# 10-301/601: Introduction to Machine Learning Lecture 6 – Perceptron

**Geoff Gordon** 

with thanks to Henry Chai & Matt Gormley

Q & A:

Can we use k-NN with categorical features?

#### Q & A:

Can we use k-NN with categorical features?

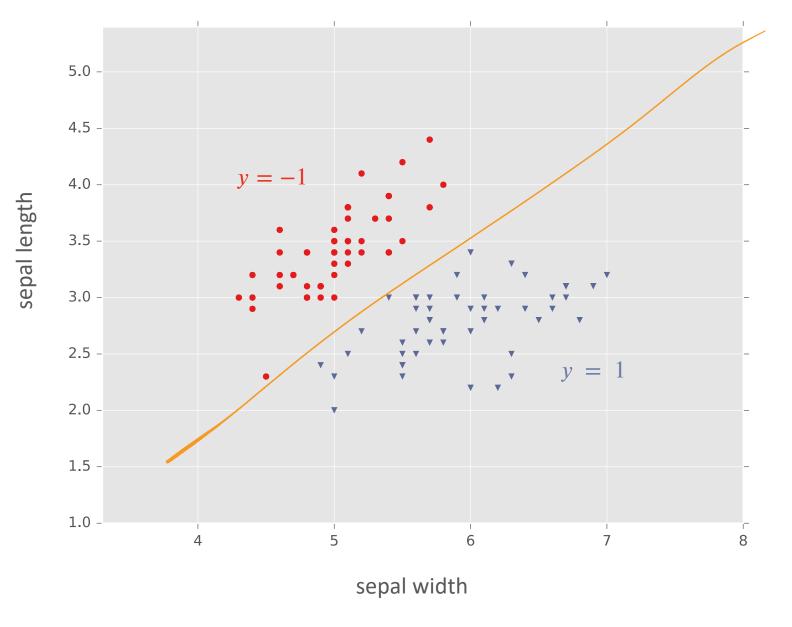
- Yes! We can either:
  - 1. Convert categorical features into binary ones:

erol		Abnormal Cholesterol?
al	 I	0
	1	0
	0	I

Use a distance metric that works over categorical features, e.g., the Hamming distance:  $d(x,x') = \sum_{i=1}^{M} \mathbb{I}(x_i \neq x_i')$ 

$$d(x, x') = \sum_{i=1}^{M} \mathbb{I}(x_i \neq x_i')$$

#### Fisher Iris Dataset



#### Linear Algebra Review

Notation: in this we use column vectors by default, i.e.,

$$\boldsymbol{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_D \end{bmatrix} \text{ and } \boldsymbol{a}^T = \begin{bmatrix} a_1 & a_2 & \cdots & a_D \end{bmatrix}$$

The dot product between two D-dimensional vectors is

$$\boldsymbol{a}^T \boldsymbol{b} = \begin{bmatrix} a_1 & a_2 & \cdots & a_D \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_D \end{bmatrix} = \sum_{d=1}^D a_d b_d$$

#### Linear Algebra Review

- Dot products represent linear functions:  $\mathbf{w}^{\mathsf{T}}\mathbf{x}$  is linear in  $\mathbf{x}$  for fixed  $\mathbf{w}$  (and all linear functions of  $\mathbf{x}$  that pass through origin can be written this way)
  - •include intercept  $w^{\top}x + b$  to not hit origin
- The L2-norm of  $\boldsymbol{a} = \|\boldsymbol{a}\|_2 = \sqrt{\boldsymbol{a}^T \boldsymbol{a}}$
- Two vectors are *orthogonal* iff  $\mathbf{a}^T \mathbf{b} = 0$

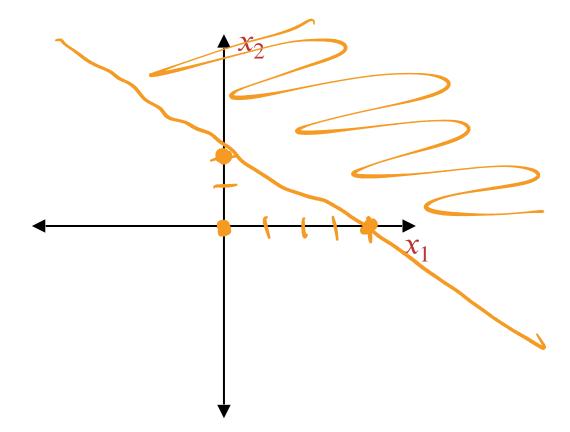
#### Geometry Warm-up

1. On the axes below, draw the region corresponding to

$$w_1 x_1 + w_2 x_2 + b > 0$$

where  $w_1 = 1$ ,  $w_2 = 2$  and b = -4.

2. Then draw the vector 
$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$



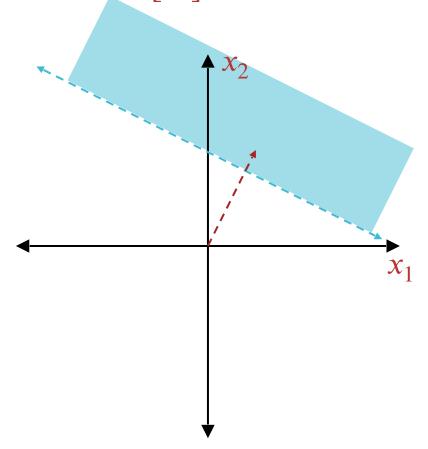
#### Geometry Warm-up

1. On the axes below, draw the region corresponding to

$$w_1 x_1 + w_2 x_2 + b > 0$$

where  $w_1 = 1$ ,  $w_2 = 2$  and b = -4.

2. Then draw the vector  $\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$ 



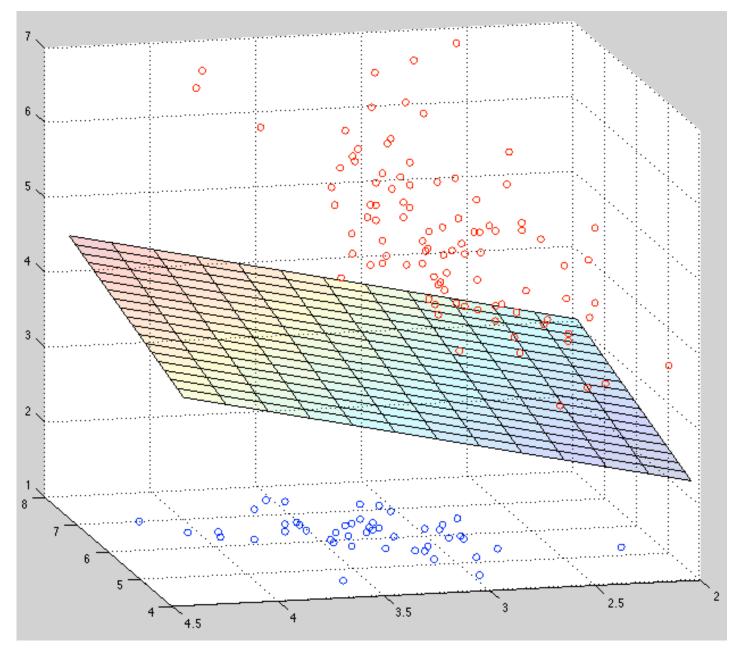
### Linear Decision Boundaries

- In 2 dimensions,  $w_1x_1 + w_2x_2 + b = 0$  defines a *line*
- •In 3 dimensions,  $w_1x_1 + w_2x_2 + w_3x_3 + b = 0$  defines a plane
- •In 4+ dimensions,  $\mathbf{w}^T \mathbf{x} + \mathbf{b} = 0$  defines a hyperplane
  - The vector  $\boldsymbol{w}$  is always orthogonal to this hyperplane and always points in the direction where  $\boldsymbol{w}^T\boldsymbol{x} + b > 0$ !
- A hyperplane creates two halfspaces:

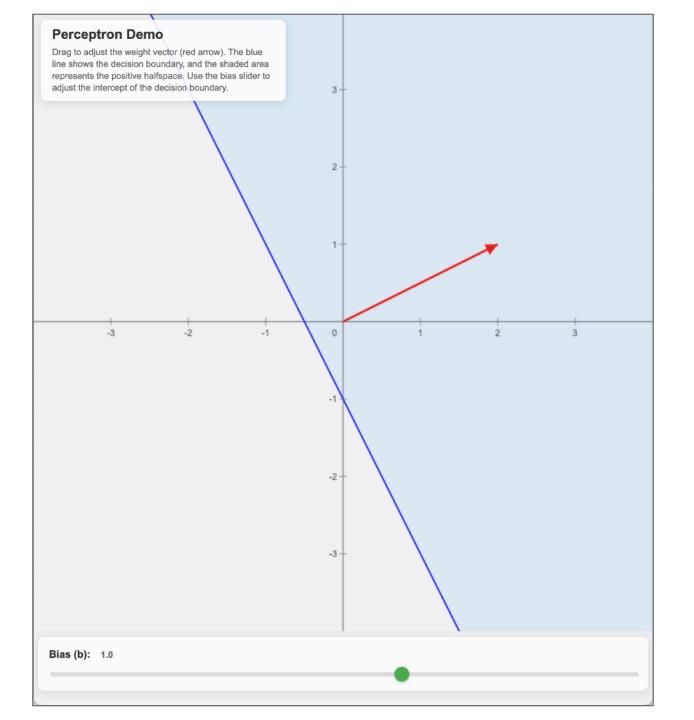
• 
$$\mathcal{S}_+ = \{ \mathbf{x} \colon \mathbf{w}^T \mathbf{x} + b > 0 \}$$
 or all  $\mathbf{x}$  s.t.  $\mathbf{w}^T \mathbf{x} + b$  is positive

$$\bullet S_- = \{ \mathbf{x} : \mathbf{w}^T \mathbf{x} + b < 0 \}$$
 or all  $\mathbf{x}$  s.t.  $\mathbf{w}^T \mathbf{x} + b$  is negative

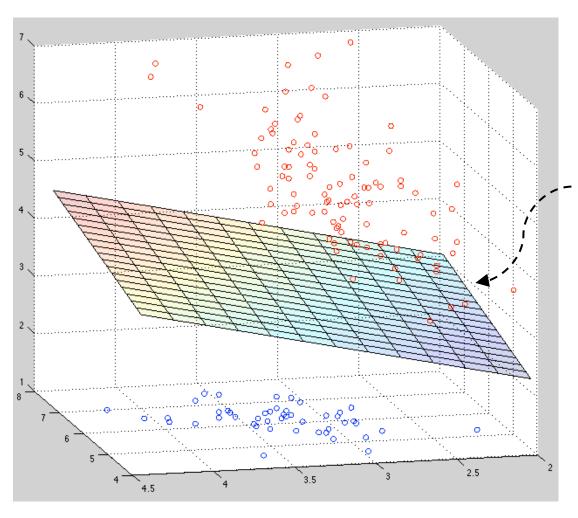
## Linear Decision Boundaries: Example



## Interactive decision boundary



## Learning a linear classifier

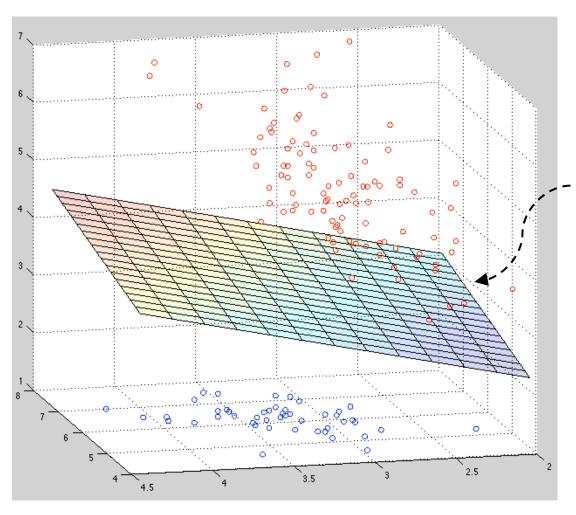


Goal: learn classifiers of the form

$$h(x) = sign(w^T x + b)$$
(assuming

$$y \in \{-1, +1\}$$

### Learning a linear classifier



Goal: learn classifiers of the form

$$h(x) = \operatorname{sign}(w^T x + b)$$
(assuming

$$y \in \{-1, +1\}$$

Key question: how do we learn the parameters, **w**, b?

#### Online Learning

- So far, we've been learning in the batch setting, where we have access to the entire training dataset at once
- •A common alternative is the *online* setting, where examples arrive gradually and we learn continuously
- Examples of online learning:

autocorrect — personalization recommender system performance turing stock prediction

## Online Learning: Setup

- For t = 1, 2, 3, ...
  - Receive an unlabeled example,  $\mathbf{x}^{(t)}$
  - Predict its label,  $\hat{y} = h_{w,b}(x^{(t)})$
  - Observe its true label,  $y^{(t)}$
  - Pay a penalty if we made a mistake,  $\hat{y} \neq y^{(t)}$
  - Update the parameters, w and b

Goal: minimize the number of mistakes made

# (Online) Perceptron Learning Algorithm

Initialize the weight vector and intercept to all zeros:

$$w = [0 \ 0 \ \cdots \ 0] \text{ and } b = 0$$

- For t = 1, 2, 3, ...
  - Receive an unlabeled example,  $\mathbf{x}^{(t)}$

Predict its label, 
$$\hat{y} = \text{sign}(\boldsymbol{w}^T \boldsymbol{x} + b) = \begin{cases} +1 \text{ if } \boldsymbol{w}^T \boldsymbol{x} + b \ge 0 \\ -1 \text{ otherwise} \end{cases}$$

- Observe its true label,  $y^{(t)}$
- If we misclassified a positive example  $(y^{(t)} = +1, \hat{y} = -1)$ :

$$\boldsymbol{w} \leftarrow \boldsymbol{w} + \boldsymbol{x}^{(t)}$$

$$b \leftarrow b + 1$$

• If we misclassified a negative example  $(y^{(t)} = -1, \hat{y} = +1)$ :

$$\bullet w \leftarrow w - x^{(t)}$$

$$\bullet b \leftarrow b - 1$$

#### Notational hack

Initialize the weight vector and intercept to all zeros:

$$w = [0 \ 0 \ \cdots \ 0] \text{ and } b = 0$$

- For t = 1, 2, 3, ...
  - Receive an unlabeled example,  $\mathbf{x}^{(t)}$

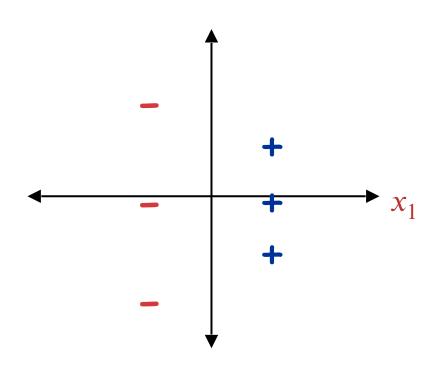
Predict its label, 
$$\hat{y} = \text{sign}(\boldsymbol{w}^T \boldsymbol{x} + b) = \begin{cases} +1 \text{ if } \boldsymbol{w}^T \boldsymbol{x} + b \ge 0 \\ -1 \text{ otherwise} \end{cases}$$

- Observe its true label,  $y^{(t)}$
- If we misclassified an example  $(y^{(t)} \neq \hat{y})$ :

$$\bullet w \leftarrow w + y^{(t)} x^{(t)}$$

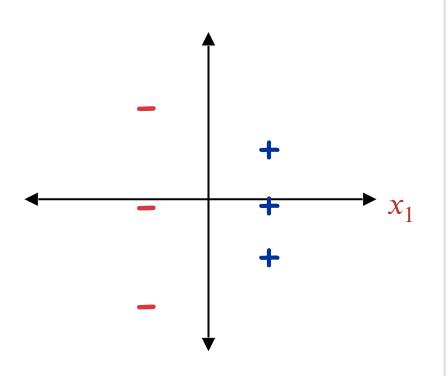
$$b \leftarrow b + y^{(t)}$$

$x_1$	$x_2$	ŷ	у	Mistake?
-1	2	+	_	Yes



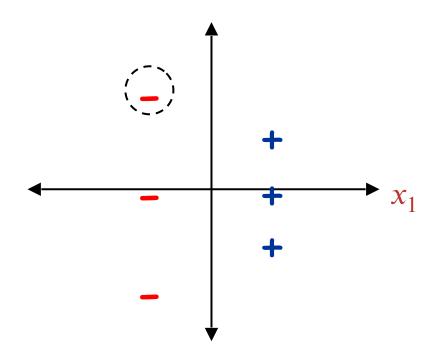
$x_1$	$x_2$	ŷ	y	Mistake?
-	2	+	_	Yes

$$\boldsymbol{w} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

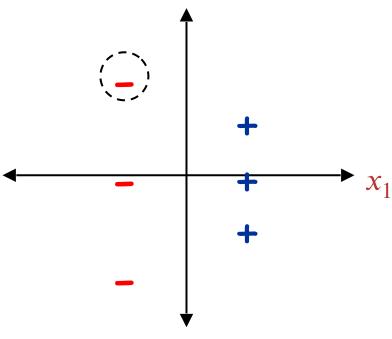


$x_1$	$x_2$	ŷ	y	Mistake?
-	2	+	_	Yes

$$\boldsymbol{w} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$



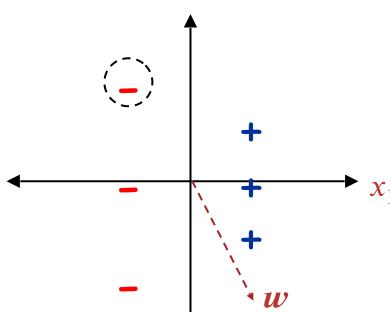
$x_1$	$x_2$	ŷ	у	Mistake?
-1	2	+	_	Yes



 $\chi_2$ 

$$\mathbf{w} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$
$$\mathbf{w} \leftarrow \mathbf{w} + y^{(1)}\mathbf{x}^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} -1 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

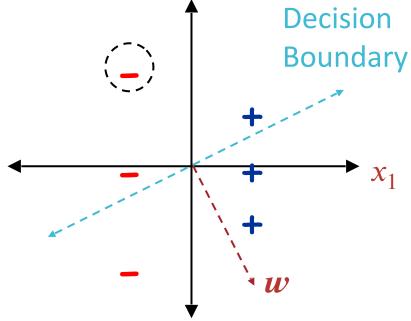
$x_1$	$x_2$	ŷ	y	Mistake?
-1	2	+	_	Yes



 $\chi_2$ 

$$\boldsymbol{w} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$
$$\boldsymbol{w} \leftarrow \boldsymbol{w} + y^{(1)} \boldsymbol{x}^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} -1 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

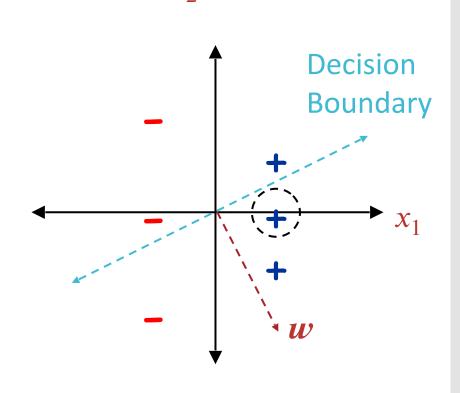
$x_1$	$x_2$	ŷ	y	Mistake?
-1	2	+	_	Yes



$$\mathbf{w} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$
$$\mathbf{w} \leftarrow \mathbf{w} + y^{(1)} \mathbf{x}^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} -1 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

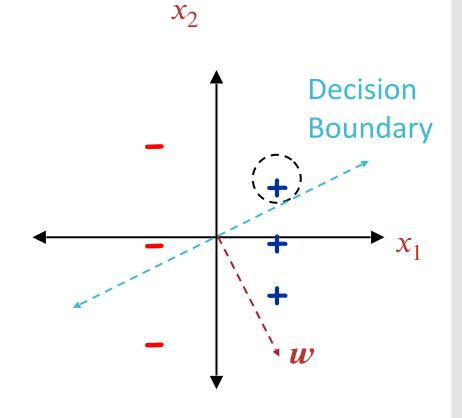
$x_1$	$x_2$	ŷ	y	Mistake?
-1	2	+	_	Yes
	0	+	+	No

$$\boldsymbol{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

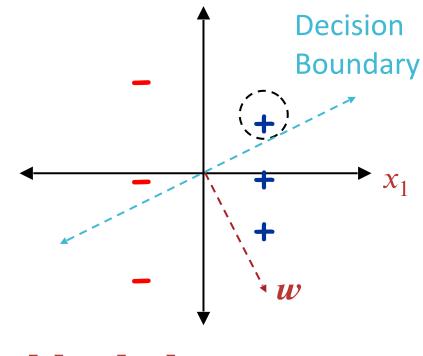


$x_1$	$x_2$	ŷ	y	Mistake?
-	2	+	_	Yes
	0	+	+	No
	I	_	+	Yes

$$\boldsymbol{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

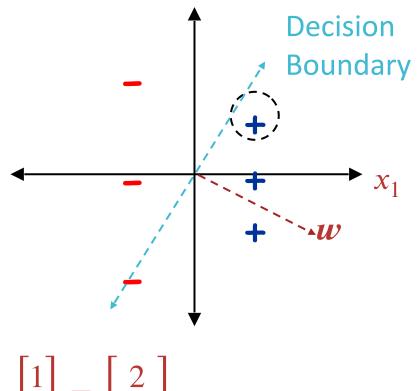


$x_1$	$x_2$	ŷ	y	Mistake?
-1	2	+	_	Yes
1	0	+	+	No
- 1	I	_	+	Yes



$$\boldsymbol{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$
$$\boldsymbol{w} \leftarrow \boldsymbol{w} + y^{(3)} \boldsymbol{x}^{(3)} = \begin{bmatrix} 1 \\ -2 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$

$x_1$	$x_2$	ŷ	y	Mistake?
-1	2	+	_	Yes
1	0	+	+	No
1	l	_	+	Yes

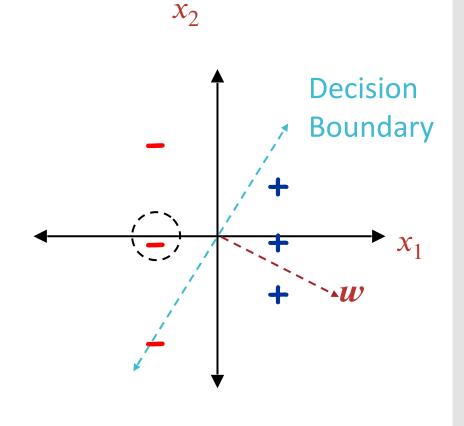


$$\boldsymbol{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

$$\boldsymbol{w} \leftarrow \boldsymbol{w} + y^{(3)} \boldsymbol{x}^{(3)} = \begin{bmatrix} 1 \\ -2 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$

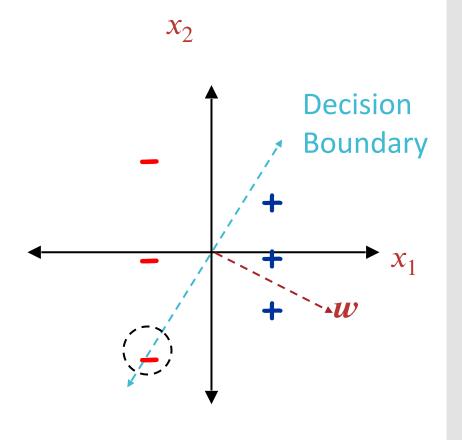
$x_1$	$x_2$	ŷ	y	Mistake?
-1	2	+	_	Yes
I	0	+	+	No
I	I	_	+	Yes
-1	0	_	_	No

$$\boldsymbol{w} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$

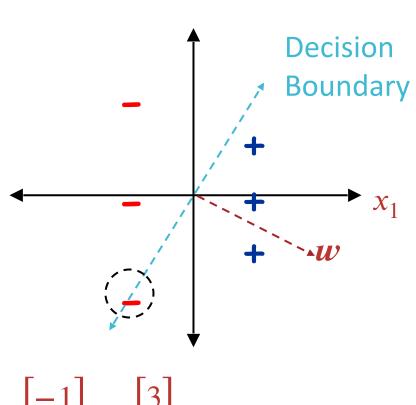


$x_1$	$x_2$	ŷ	у	Mistake?
-1	2	+	_	Yes
- 1	0	+	+	No
- 1	I	_	+	Yes
-1	0	_	_	No
-1	-2	+	_	Yes

$$\boldsymbol{w} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$



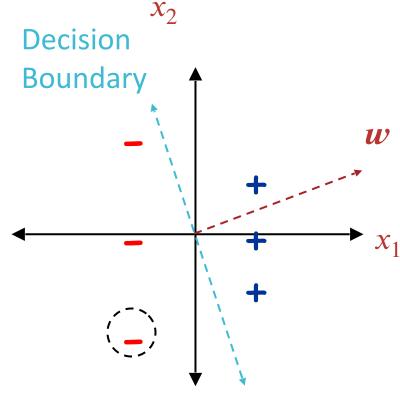
$x_1$	$x_2$	ŷ	у	Mistake?
<b>–</b> I	2	+	_	Yes
1	0	+	+	No
1		_	+	Yes
-1	0	_	_	No
-1	-2	+	_	Yes



$$\mathbf{w} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$

$$\mathbf{w} \leftarrow \mathbf{w} + y^{(5)} \mathbf{x}^{(5)} = \begin{bmatrix} 2 \\ -1 \end{bmatrix} - \begin{bmatrix} -1 \\ -2 \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$$

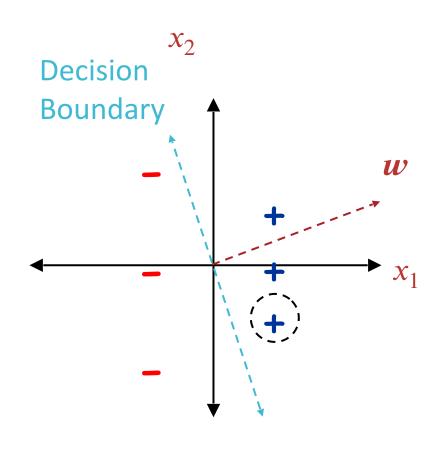
$x_1$	$x_2$	ŷ	y	Mistake?
-1	2	+	_	Yes
I	0	+	+	No
I	I	_	+	Yes
<b>–</b> I	0	_	_	No
<b>—</b>	-2	+	_	Yes



$$\boldsymbol{w} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$
$$\boldsymbol{w} \leftarrow \boldsymbol{w} + y^{(5)} \boldsymbol{x}^{(5)} = \begin{bmatrix} 2 \\ -1 \end{bmatrix} - \begin{bmatrix} -1 \\ -2 \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$$

$x_1$	$x_2$	ŷ	y	Mistake?
-	2	+	_	Yes
I	0	+	+	No
I	I	_	+	Yes
-1	0	_	_	No
-1	-2	+	_	Yes
I	-	+	+	No

$$\boldsymbol{w} = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$$



## Thought experiment

- [No bias term for now]
- What if we scaled up every training example?

•
$$x^{(i)} \rightarrow 2x^{(i)}$$
 for all  $i$ 



## Thought experiment

• What if we scaled up *half* of the training examples?

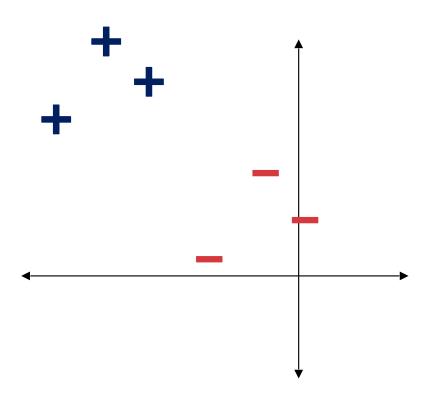
 $ullet x^{(i)} o 2x^{(i)}$  for every **even** value of i

Scale a feature Societ a feature S

Laure important

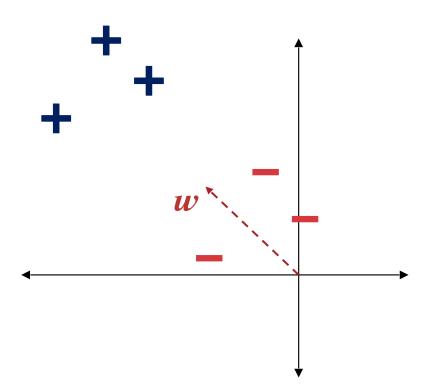
## Updating the Intercept

 The intercept shifts the decision boundary off the origin

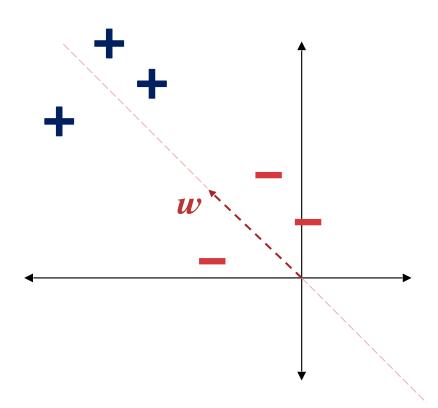


## Updating the Intercept

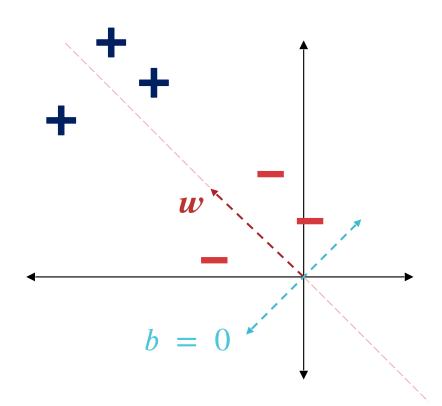
 The intercept shifts the decision boundary off the origin



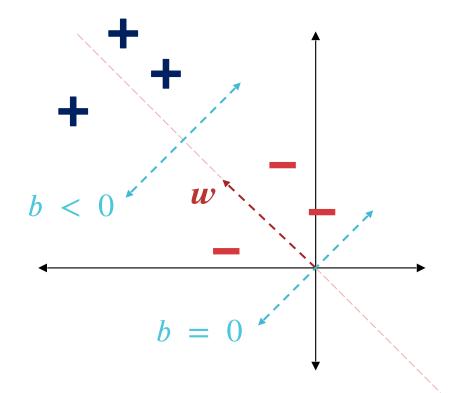
 The intercept shifts the decision boundary off the origin



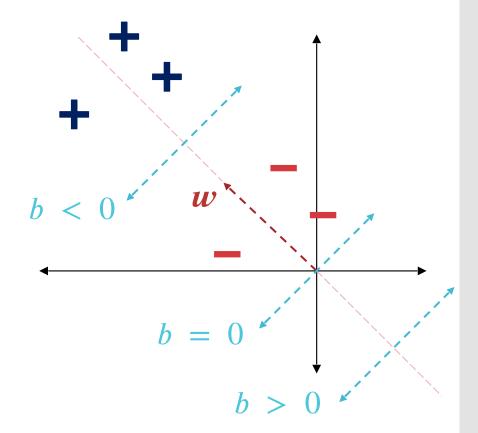
 The intercept shifts the decision boundary off the origin



- The intercept shifts the decision boundary off the origin
  - Increasing b shifts
     the decision
     boundary towards
     the negative side



- The intercept shifts the decision boundary off the origin
  - Increasing b shifts
     the decision
     boundary towards
     the negative side
  - Decreasing b shifts
     the decision
     boundary towards
     the positive side



#### Poll Question 1

poll. m/couse.og

• **True or False**: Unlike Decision Trees and k-Nearest Neighbors, the Perceptron learning algorithm does not suffer from overfitting because it does not have any hyperparameters that could be over-tuned on a validation dataset.

A: Falk C: The

B: fox. C

### Notational hack #2

• If we add a 1 to the beginning of every example e.g.,

$$\mathbf{x}' = \begin{vmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_D \end{vmatrix}$$

• ... we can just fold the intercept into the weight vector!

$$\theta = \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix} \rightarrow \theta^T \mathbf{x}' = \mathbf{w}^T \mathbf{x} + b$$

### Notational hack #2

•If we add a 1 to the beginning of every example e.g.,

$$\mathbf{x}' = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}$$

• ... we can just fold the intercept into the weight vector!

$$\theta = \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix} \rightarrow \theta^T \mathbf{x}' = \mathbf{w}^T \mathbf{x} + b$$

$$\vdots$$
and powerfaces

... and now if we update  $\theta \to \theta + yx'$ , that's the same as  $w \to w + yx$  and  $b \to b + y$ 

# (Online) Perceptron Learning Algorithm

• Initialize the parameters to all zeros:

$$\theta = [0 \quad 0 \quad \cdots \quad 0]$$

- For t = 1, 2, 3, ...
  - Receive an unlabeled example,  $x'^{(t)}$

Predict its label, 
$$\hat{y} = \operatorname{sign}(\theta^T x^{\prime(t)}) = \begin{cases} +1 \text{ if } \theta^T x^{\prime(t)} \ge 0 \\ -1 \text{ otherwise} \end{cases}$$

- Observe its true label,  $y^{(t)}$
- If we misclassified an example  $(y^{(t)} \neq \hat{y})$ :

$$\bullet \theta \leftarrow \theta + y^{(t)} x'^{(t)}$$

# (Online) Perceptron Learning Algorithm

• Initialize the parameters to all zeros:

$$\theta = \begin{bmatrix} 0 & 0 & \cdots & 0 \end{bmatrix}$$

- For t = 1, 2, 3, ...
  - Receive an unlabeled example,  $x'^{(t)} \leftarrow ---$

Predict its label, 
$$\hat{y} = \text{sign}(\theta^T x^{\prime(t)}) = \begin{cases} +1 \text{ if } \theta^T x^{\prime(t)} \ge 0 \\ -1 \text{ otherwise} \end{cases}$$

- Observe its true label,  $y^{(t)}$
- If we misclassified an example  $(y^{(t)} \neq \hat{y})$ :

$$\bullet \theta \leftarrow \theta + y^{(t)} x'^{(t)}$$

1 prepended to  $\mathbf{x}^{(t)}$ 

# (Online) Perceptron Learning Algorithm

• Initialize the parameters to all zeros:

$$\theta = [0 \quad 0 \quad \cdots \quad 0]$$

- For t = 1, 2, 3, ...
  - Receive an unlabeled example,  $x'^{(t)}$

Predict its label, 
$$\hat{y} = \text{sign}(\theta^T x^{\prime(t)}) = \begin{cases} +1 \text{ if } \theta^T x^{\prime(t)} \ge 0 \\ -1 \text{ otherwise} \end{cases}$$

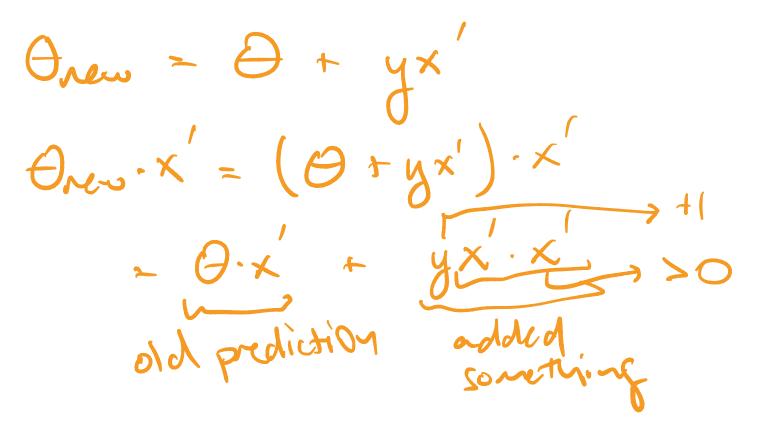
- Observe its true label,  $y^{(t)}$
- If we misclassified an example  $(y^{(t)} \neq \hat{y})$ :

$$\bullet \theta \leftarrow \theta + y^{(t)} x'^{(t)}$$

Automatically handles updating the intercept

(Online)
Perceptron
Learning
Algorithm:
Intuition

•Suppose  $(x, y) \in \mathcal{D}$  is a misclassified training example and y = +1 (the y = -1 case is similar)



(Online)
Perceptron
Learning
Algorithm:
Inductive Bias

don't change too fast linear decision bol. outlier sensitiel - higher scale ex continuos ests represent bd. important higher scale fontures mor insportant

# (Batch) Perceptron Learning Algorithm

•Input: 
$$\mathcal{D} = \{ (\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), ..., (\mathbf{x}^{(N)}, y^{(N)}) \}$$

• Initialize the parameters to all zeros:

$$\theta = [0 \quad 0 \quad \cdots \quad 0]$$

- While NOT CONVERGED
  - For  $t \in \{1,...,N\}$  [optionally: permute each epoch]
    - Predict the label of  $\mathbf{x}^{\prime(t)}$ ,  $\hat{\mathbf{y}} = \operatorname{sign}\left(\mathbf{\theta}^T \mathbf{x}^{\prime(t)}\right)$
    - •Observe its true label,  $y^{(t)}$
    - If we misclassified  $\mathbf{x}^{\prime(t)}$  ( $\mathbf{y}^{(t)} \neq \mathbf{\hat{y}}$ ):

$$\bullet \theta \leftarrow \theta + y^{(t)} x'^{(t)}$$

# (Batch) Perceptron Learning Algorithm

•Input: 
$$\mathcal{D} = \{ (\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), ..., (\mathbf{x}^{(N)}, y^{(N)}) \}$$

• Initialize the parameters to all zeros:

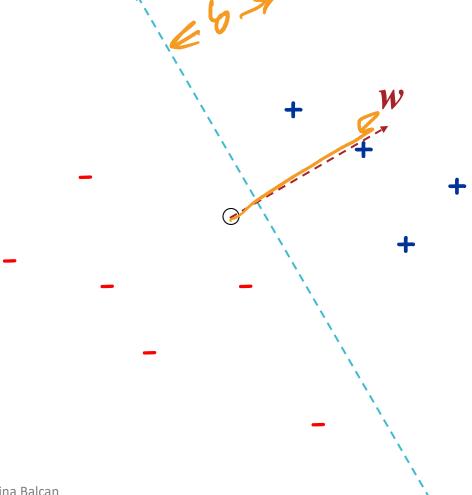
$$\theta = [0 \ 0 \ \cdots \ 0]$$

- While NOT CONVERGED ← what does this mean?
  - For  $t \in \{1,...,N\}$  [optionally: permute each epoch]
    - Predict the label of  $\mathbf{x}^{\prime(t)}$ ,  $\hat{\mathbf{y}} = \operatorname{sign}\left(\mathbf{\theta}^T \mathbf{x}^{\prime(t)}\right)$
    - •Observe its true label,  $y^{(t)}$
    - If we misclassified  $\mathbf{x}^{\prime(t)}$  ( $\mathbf{y}^{(t)} \neq \mathbf{\hat{y}}$ ):

$$\bullet \theta \leftarrow \theta + y^{(t)} x'^{(t)}$$

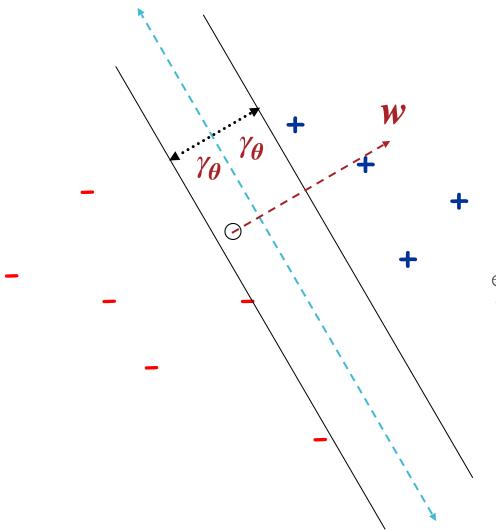
### Linearly separable

•A dataset  $\mathscr{D}$  is *linearly separable* if  $\exists$  a linear decision boundary  $\theta = (w, b)$  that perfectly classifies all examples in  $\mathscr{D}$ 



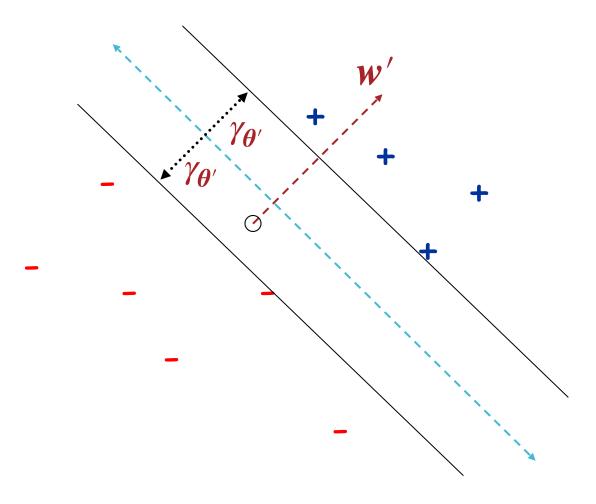
#### Margin

• The margin,  $\gamma_{\theta}$ , of any separator  $\theta$  is the distance of the closest example in  $\mathscr D$  to that separator



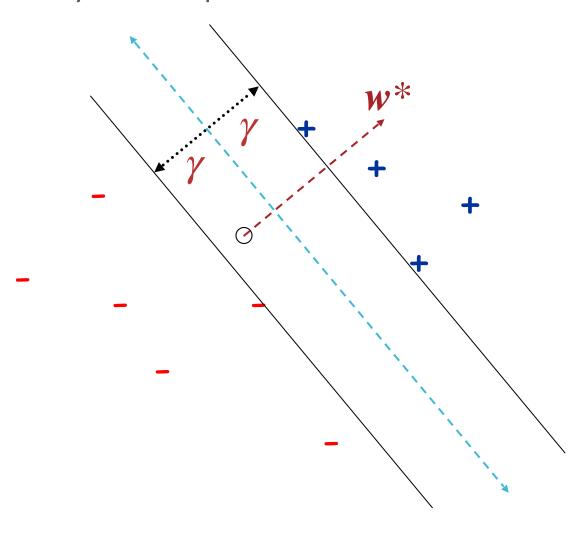
expand a slab around the separator until it touches an example on one side or the other Margin

Different separators can have different margins



#### Margin

•The margin,  $\gamma$ , of the entire dataset  $\mathscr D$  is the largest margin of any linear separator



#### Perceptron Mistake Bound

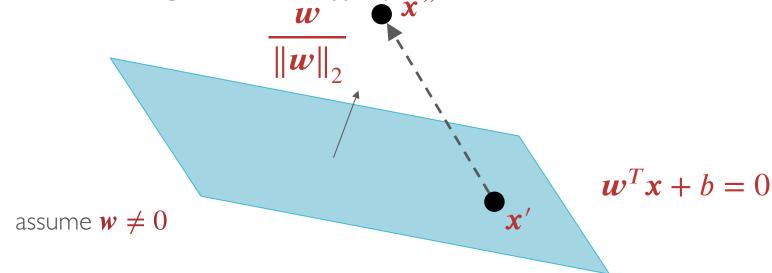
- Theorem: if the examples seen by the Perceptron Learning Algorithm (either online or batch)
  - 1. lie in a ball of radius R (centered around the origin)
  - 2. have a margin of  $\gamma > 0$

then the algorithm makes at most  $(R/\gamma)^2$  mistakes total!

• Key Takeaway: if the training dataset is linearly separable, the batch Perceptron Learning Algorithm will converge (i.e., stop making mistakes on the training dataset or achieve 0 training error) in a finite number of steps!

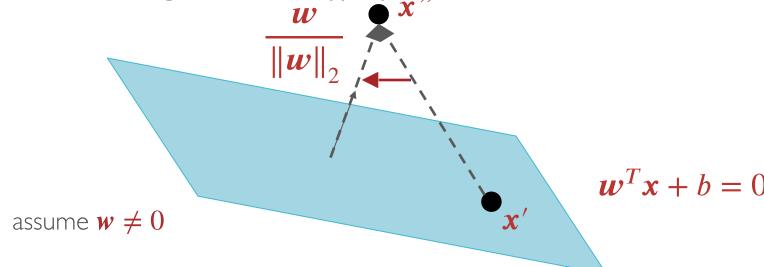
- •Let x' be an arbitrary point on the hyperplane  $w^Tx + b = 0$  and let x" be an arbitrary point
- •The distance between  $\mathbf{x}$  " and  $\mathbf{w}^T\mathbf{x} + b = 0$  is equal to the magnitude of the projection of  $\mathbf{x}$ "  $\mathbf{x}'$  onto  $\frac{\mathbf{w}}{\|\mathbf{w}\|_2}$ , the unit

vector orthogonal to the hyperplane



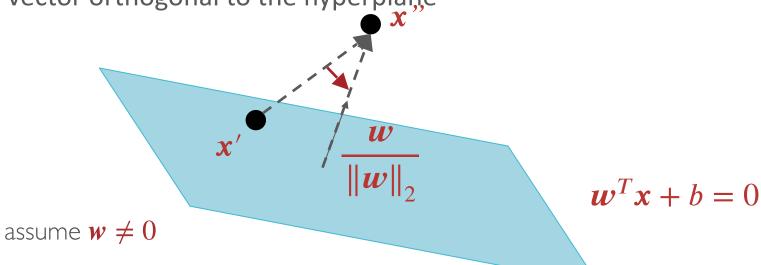
- •Let  $\mathbf{x}'$  be an arbitrary point on the hyperplane  $\mathbf{w}^T\mathbf{x} + b = 0$  and let  $\mathbf{x}$ " be an arbitrary point
- •The distance between x " and  $w^Tx + b = 0$  is equal to the magnitude of the projection of x" x' onto  $\frac{w}{\|w\|_2}$ , the unit

vector orthogonal to the hyperplane



- •Let  $\mathbf{x}'$  be an arbitrary point on the hyperplane  $\mathbf{w}^T\mathbf{x} + b = 0$  and let  $\mathbf{x}$ " be an arbitrary point
- The distance between x " and  $w^Tx + b = 0$  is equal to the magnitude of the projection of x " -x " onto  $\frac{w}{\|w\|_2}$ , the unit

vector orthogonal to the hyperplane



- •Let x' be an arbitrary point on the hyperplane and let x" be an arbitrary point
- The distance between x " and  $w^Tx + b = 0$  is equal to the magnitude of the projection of x " -x " onto  $\frac{w}{\|w\|_2}$ ,

the unit vector orthogonal to the hyperplane

$$\left| \frac{\boldsymbol{w}^{T}(\boldsymbol{x}" - \boldsymbol{x}')}{\|\boldsymbol{w}\|_{2}} \right| = \frac{\left| \boldsymbol{w}^{T}\boldsymbol{x}" - \boldsymbol{w}^{T}\boldsymbol{x}' \right|}{\|\boldsymbol{w}\|_{2}} = \frac{\left| \boldsymbol{w}^{T}\boldsymbol{x}" + b \right|}{\|\boldsymbol{w}\|_{2}}$$

#### Perceptron Learning Objectives

You should be able to...

- Explain the difference between online learning and batch learning
- Implement the perceptron algorithm for binary classification [CIML]
- Determine whether the perceptron algorithm will converge based on properties of the dataset, and the limitations of the convergence guarantees
- Describe the inductive bias of perceptron and the limitations of linear models
- Draw the decision boundary of a linear model
- Identify whether a dataset is linearly separable or not
- Defend the use of a bias term in perceptron (shifting points after projection onto weight vector)