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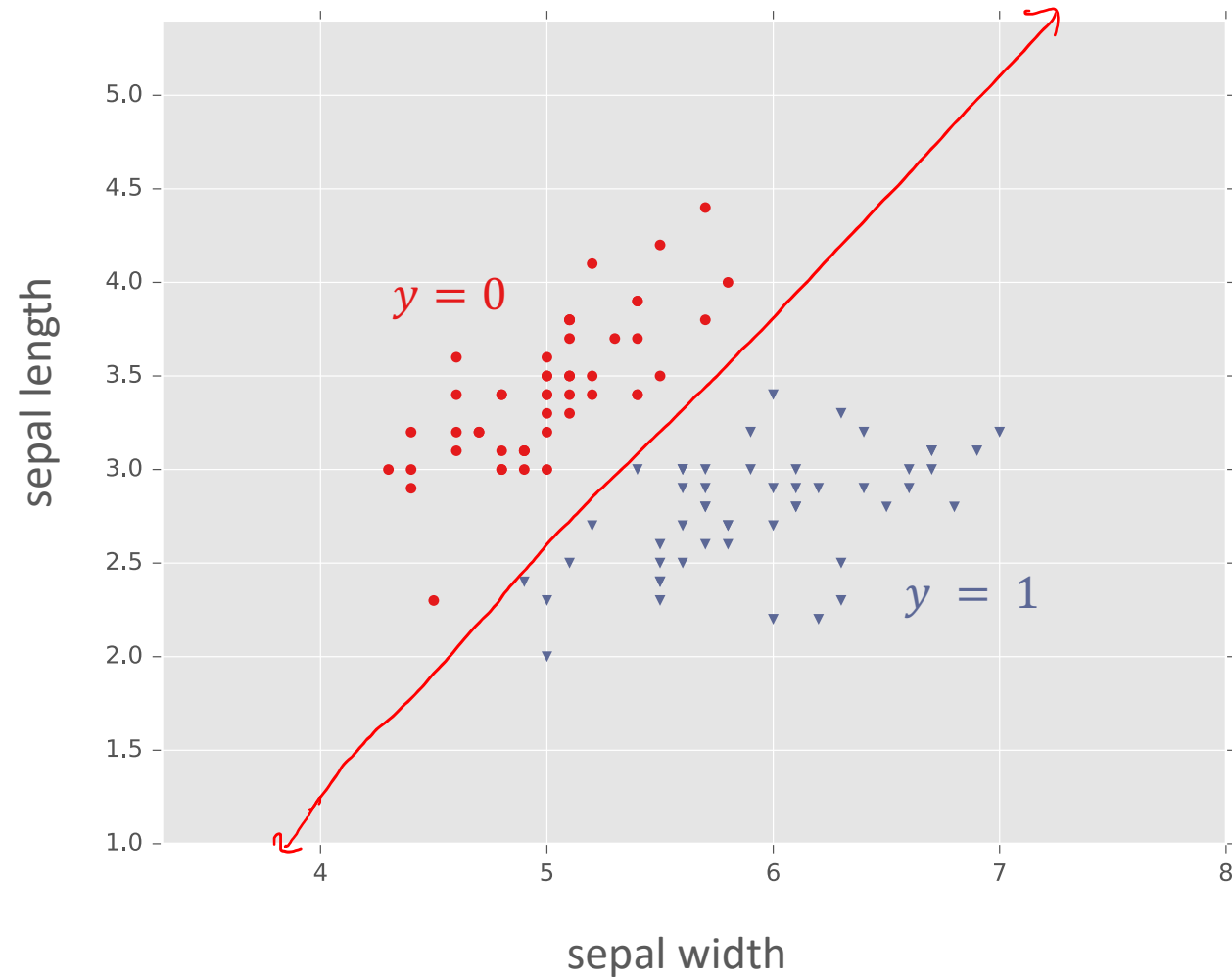
Henry Chai

5/24/23

# Front Matter

- Announcements:

  - PA1 released 5/18, due 5/25 (tomorrow) at 11:59 PM

  - PA2 released 5/25 (tomorrow), due 6/01 at 11:59 PM

  - No lecture or OH on Memorial Day (5/29);

    please plan accordingly!

- Recommended Readings:

  - Mitchell, Chapter 4.4

# Recall: Fisher Iris Dataset



$y = 0$

$y = 1$

sepal length

sepal width

Figure courtesy of Matt Gormley

# Linear Algebra Review

- Notation: in this class vectors will be assumed to be 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

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

- The $L2$-norm of $\boldsymbol{a} = \|\boldsymbol{a}\|_2 = \sqrt{\boldsymbol{a}^T \boldsymbol{a}}$

$$\sqrt{a^T a} = \sqrt{\sum_{d=1}^{D} a_d^2}$$

$$\|a\|_2^2$$

- Two vectors are *orthogonal* iff

$$\boldsymbol{a}^T \boldsymbol{b} = 0$$

$$e.g. \begin{bmatrix} 0 \\ 1 \end{bmatrix} \sim \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$
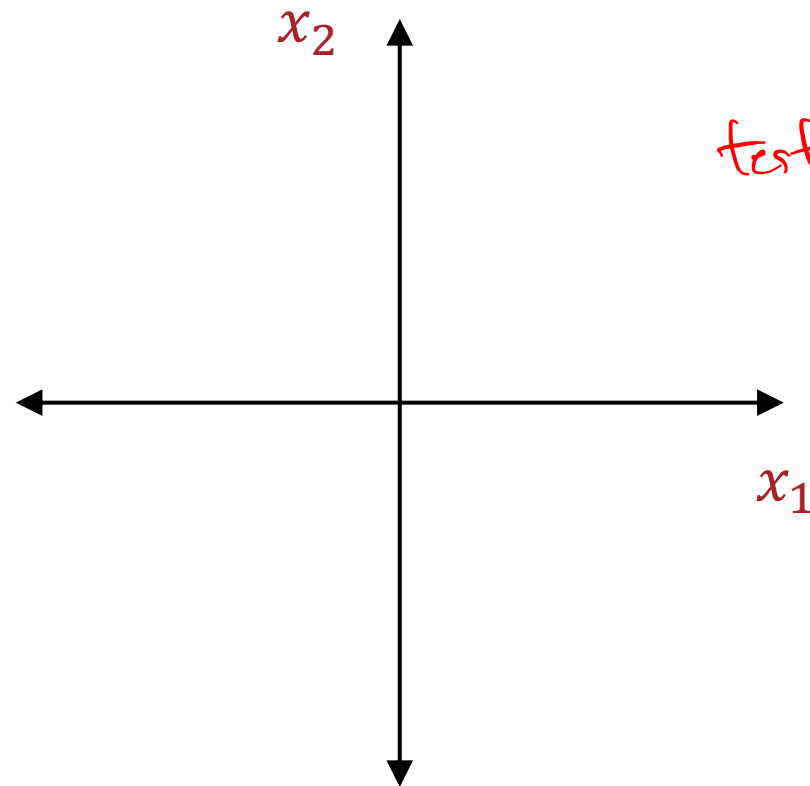
# Geometry Warm-up

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 $w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$

rewrite

$x_2 = mx_1 + b$

set $x_1 = 0$

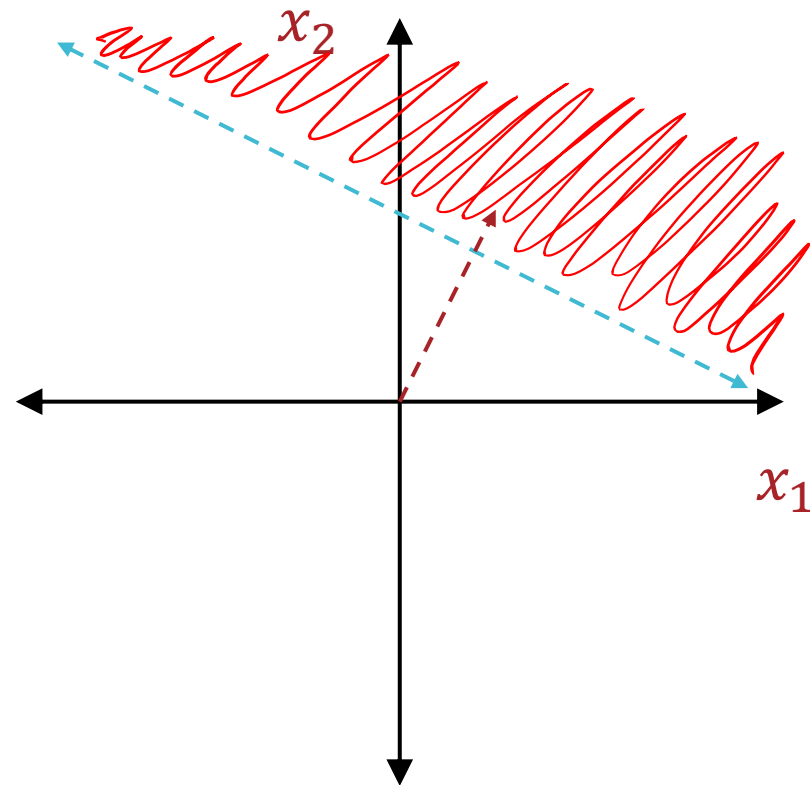$x_2 = 0$

test $(0, 0)$

$(0, 0)$ not in region

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 $b = -4$.

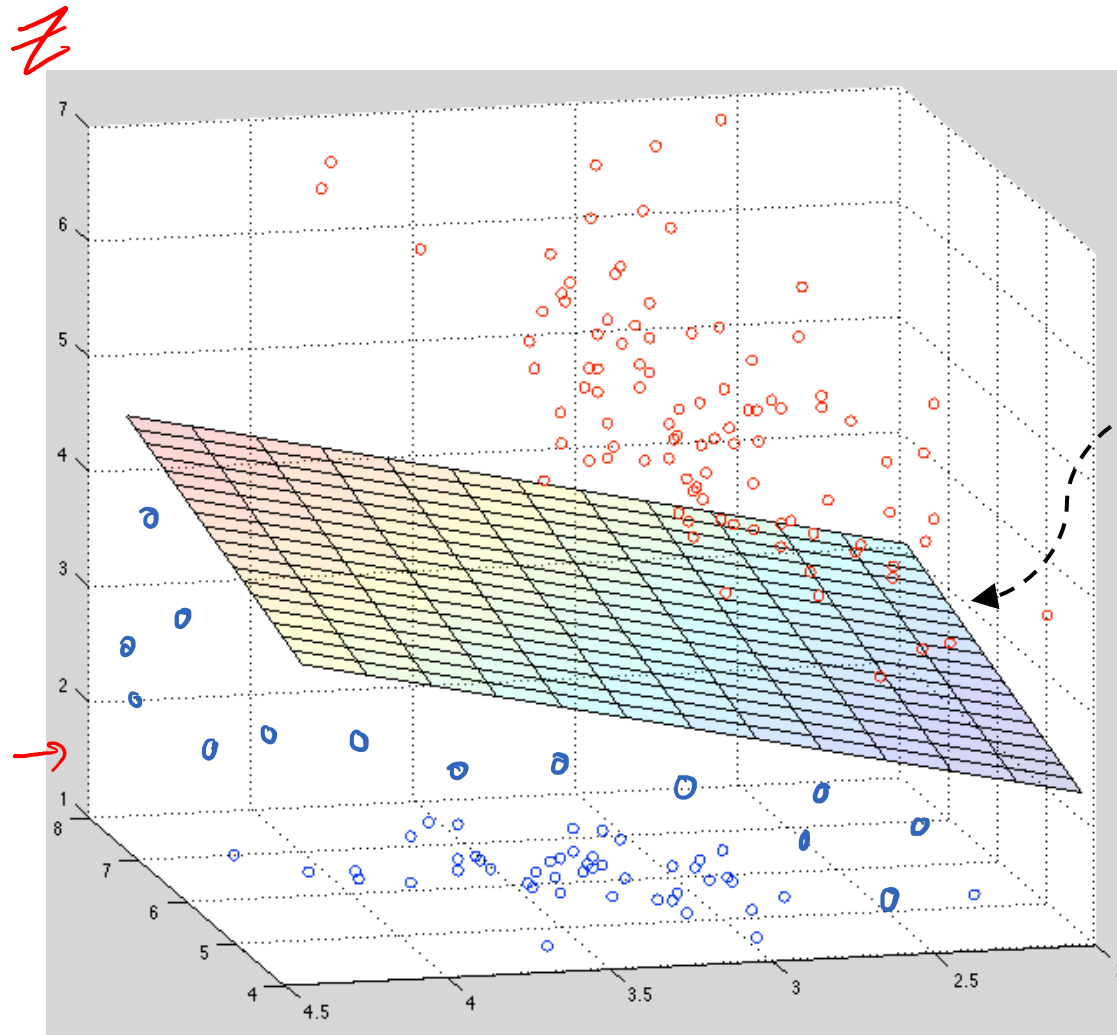2. Then draw the vector $w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$

$$\left( \sum_{d=1}^{D} w_d x_d \right) + b$$

# Linear Decision Boundaries

- In 2 dimensions, $w_1 x_1 + w_2 x_2 + b = 0$ defines a *line*

- In 3 dimensions, $w_1 x_1 + w_2 x_2 + w_3 x_3 + b = 0$ defines a *plane*

- In 4+ dimensions, $\boldsymbol{w}^T \boldsymbol{x} + b = 0$ 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*:

  - $\mathcal{S}_+ = \{\boldsymbol{x} : \boldsymbol{w}^T \boldsymbol{x} + b > 0\}$ or all $\boldsymbol{x}$ s.t. $\boldsymbol{w}^T \boldsymbol{x} + b$ is positive

  - $\mathcal{S}_- = \{\boldsymbol{x} : \boldsymbol{w}^T \boldsymbol{x} + b < 0\}$ or all $\boldsymbol{x}$ s.t. $\boldsymbol{w}^T \boldsymbol{x} + b$ is negative

# Linear Decision Boundaries: Example



Figure courtesy of Matt Gormley

Goal: learn classifiers of the form $h(x) = \text{sign}(w^T x + b)$ (assuming $y \in \{-1, +1\}$)

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

and $b$

# 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:

  - Predicting stock prices

  - Recommender systems

  - Medical diagnosis

  - Robotics

# Online Learning: Setup

- For $t = 1, 2, 3, \ldots$

  - Receive an unlabeled example, $\boldsymbol{x}^{(t)}$

  - Predict its label, $\hat{y} = h_{\boldsymbol{w}, b}\left(\boldsymbol{x}^{(t)}\right)$

  - Observe its true label, $y^{(t)}$

  - Pay a penalty if we made a mistake, $\hat{y} \neq y^{(t)}$

  - Update the parameters, $\boldsymbol{w}$ and $b$

- Goal: minimize the number of mistakes made

**(Online) Perceptron Learning Algorithm**

- Initialize the weight vector and intercept to all zeros:

$$\boldsymbol{w} = \begin{bmatrix} 0 & 0 & \cdots & 0 \end{bmatrix} \text{ and } b = 0$$

- For $t = 1, 2, 3, \ldots$

  - Receive an unlabeled example, $\boldsymbol{x}^{(t)}$

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

  - 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 ($y^{(t)} = -1, \hat{y} = +1$):
    - $\boldsymbol{w} \leftarrow \boldsymbol{w} - \boldsymbol{x}^{(t)}$
    - $b \leftarrow b - 1$

# (Online) Perceptron Learning Algorithm

- Initialize the weight vector and intercept to all zeros:

$$\boldsymbol{w} = \begin{bmatrix} 0 & 0 & \cdots & 0 \end{bmatrix} \text{ and } b = 0$$

- For $t = 1, 2, 3, \dots$
  - Receive an unlabeled example, $\boldsymbol{x}^{(t)}$
  - Predict its label, $\hat{y} = \text{sign}(\boldsymbol{w}^T \boldsymbol{x} + b) = \begin{cases} +1 \text{ if } \boldsymbol{w}^T \boldsymbol{x} + b \geq 0 \\ -1 \text{ otherwise} \end{cases}$
  - Observe its true label, $y^{(t)}$
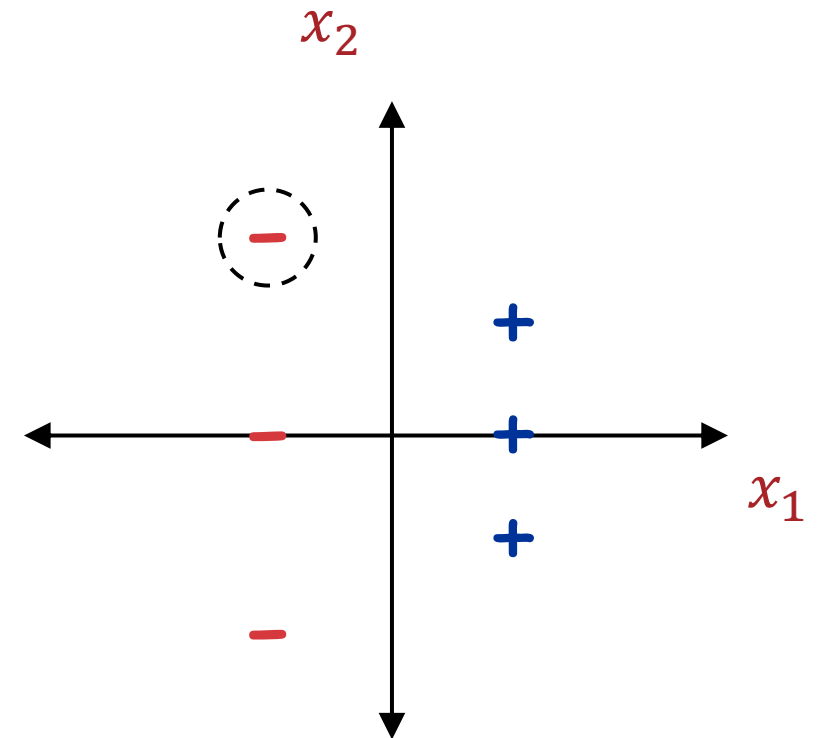  - If we misclassified an example ($y^{(t)} \neq \hat{y}$):
    - $\boldsymbol{w} \leftarrow \boldsymbol{w} + y^{(t)} \boldsymbol{x}^{(t)}$
    - $b \leftarrow b + y^{(t)}$

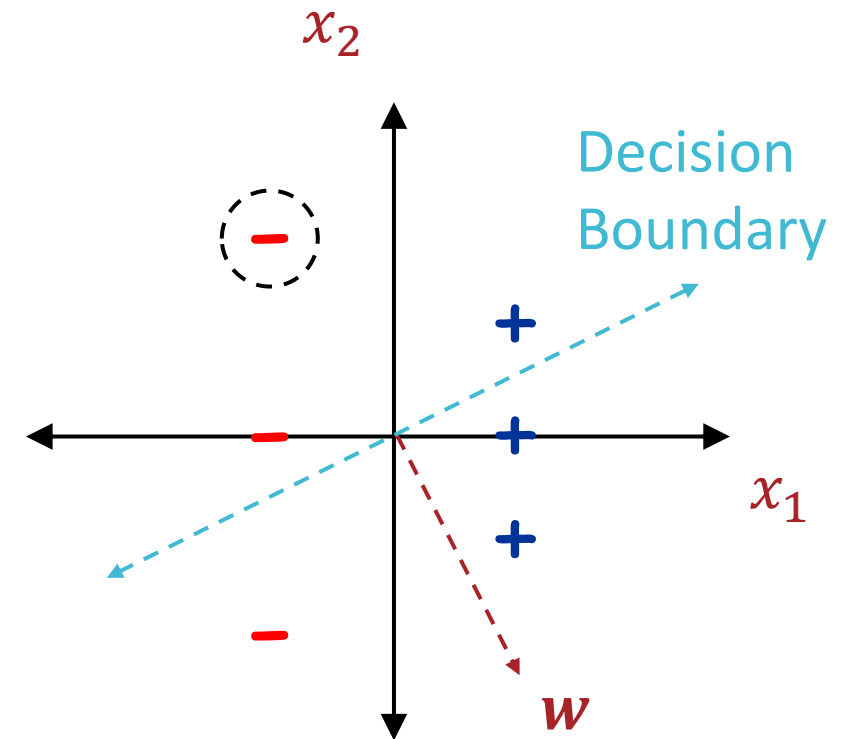# (Online) Perceptron Learning Algorithm: Example (no Intercept)

| $x_1$ | $x_2$ | $\widehat{y}$ | $y$ | Mistake? |
|-------|-------|-------|-----|----------|
| $-1$ | $2$ | $+$ | $-$ | Yes |

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

Example courtesy of Nina Balcan

(Online) Perceptron Learning Algorithm: Example (no Intercept)

| $x_1$ | $x_2$ | $\widehat{y}$ | $y$ | Mistake? |
|---|---|---|---|---|
| $-1$ | $2$ | $+$ | $-$ | Yes |

Decision Boundary

$x_2$

$x_1$

$\boldsymbol{w}$

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

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

Example courtesy of Nina Balcan

# (Online) Perceptron Learning Algorithm: Example (no Intercept)

| $x_1$ | $x_2$ | $\hat{y}$ | $y$ | Mistake? |
|---|---|---|---|---|
| $-1$ | $2$ | $+$ | $-$ | Yes |
| $1$ | $0$ | $+$ | $+$ | No |

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

$x_2$

Decision Boundary

$x_1$

$w$

Example courtesy of Nina Balcan

# (Online) Perceptron Learning Algorithm: Example (no Intercept)

| $x_1$ | $x_2$ | $\widehat{y}$ | $y$ | Mistake? |
|---|---|---|---|---|
| $-1$ | 2 | $+$ | $-$ | Yes |
| 1 | 0 | $+$ | $+$ | No |
| 1 | 1 | $-$ | $+$ | Yes |

$x_2$

Decision Boundary

$x_1$

$\boldsymbol{w}$

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

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

Example courtesy of Nina Balcan

## (Online) Perceptron Learning Algorithm: Example (no Intercept)

| $x_1$ | $x_2$ | $\widehat{y}$ | $y$ | Mistake? |
|-------|-------|------|-----|----------|
| $-1$ | 2 | $+$ | $-$ | Yes |
| 1 | 0 | $+$ | $+$ | No |
| 1 | 1 | $-$ | $+$ | Yes |



$x_2$

Decision Boundary

$x_1$

$w$

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

$$w \leftarrow w + y^{(3)} x^{(3)} = \begin{bmatrix} 1 \\ -2 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$
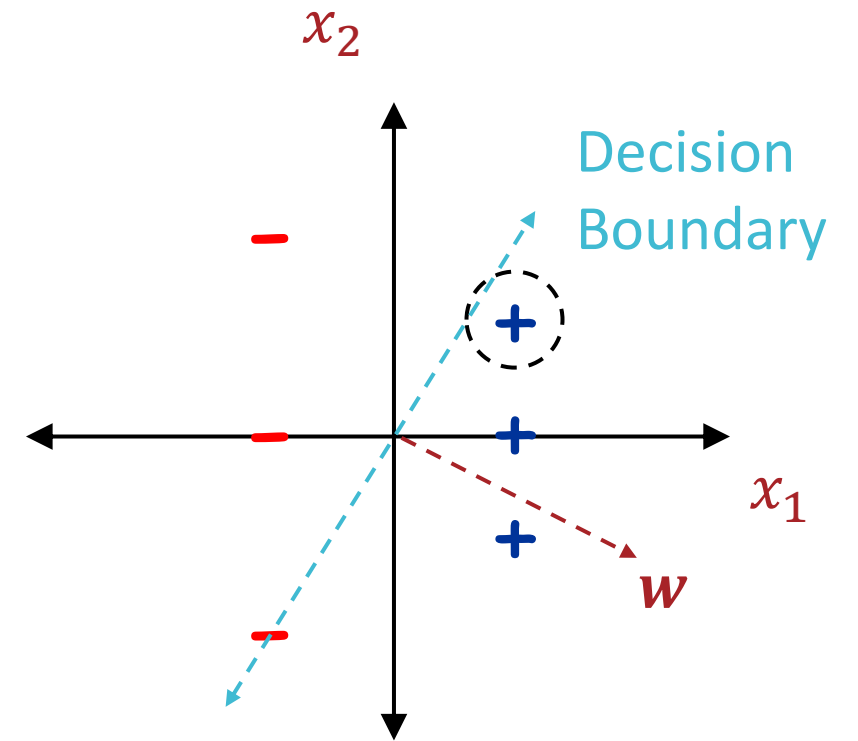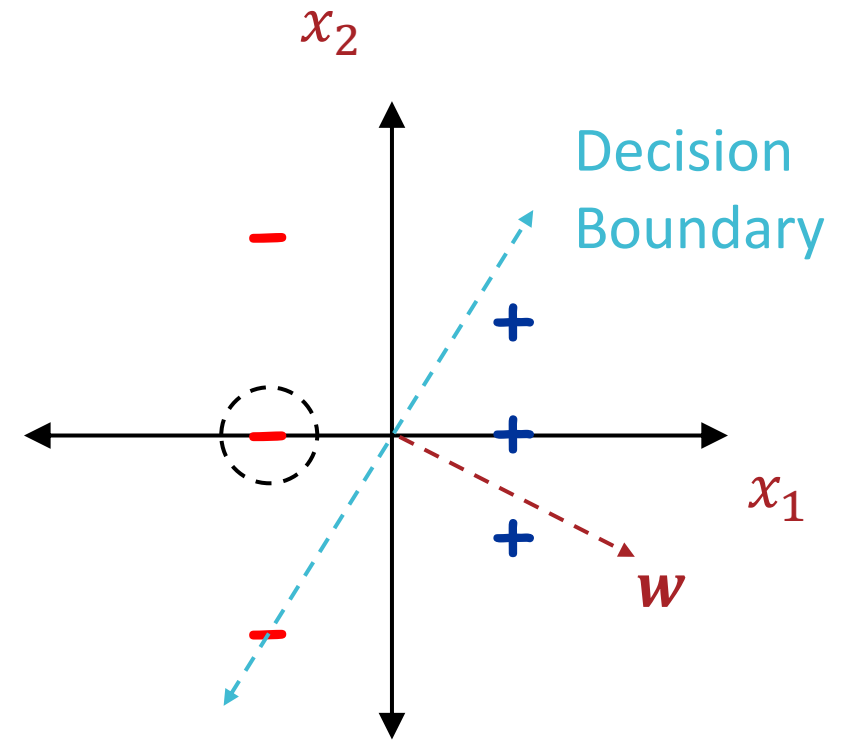
# (Online) Perceptron Learning Algorithm: Example (no Intercept)

| $x_1$ | $x_2$ | $\hat{y}$ | $y$ | Mistake? |
|-------|-------|-----------|-----|----------|
| $-1$ | 2 | $+$ | $-$ | Yes |
| 1 | 0 | $+$ | $+$ | No |
| 1 | 1 | $-$ | $+$ | Yes |
| $-1$ | 0 | $-$ | $-$ | No |

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

$x_2$

Decision Boundary

$x_1$

$w$

Example courtesy of Nina Balcan

## (Online) Perceptron Learning Algorithm: Example (no Intercept)

| $x_1$ | $x_2$ | $\hat{y}$ | $y$ | Mistake? |
|-------|-------|-----------|-----|----------|
| $-1$ | 2 | $+$ | $-$ | Yes |
| 1 | 0 | $+$ | $+$ | No |
| 1 | 1 | $-$ | $+$ | Yes |
| $-1$ | 0 | $-$ | $-$ | No |
| $-1$ | $-2$ | $+$ | $-$ | Yes |

$x_2$

Decision Boundary

$x_1$

$\boldsymbol{w}$

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

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

Example courtesy of Nina Balcan

# (Online) Perceptron Learning Algorithm: Example (no Intercept)

| $x_1$ | $x_2$ | $\widehat{y}$ | $y$ | Mistake? |
|:---:|:---:|:---:|:---:|:---:|
| $-1$ | $2$ | $+$ | $-$ | Yes |
| $1$ | $0$ | $+$ | $+$ | No |
| $1$ | $1$ | $-$ | $+$ | Yes |
| $-1$ | $0$ | $-$ | $-$ | No |
| $-1$ | $-2$ | $+$ | $-$ | Yes |

Decision Boundary
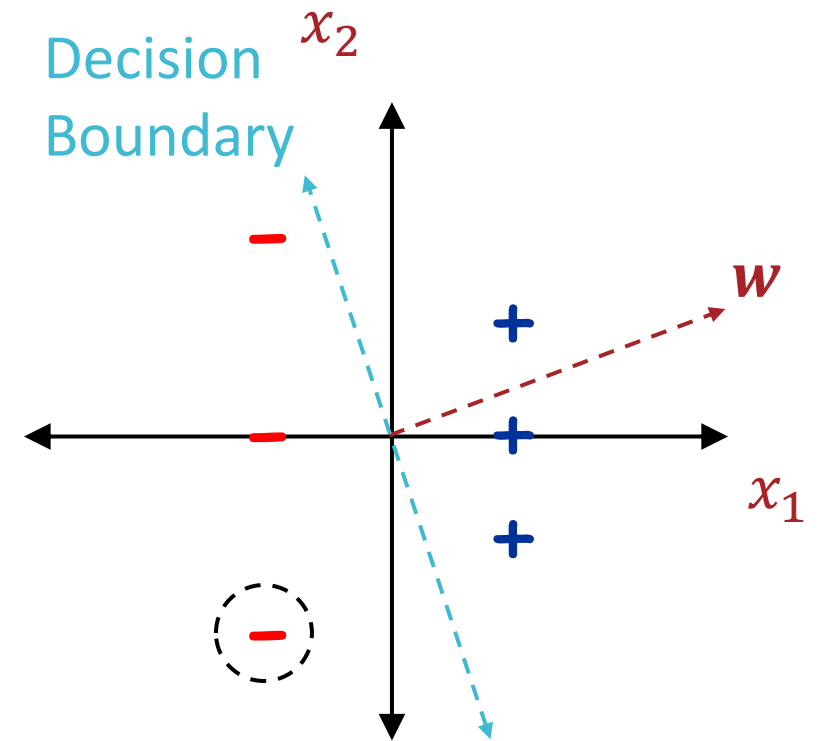
$x_2$

$w$

$x_1$

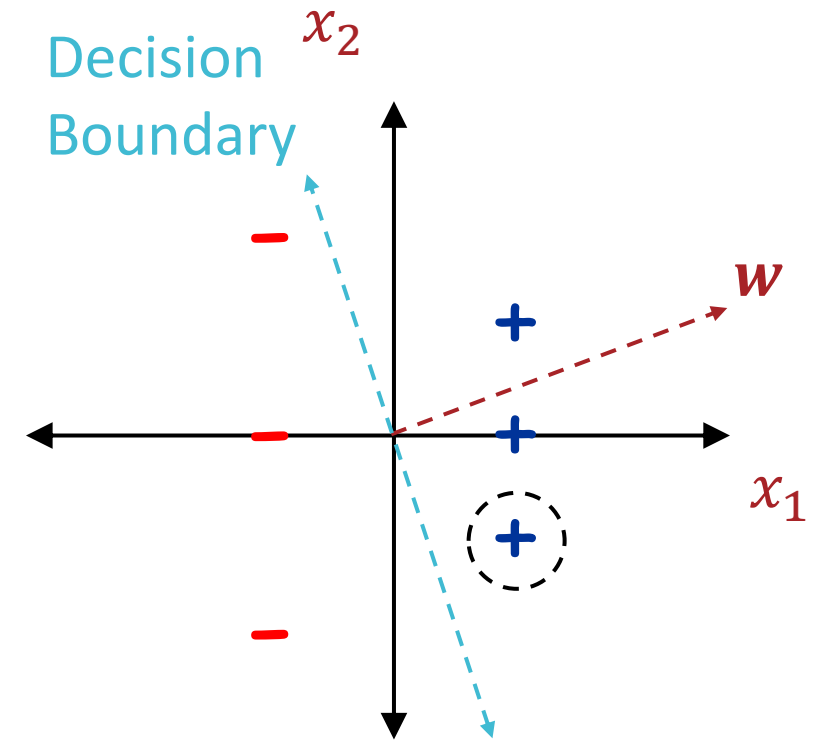$$w = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$

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

Example courtesy of Nina Balcan

# (Online) Perceptron Learning Algorithm: Example (no Intercept)

| $x_1$ | $x_2$ | $\widehat{y}$ | $y$ | Mistake? |
|---|---|---|---|---|
| −1 | 2 | + | − | Yes |
| 1 | 0 | + | + | No |
| 1 | 1 | − | + | Yes |
| −1 | 0 | − | − | No |
| −1 | −2 | + | − | Yes |
| 1 | −1 | + | + | No |

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

Decision Boundary

$x_2$

$w$

$x_1$

Example courtesy of Nina Balcan

# Updating the Intercept

- 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

## Notational Hack

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

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

- … we can just fold the intercept into the weight vector!

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

# (Online) Perceptron Learning Algorithm

- Initialize the weight vector and intercept to all zeros:

$$\boldsymbol{w} = \begin{bmatrix} 0 & 0 & \cdots & 0 \end{bmatrix} \text{ and } b = 0$$

- For $t = 1, 2, 3, \ldots$

  - Receive an unlabeled example, $\boldsymbol{x}^{(t)}$

  *instead use* $\boldsymbol{w}^T x + 2b$

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

  - Observe its true label, $y^{(t)}$

  - If we misclassified an example ($y^{(t)} \neq \hat{y}$):

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

# (Online) Perceptron Learning Algorithm

- Initialize the parameters to all zeros:

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

- For $t = 1, 2, 3, \ldots$

  - Receive an unlabeled example, $\boldsymbol{x}^{(t)}$

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

  - Predict its label, $\hat{y} = \mathrm{sign}\left(\boldsymbol{\theta}^T \boldsymbol{x}'^{(t)}\right) = \begin{cases} +1 \text{ if } \boldsymbol{\theta}^T \boldsymbol{x}'^{(t)} \geq 0 \\ -1 \text{ otherwise} \end{cases}$

  - Observe its true label, $y^{(t)}$

  - If we misclassified an example ($y^{(t)} \neq \hat{y}$):

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

Automatically handles updating the intercept

# Perceptron Learning Algorithm: Intuition

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

$$\rightarrow \Theta^T x \text{ is negative}$$

$$\rightarrow \Theta_{new} = \Theta + yx = \Theta + x$$

$$\rightarrow \Theta_{new}^T x = (\Theta + x)^T x$$

$$= (\Theta^T + x^T) x \nearrow \sum_{d=0}^{D} x_d^2$$

$$= \Theta^T x + x^T x$$

which is "less negative" than $\Theta^T x$

$\rightarrow$ a similar thing holds for mistakes on negative points

# (Online) Perceptron Learning Algorithm: Inductive Bias

— The decision boundary is linear and correcting recent mistakes is the priority (even over potentially misclassifying previously seen data points)

## (Online) Perceptron Learning Algorithm

- Initialize the parameters to all zeros:

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

- For $t = 1, 2, 3, \ldots$

  - Receive an unlabeled example, $\boldsymbol{x}^{(t)}$

  - Predict its label, $\hat{y} = \text{sign}\left(\boldsymbol{\theta}^T \boldsymbol{x}'^{(t)}\right)$

  - Observe its true label, $y^{(t)}$

  - If we misclassified an example ($y^{(t)} \neq \hat{y}$):

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

# (Batch) Perceptron Learning Algorithm

- Input: $\mathcal{D} = \{(\boldsymbol{x}^{(1)}, y^{(1)}), (\boldsymbol{x}^{(2)}, y^{(2)}), \dots, (\boldsymbol{x}^{(N)}, y^{(N)})\}$

- Initialize the parameters to all zeros:

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

- While NOT CONVERGED

  - For $t \in \{1, \dots, N\}$

    - Predict the label of $\boldsymbol{x'}^{(t)}$, $\hat{y} = \mathrm{sign}\left(\boldsymbol{\theta}^T \boldsymbol{x'}^{(t)}\right)$

    - Observe its true label, $y^{(t)}$

    - If we misclassified $\boldsymbol{x'}^{(t)}$ ($y^{(t)} \neq \hat{y}$):

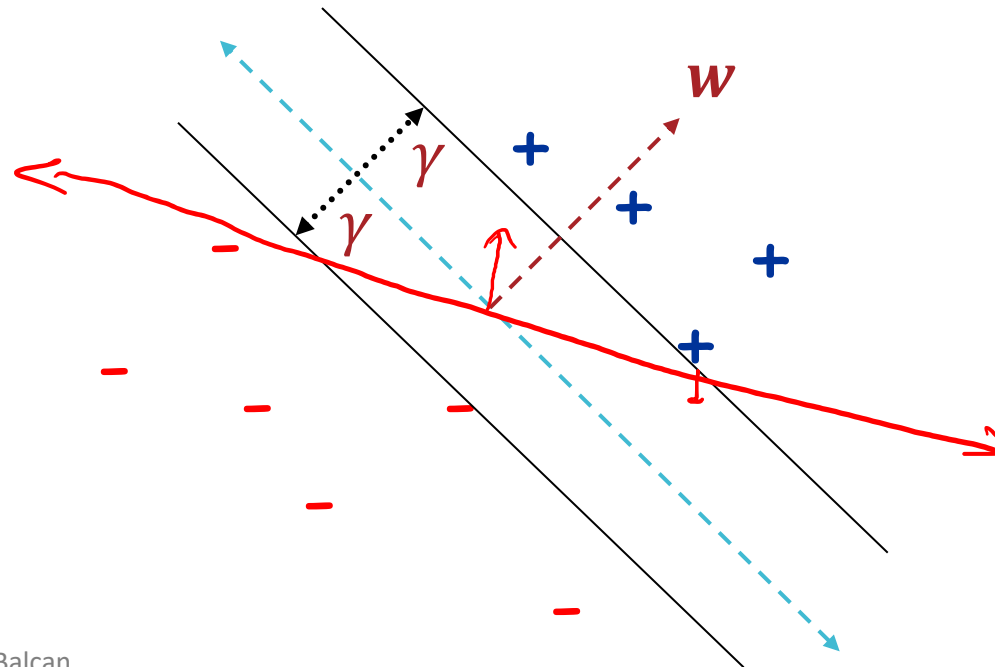      - $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + y^{(t)} \boldsymbol{x'}^{(t)}$

# True or False: The parameter vector $w$ learned by the batch Perceptron Learning Algorithm can be written as a linear combination of the examples, i.e.,

$$w = c_1 x^{(1)} + c_2 x^{(2)} + \ldots + c_N x^{(N)}$$

True

False

# Perceptron Mistake Bound

- Definitions:
  - A dataset $\mathcal{D}$ is *linearly separable* if $\exists$ a linear decision boundary that perfectly classifies the examples in $\mathcal{D}$
  - The margin, $\gamma$, of a dataset $\mathcal{D}$ is the greatest possible distance between a linear separator and the closest example in $\mathcal{D}$ to that linear separator
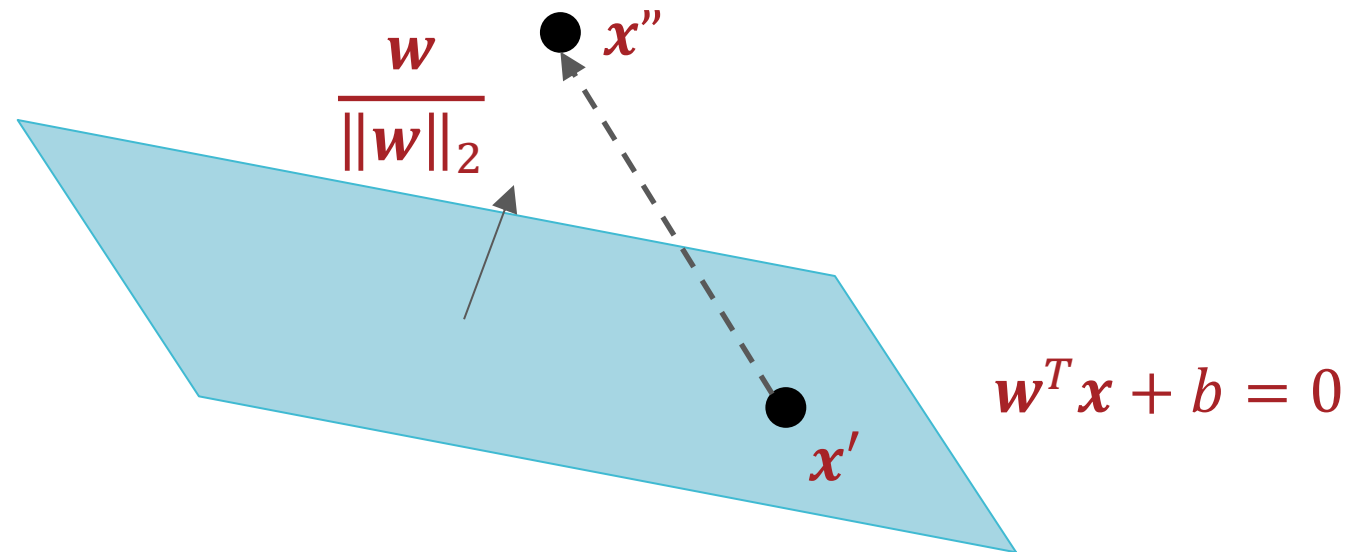
Figure courtesy of Nina Balcan

## Perceptron Mistake Bound

- Theorem: if the examples seen by the Perceptron Learning Algorithm (online and batch)

  1. lie in a ball of radius $R$ (centered around the origin)

  2. have a margin of $\gamma$

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

- 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!
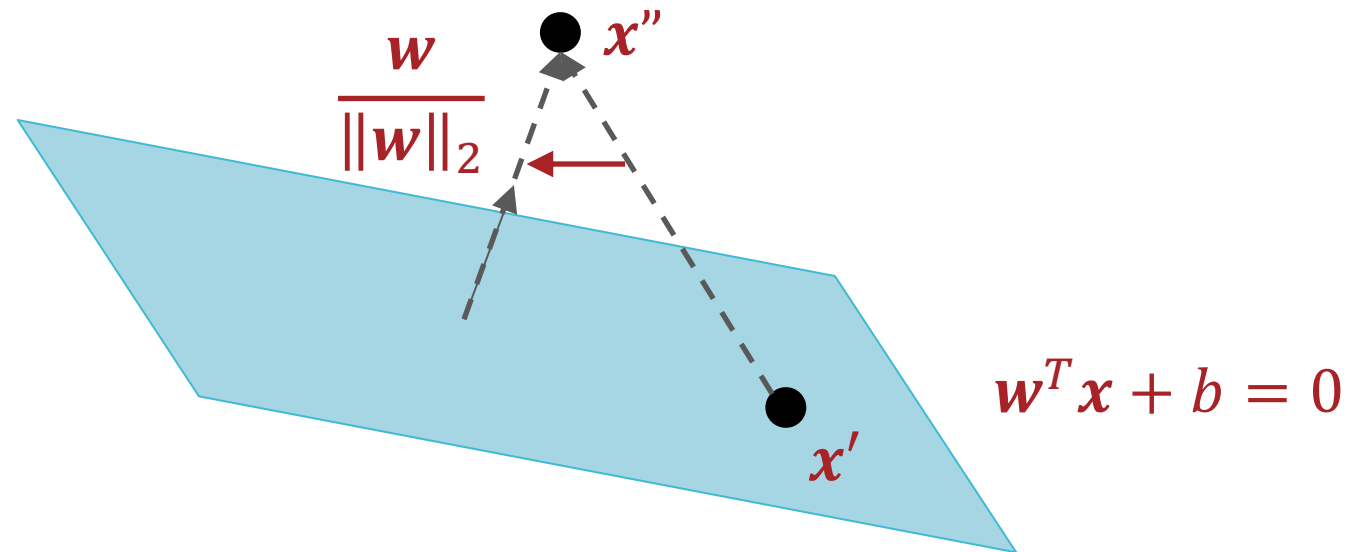
# Computing the Margin

- Let $\boldsymbol{x}'$ be an arbitrary point on the hyperplane $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ and let $\boldsymbol{x}''$ be an arbitrary point

- The distance between $\boldsymbol{x}''$ and $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ is equal to the magnitude of the projection of $\boldsymbol{x}'' - \boldsymbol{x}'$ onto $\dfrac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2}$, the unit vector orthogonal to the hyperplane

$$\boldsymbol{x}''$$

$$\frac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2}$$

$$\boldsymbol{w}^T\boldsymbol{x} + b = 0$$

$$\boldsymbol{x}'$$

# Computing the Margin

- Let $\boldsymbol{x}'$ be an arbitrary point on the hyperplane $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ and let $\boldsymbol{x}''$ be an arbitrary point

- The distance between $\boldsymbol{x}''$ and $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ is equal to the magnitude of the projection of $\boldsymbol{x}'' - \boldsymbol{x}'$ onto $\frac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2}$, the unit vector orthogonal to the hyperplane



$$\frac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2}$$

$$\boldsymbol{x}''$$

$$\boldsymbol{x}'$$

$$\boldsymbol{w}^T\boldsymbol{x} + b = 0$$

# Computing the Margin

- Let $\boldsymbol{x}'$ be an arbitrary point on the hyperplane $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ and let $\boldsymbol{x}''$ be an arbitrary point

- The distance between $\boldsymbol{x}''$ and $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ is equal to the magnitude of the projection of $\boldsymbol{x}'' - \boldsymbol{x}'$ onto $\dfrac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2}$, the unit vector orthogonal to the hyperplane



$$\boldsymbol{x}''$$

$$\dfrac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2}$$

$$\boldsymbol{x}'$$

$$\boldsymbol{w}^T\boldsymbol{x} + b = 0$$

# Computing the Margin

- Let $\boldsymbol{x}'$ be an arbitrary point on the hyperplane and let $\boldsymbol{x}''$ be an arbitrary point

- The distance between $\boldsymbol{x}''$ and $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ is equal to the magnitude of the projection of $\boldsymbol{x}'' - \boldsymbol{x}'$ onto $\dfrac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2}$, the unit vector orthogonal to the hyperplane

$$\frac{|\boldsymbol{w}^T x + b|}{\|w\|_2}$$

# Key Takeaways

- Batch vs. online learning

- Perceptron learning algorithm for binary classification

- Impact of the bias term in perceptron

- Inductive bias of perceptron

- Convergence properties, guarantees and limitations for the batch Perceptron learning algorithm