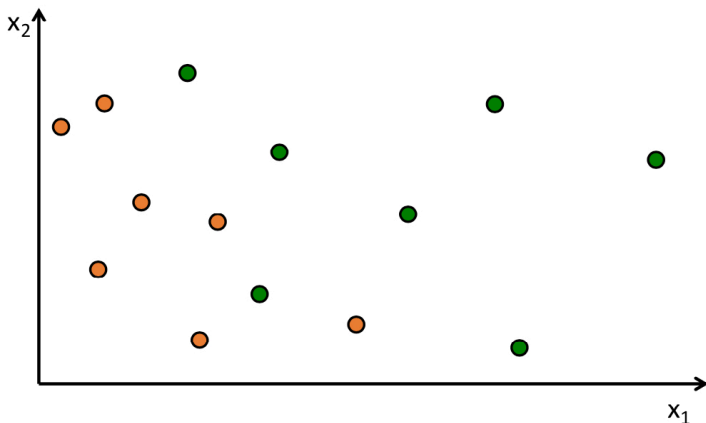
Key issues in machine learning

- Modeling
 - ▶ How to formulate the problem ?
- Representation
 - ▶ What is the input/output space ?
 - ▶ What is the model/ hypothesis space?
- Algorithms
 - ▶ How to find the best hypothesis?

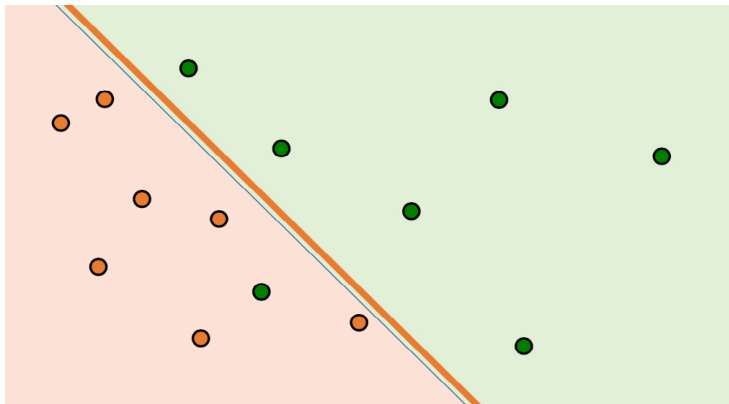
Outline

- 1 Perceptron
 - Setup for binary classification
- 2 Perceptron learning
- 3 Convergence of the Perceptron Learning Algorithm
- 4 Variants of perceptron
- 5 What we have learned

Training data for binary classification



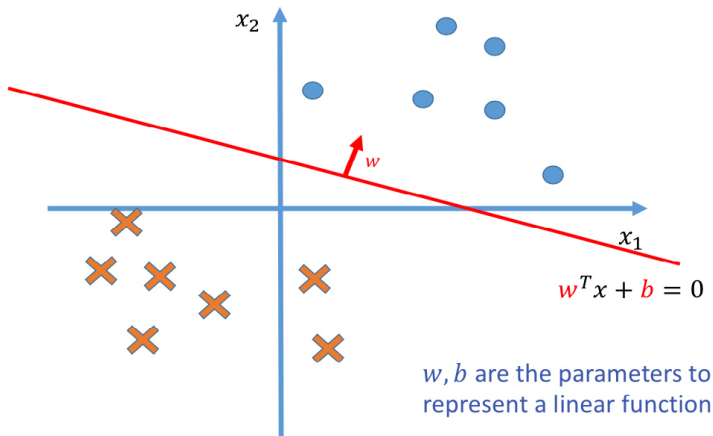
Linear classifier



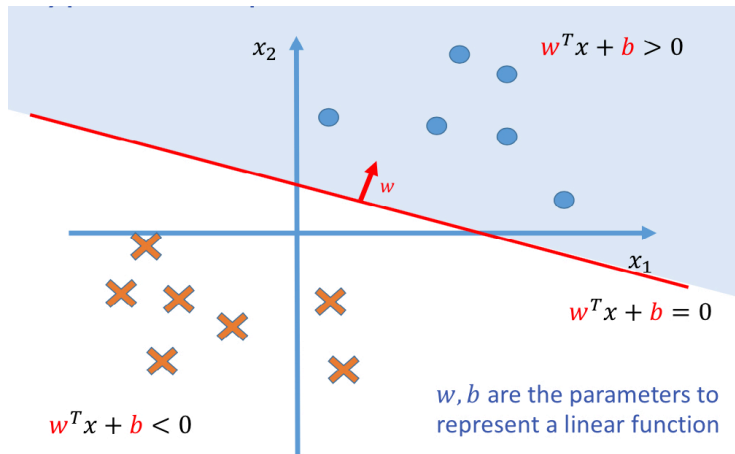
Prediction using a linear classifier



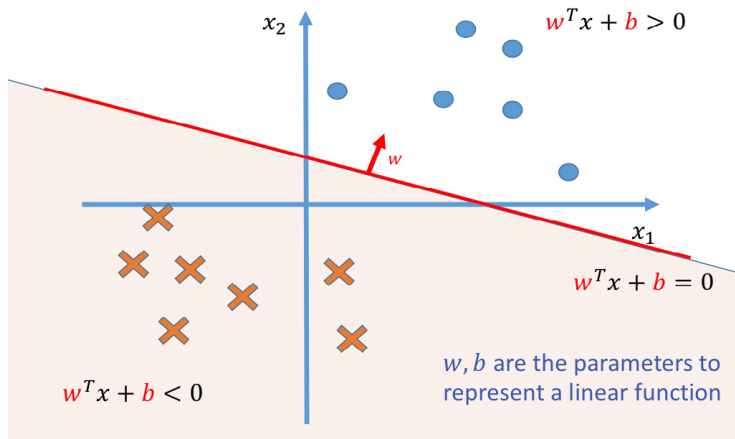
Representing a linear classifier



Representing a linear classifier



Classification using a linear classifier



Perceptron learning

Binary classification

- Instance (**feature vectors**): $\mathbf{x} \in \mathbb{R}^D$
- **Label**: $y \in \{-1, +1\}$
- **Model/Hypotheses**:
 $H = \{h | h : \mathbb{X} \rightarrow \mathbb{Y}, h(\mathbf{x}) = \text{sign}(\sum_{d=1}^D w_d x_d + b)\}.$
- Learning goal: $\hat{y} = h(\mathbf{x})$
 - ▶ Learn w_1, \dots, w_D, b .
 - ▶ **Parameters**: w_1, \dots, w_D, b .
 - ▶ w : **weights**, b : **bias**

Perceptron predict

- Input: $\mathbf{x} \in \mathbb{R}^D$, $\mathbf{w} \in \mathbb{R}^D$, $b \in \mathbb{R}$.

$$a = \sum_{d=1}^D w_d x_d + b = \mathbf{w}^T \mathbf{x} + b$$

$$\hat{y} = \text{sign}(a)$$

- Output: \hat{y} .
- $\sum_{d=1}^D w_d x_d + b = \mathbf{w}^T \mathbf{x} + b = 0$: hyperplane in D dimensions with parameters (\mathbf{w}, b) .
- \mathbf{w} : weights, b : bias
- a : activation
- $\text{sign}(\sum_{d=1}^D w_d x_d + b)$: Linear Threshold Unit (LTU)

Hyperplanes through the origin

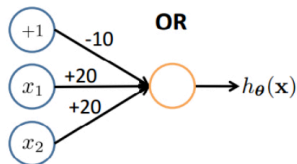
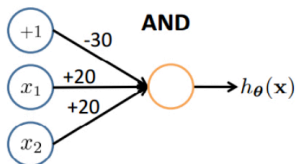
Consider \mathbf{x} that satisfies $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$. These \mathbf{x} define a hyperplane in D dimensions.

We can always write this as a hyperplane passing through the origin in $D + 1$ dimensions.

$$\begin{aligned}\tilde{\mathbf{x}} &\equiv \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_D \end{pmatrix} \quad \tilde{\mathbf{w}} \equiv \begin{pmatrix} b \\ w_1 \\ \vdots \\ w_D \end{pmatrix} \\ \tilde{g}(\tilde{\mathbf{x}}) &= \tilde{\mathbf{w}}^T \tilde{\mathbf{x}} \\ &= \sum_{d=1}^D w_d x_d + b \\ &= g(\mathbf{x})\end{aligned}$$

For simplicity, I may write $\tilde{\mathbf{w}}$ and $\tilde{\mathbf{x}}$ as \mathbf{w} and \mathbf{x} when there is no confusion

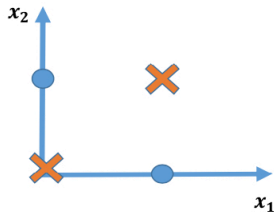
Representing Boolean functions



Representing Boolean functions

Can linear model represent XOR?

x_1	x_2	y
0	0	0
1	0	1
0	1	1
1	1	0



Learning a linear classifier

Several algorithms

- Perceptron
- Logistic regression
- (Linear) Support Vector Machines

Based on different assumptions

Outline

- 1 Perceptron
- 2 Perceptron learning**
- 3 Convergence of the Perceptron Learning Algorithm
- 4 Variants of perceptron
- 5 What we have learned

Perceptron learning

Learning by making mistakes

mistake
+
correction
=
learning

Perceptron learning

If we have only one training example (\mathbf{x}_n, y_n) .

Assume $b = 0$.

How can we change \mathbf{w} such that

$$y_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$$

Two cases

- If $y_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$, do nothing.
- If $y_n \neq \text{sign}(\mathbf{w}^T \mathbf{x}_n)$,

$$\mathbf{w}^{\text{NEW}} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$

Perceptron learning

If we have only one training example (\mathbf{x}_n, y_n) .

Assume $b = 0$.

How can we change \mathbf{w} such that

$$y_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$$

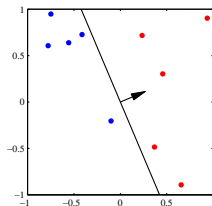
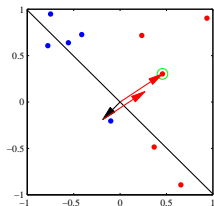
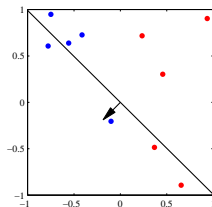
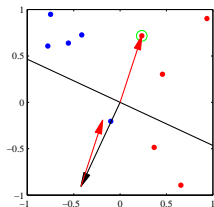
Another way of saying the same thing

- $a = \mathbf{w}^T \mathbf{x}_n$
- If $y_n a > 0$, do nothing.
- If $y_n a \leq 0$,

$$\mathbf{w}^{\text{NEW}} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$

Example of perceptron update

Red is $+1$, Blue is -1



Why would it work?

If $y_n a \leq 0$, then

$$y_n(\mathbf{w}^T \mathbf{x}_n) \leq 0$$

Why would it work?

If $y_n a \leq 0$, then

$$y_n(\mathbf{w}^T \mathbf{x}_n) \leq 0$$

What would happen if we change to new $\mathbf{w}^{\text{NEW}} = \mathbf{w} + y_n \mathbf{x}_n$?

$$y_n[(\mathbf{w} + y_n \mathbf{x}_n)^T \mathbf{x}_n] = y_n \mathbf{w}^T \mathbf{x}_n + y_n^2 \mathbf{x}_n^T \mathbf{x}_n$$

Why would it work?

If $y_n a \leq 0$, then

$$y_n(\mathbf{w}^T \mathbf{x}_n) \leq 0$$

What would happen if we change to new $\mathbf{w}^{\text{NEW}} = \mathbf{w} + y_n \mathbf{x}_n$?

$$y_n[(\mathbf{w} + y_n \mathbf{x}_n)^T \mathbf{x}_n] = y_n \mathbf{w}^T \mathbf{x}_n + y_n^2 \mathbf{x}_n^T \mathbf{x}_n$$

We are adding a positive number, so it is possible that

$$y_n(\mathbf{w}^{\text{NEW}T} \mathbf{x}_n) > 0$$

i.e., we are more likely to classify correctly

Perceptron learning

Iteratively solving one case at a time

- REPEAT
- Pick a data point \mathbf{x}_n
- Compute $a = \mathbf{w}^T \mathbf{x}_n$ using the *current* \mathbf{w}
- If $ay_n > 0$, do nothing. Else,

$$\mathbf{w} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$

- UNTIL converged.

Perceptron training/learning

$data = N$ **samples/instances**: $= \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$

Algorithm 1 PerceptronTrain ($data, maxIter$)

```
1:  $\mathbf{w} \leftarrow \mathbf{0}$ 
2: for  $iter = 1 \dots MaxIter$  do
3:   for  $(\mathbf{x}, y) \in data$  do
4:      $a \leftarrow \mathbf{w}^T \mathbf{x}$ 
5:     if  $ay \leq 0$  then
6:        $\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$ 
7:     end if
8:   end for
9: end for
10: return  $\mathbf{w}$ 
```

Prediction: $sign(\mathbf{w}^T \mathbf{x})$.

Design decisions

- *MaxIter*: Hyperparameter
- How to loop over the data?
 - ▶ Constant.
 - ▶ Permuting once
 - ▶ Permuting in each iteration

Properties of perceptron learning

- This is an **online** algorithm – looks at one instance at a time.
- Does the algorithm terminate (**convergence**)?

Properties of perceptron learning

- This is an **online** algorithm – looks at one instance at a time.
- Does the algorithm terminate (**convergence**)?
 - ▶ If training data is **not linearly separable**, the algorithm does not converge.
 - ▶ If the training data is linearly separable, the algorithm stops in a finite number of steps (**converges**).

Properties of perceptron learning

- This is an **online** algorithm – looks at one instance at a time.
- Does the algorithm terminate (**convergence**)?
 - ▶ If training data is **not linearly separable**, the algorithm does not converge.
 - ▶ If the training data is linearly separable, the algorithm stops in a finite number of steps (**converges**).
- How long to convergence ?
 - ▶ Depends on the difficulty of the problem (**margin**).

Outline

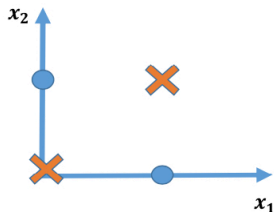
- 1 Perceptron
- 2 Perceptron learning
- 3 Convergence of the Perceptron Learning Algorithm
- 4 Variants of perceptron
- 5 What we have learned

Perceptron learnability

Perceptron cannot learn what it cannot represent

- Only linearly separable functions (Minsky and Papaert 1969).
- Parity function (XOR) cannot be learned.

x_1	x_2	y
0	0	0
1	0	1
0	1	1
1	1	0



Convergence

Convergence theorem

- If the data is linearly separable, the perceptron algorithm will converge after making mistakes that depend on the difficulty of the problem (margin).

Cycling theorem

- If the training data is not linearly separable, then the learning algorithm will eventually repeat the same set of weights and enter an infinite loop

Margin

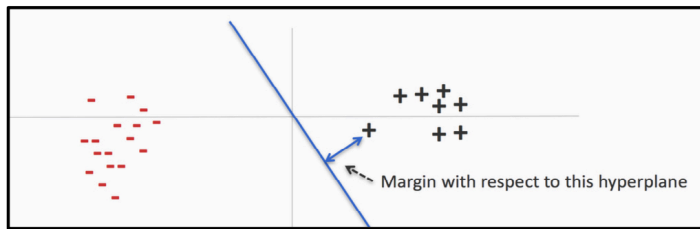
- The margin of a separating hyperplane for a dataset is the distance between the hyperplane and the data point nearest to it.
- The margin of a data set is the maximum margin possible for that dataset using any weight vector.

Margin

$$\text{Margin}(\mathcal{D}, \mathbf{w}) = \begin{cases} \min_{(\mathbf{x}, y) \in \mathcal{D}} y \mathbf{w}^T \mathbf{x} & \text{for a separating hyperplane } \mathbf{w} \\ -\infty & \text{else} \end{cases}$$

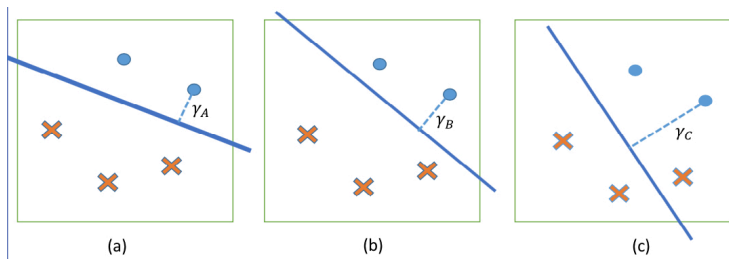
$$\text{Margin}(\mathcal{D}) = \sup_{\mathbf{w}} \text{Margin}(\mathcal{D}, \mathbf{w})$$

Margin



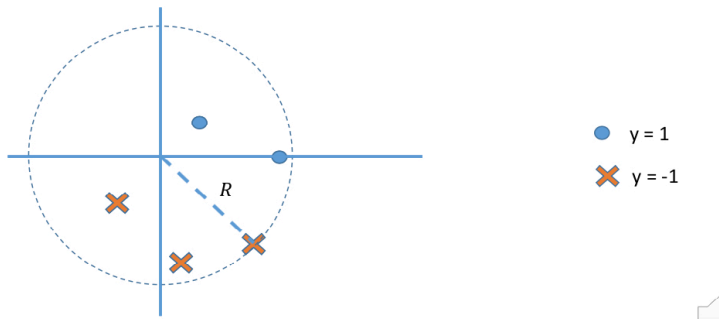
Margin

Which γ is the margin of the data?



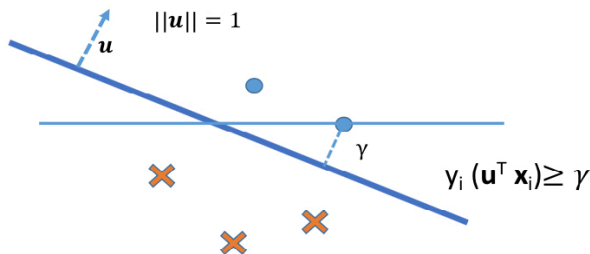
The Mistake Bound Theorem (Novikoff 1962, Block 1962)

- Let $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ be a sequence of training examples such that $\|\mathbf{x}_n\|_2 \leq R$ and label $y_n \in \{-1, +1\}$.



The Mistake Bound Theorem (Novikoff 1962, Block 1962)

- Let $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ be a sequence of training examples such that $\|\mathbf{x}_n\|_2 \leq R$ and label $y_n \in \{-1, +1\}$.
- Suppose there exists a unit vector $\mathbf{u} \in \mathbb{R}^D$ such that for some $\gamma > 0$, we have $y_n \mathbf{u}^T \mathbf{x}_n \geq \gamma$.



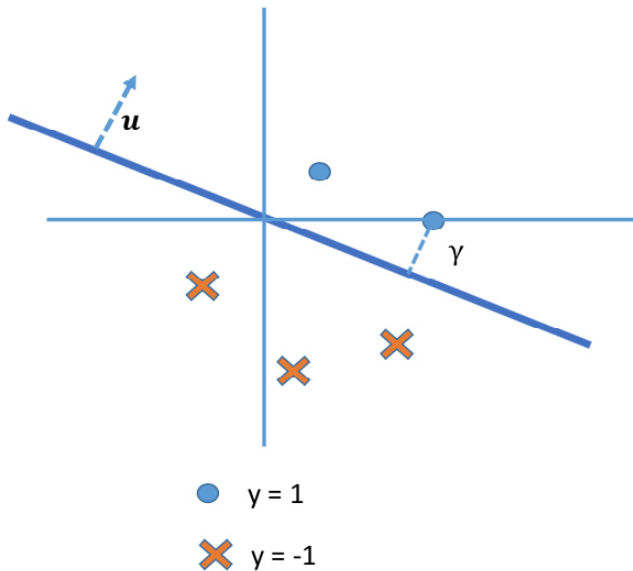
The Mistake Bound Theorem (Novikoff 1962, Block 1962)

- Let $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ be a sequence of training examples such that $\|\mathbf{x}_n\|_2 \leq R$ and label $y_n \in \{-1, +1\}$.
- Suppose there exists a unit vector $\mathbf{u} \in \mathbb{R}^D$ such that for some $\gamma > 0$, we have $y_n \mathbf{u}^T \mathbf{x}_n \geq \gamma$.
- Then the Perceptron algorithm will make at most $\frac{R^2}{\gamma^2}$ mistakes on the training sequence.

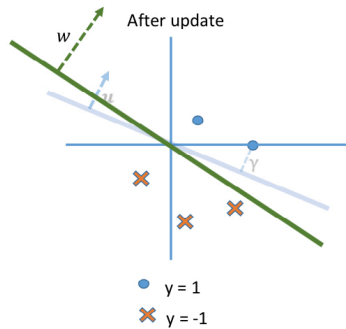
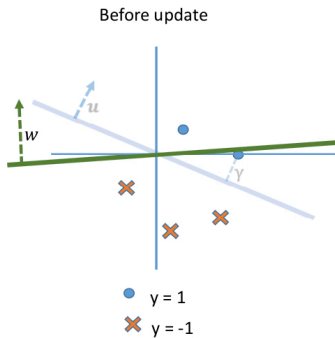
The Mistake Bound Theorem (Novikoff 1962, Block 1962)

- If the data is separable....
- then the perceptron algorithm will find a separating hyperplane after making a finite number of mistakes.

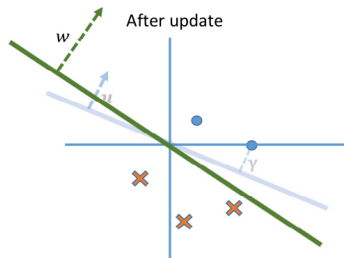
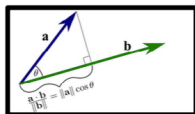
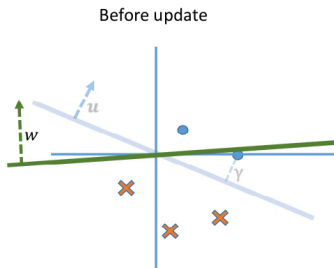
Intuition



Intuition

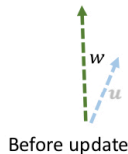
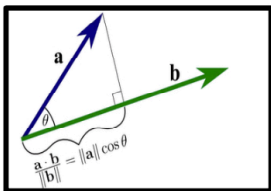


Intuition



Intuition

- After update, $\mathbf{u}^\top \mathbf{w}_{t+1}$ is larger than $\mathbf{u}^\top \mathbf{w}_t$.
 - ▶ After t mistakes, $\mathbf{u}^\top \mathbf{w}_t \geq t\gamma$.
- The size of $\|\mathbf{w}_{t+1}\|$ may increase but not too much.
 - ▶ After t mistakes, $\|\mathbf{w}_t\|^2 \leq tR^2$.



Proof (Preliminaries)

Setting

- Initial weight vector $\mathbf{w}_0 = \mathbf{0}$.
- All training examples are contained in a ball of size R .
 $\|\mathbf{x}_n\| \leq R$.
- The training data is separable by a margin γ using a unit vector \mathbf{u} .
 $y_n \mathbf{u}^T \mathbf{x}_n \geq \gamma$.

Proof (1/3)

Claim 1: After t mistakes, $\mathbf{u}^T \mathbf{w}_t \geq t\gamma$.

$$\begin{aligned}\mathbf{u}^T \mathbf{w}_{t+1} &= \mathbf{u}^T (\mathbf{w}_t + y_n \mathbf{x}_n) \\ &\geq \mathbf{u}^T \mathbf{w}_t + \gamma\end{aligned}$$

Because $\mathbf{w}_0 = 0$, simple induction gives us: $\mathbf{u}^T \mathbf{w}_t \geq t\gamma$.

```
a ← wTx
if ay ≤ 0 then
  w ← w + yx
```

Proof (1/3)

Claim 1: After t mistakes, $\mathbf{u}^T \mathbf{w}_t \geq t\gamma$.

$$\begin{aligned}\mathbf{u}^T \mathbf{w}_{t+1} &= \mathbf{u}^T (\mathbf{w}_t + y_n \mathbf{x}_n) \\ &\geq \mathbf{u}^T \mathbf{w}_t + \gamma\end{aligned}$$

Because $\mathbf{w}_0 = 0$, simple induction gives us: $\mathbf{u}^T \mathbf{w}_t \geq t\gamma$.

```
a ← wTx
if ay ≤ 0 then
  w ← w + yx
```

The inner product between the true underlying model and the current model is non-decreasing after each update. **This could be because the directions of \mathbf{w} and \mathbf{u} align or because the length of \mathbf{w} increases.**

Proof (2/3)

Claim 2: After t mistakes, $\|\mathbf{w}_t\|^2 \leq tR^2$

$$\begin{aligned}\|\mathbf{w}_{t+1}\|^2 &= \|\mathbf{w}_t + y_n \mathbf{x}_n\|^2 \\ &= \|\mathbf{w}_t\|^2 + 2\mathbf{w}_t^\top (y_n \mathbf{x}_n) + \|y_n \mathbf{x}_n\|^2\end{aligned}$$

Proof (2/3)

Claim 2: After t mistakes, $\|\mathbf{w}_t\|^2 \leq tR^2$

$$\begin{aligned}\|\mathbf{w}_{t+1}\|^2 &= \|\mathbf{w}_t + y_n \mathbf{x}_n\|^2 \\ &= \|\mathbf{w}_t\|^2 + 2\mathbf{w}_t^\top (y_n \mathbf{x}_n) + \|y_n \mathbf{x}_n\|^2 \\ &= \|\mathbf{w}_t\|^2 + 2y_n (\mathbf{w}_t^\top \mathbf{x}_n) + y_n^2 \|\mathbf{x}_n\|^2\end{aligned}$$

```
a ← wTx
if ay ≤ 0 then
  w ← w + yx
```

Proof (2/3)

Claim 2: After t mistakes, $\|\mathbf{w}_t\|^2 \leq tR^2$

$$\begin{aligned}\|\mathbf{w}_{t+1}\|^2 &= \|\mathbf{w}_t + y_n \mathbf{x}_n\|^2 \\ &= \|\mathbf{w}_t\|^2 + 2\mathbf{w}_t^\top (y_n \mathbf{x}_n) + \|y_n \mathbf{x}_n\|^2 \\ &\leq \|\mathbf{w}_t\|^2 + R^2\end{aligned}$$

Because $\mathbf{w}_0 = 0$, simple induction gives us: $\|\mathbf{w}_t\|^2 \leq tR^2$.

Proof (3/3)

What we know

- ① After t mistakes, $\mathbf{u}^T \mathbf{w}_t \geq t\gamma$.
- ② After t mistakes, $\|\mathbf{w}_t\|^2 \leq tR^2$

The inner product between the true underlying model and the current model is non-decreasing after each update. This could be because the directions of \mathbf{w} and \mathbf{u} align or because the length of \mathbf{w} increases.

But the length of \mathbf{w} does not increase too much!.

Proof (3/3)

What we know

- ① After t mistakes, $\mathbf{u}^T \mathbf{w}_t \geq t\gamma$.
- ② After t mistakes, $\|\mathbf{w}_t\|^2 \leq tR^2$

$$R\sqrt{t} \geq \|\mathbf{w}_t\|$$

Proof (3/3)

What we know

- ① After t mistakes, $\mathbf{u}^T \mathbf{w}_t \geq t\gamma$.
- ② After t mistakes, $\|\mathbf{w}_t\|^2 \leq tR^2$

$$R\sqrt{t} \geq \|\mathbf{w}_t\| \geq \mathbf{u}^T \mathbf{w}_t$$

Proof (3/3)

What we know

- ① After t mistakes, $\mathbf{u}^T \mathbf{w}_t \geq t\gamma$.
- ② After t mistakes, $\|\mathbf{w}_t\|^2 \leq tR^2$

$$R\sqrt{t} \geq \|\mathbf{w}_t\| \geq \mathbf{u}^T \mathbf{w}_t \geq t\gamma$$

Proof (3/3)

What we know

- ① After t mistakes, $\mathbf{u}^T \mathbf{w}_t \geq t\gamma$.
- ② After t mistakes, $\|\mathbf{w}_t\|^2 \leq tR^2$

$$R\sqrt{t} \geq \|\mathbf{w}_t\| \geq \mathbf{u}^T \mathbf{w}_t \geq t\gamma$$

Number of mistakes $t \leq \frac{R^2}{\gamma^2}$.

Bounds the total number of mistakes!

Beyond the separable case

- Good news
 - ▶ Perceptron makes no assumptions about the data, could be even adversarial.
 - ▶ After a fixed number of mistakes, you are done. Do not need to see any more data.
- Bad news
 - ▶ Real world data is often not linearly separable.

Outline

- 1 Perceptron
- 2 Perceptron learning
- 3 Convergence of the Perceptron Learning Algorithm
- 4 Variants of perceptron**
- 5 What we have learned

Voting and averaging

- Vanilla perceptron returns final weight vector.
- Might lose good weight vectors that were learned during training.
- Aggregating the models (or weight vectors) seen during training may give better results (especially when data is not separable).

Voted perceptron

- Remember every weight vector in your sequence of updates.
- At final prediction time, each weight vector gets to vote on the label.
- The number of votes it gets is the number of iterations it survived before being updated
- Comes with strong theoretical guarantees about generalization, impractical because of storage issues

Averaged perceptron

- Instead of using all weight vectors, use the average weight vector (i.e longer surviving weight vectors get more say)
- More practical alternative and widely used

Averaged Perceptron training/learning

data = N **samples/instances**: $= \{(x_1, y_1), \dots, (x_N, y_N)\}$

Algorithm 2 AveragedPerceptronTrain (*data*, *maxIter*)

```
1:  $w \leftarrow 0$ .  $\mu \leftarrow 0$ .  
2: for  $iter = 1 \dots MaxIter$  do  
3:   for  $(x, y) \in data$  do  
4:      $a \leftarrow w^T x$   
5:     if  $ay \leq 0$  then  
6:        $w \leftarrow w + yx$   
7:     end if  
8:      $\mu \leftarrow \mu + w$   
9:   end for  
10: end for  
11: return  $\mu$ 
```

Prediction: $sign(\mu^T x)$.

Perceptron

- Extensions
 - ▶ Voting
 - ▶ Averaging
- Limitations
 - ▶ Linear separability
- Interpreting the importance of features
 - ▶ The values of weight w_d tells us the importance of feature x_d .

Outline

- 1 Perceptron
- 2 Perceptron learning
- 3 Convergence of the Perceptron Learning Algorithm
- 4 Variants of perceptron
- 5 What we have learned

Summary

- You should now be able to understand the differences between decision trees, perceptrons and nearest neighbors.
- Given data, use training, development and test splits (or cross-validation).
- Use training and development to tune hyperparameters that trades off overfitting and underfitting.
- Use test to get an estimate of generalization or accuracy on unseen data.