

10-301/601: Introduction to Machine Learning

Lecture 6 – Perceptron

Hoda Heidari, Henry Chai & Matt Gormley

2/5/24

Front Matter

- Announcements:
 - HW2 released 1/24, due 2/5 (today!) at 11:59 PM
 - HW3 released on 2/5 (today!), due 2/12 at 11:59 PM
 - HW3 is a written-only homework
 - **You may only use at most 2 late days on HW3**

Q & A:

After we do model selection using a validation dataset, should we train a final model using *both* the training and the validation datasets?

- Yes, absolutely! So really the sketch from last lecture should look something like:
 1. Split \mathcal{D} into $\mathcal{D}_{train} \cup \mathcal{D}_{val} \cup \mathcal{D}_{test}$
 2. Learn classifiers using \mathcal{D}_{train}
 3. Evaluate models using \mathcal{D}_{val} and choose the one with lowest *validation* error:
 4. **Learn a new classifier from the best model using $\mathcal{D}_{train} \cup \mathcal{D}_{val}$**
 5. Optionally, use \mathcal{D}_{test} to estimate the true error

Q & A:

Can we use k NNs with categorical features?

- Yes! We can either:

1. Convert categorical features into binary ones:

Cholesterol		Normal Cholesterol?	Abnormal Cholesterol?
Normal	→	1	0
Normal		1	0
Abnormal		0	1

2. Use a distance metric that works over categorical features e.g., the Hamming distance:

$$d(\mathbf{x}, \mathbf{x}') = \sum_{d=1}^D \mathbb{1}(x_d \neq x'_d)$$

- See HW3 for an example of this

Hyperparameter Optimization

- Given $\mathcal{D} = \mathcal{D}_{train} \cup \mathcal{D}_{val} \cup \mathcal{D}_{test}$, suppose we have multiple candidate hyperparameter settings:

$$\theta_1, \theta_2, \dots, \theta_M$$

- Learn a classifier for each setting using only \mathcal{D}_{train} :

$$h_1, h_2, \dots, h_M$$

- Evaluate each one using \mathcal{D}_{val} and choose the one with lowest *validation* error:

$$\hat{m} = \operatorname{argmin}_{m \in \{1, \dots, M\}} \operatorname{err}(h_m, \mathcal{D}_{val})$$

- Now $\operatorname{err}(h_{\hat{m}}, \mathcal{D}_{test})$ is a good estimate of $\operatorname{err}(h_{\hat{m}})$!

How to pick hyperparameter settings to try?

- Given $\mathcal{D} = \mathcal{D}_{train} \cup \mathcal{D}_{val} \cup \mathcal{D}_{test}$, suppose we have multiple candidate hyperparameter settings:

$$\theta_1, \theta_2, \dots, \theta_M$$

- Learn a classifier for each setting using only \mathcal{D}_{train} :

$$h_1, h_2, \dots, h_M$$

- Evaluate each one using \mathcal{D}_{val} and choose the one with lowest *validation* error:

$$\hat{m} = \operatorname{argmin}_{m \in \{1, \dots, M\}} \operatorname{err}(h_m, \mathcal{D}_{val})$$

- Now $\operatorname{err}(h_{\hat{m}}, \mathcal{D}_{test})$ is a good estimate of $\operatorname{err}(h_{\hat{m}})$!

General Methods for Hyperparameter Optimization

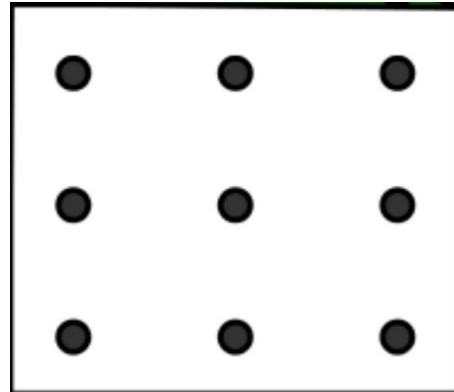
- Idea: set the hyperparameters to optimize some performance metric of the model
- Issue: if we have many hyperparameters that can all take on lots of different values, we might not be able to test all possible combinations
- Commonly used methods:
 - Grid search
 - Random search
 - Bayesian optimization (used by Google DeepMind to optimize the hyperparameters of AlphaGo: <https://arxiv.org/pdf/1812.06855v1.pdf>)
 - Evolutionary algorithms
 - Graduate-student descent

General Methods for Hyperparameter Optimization

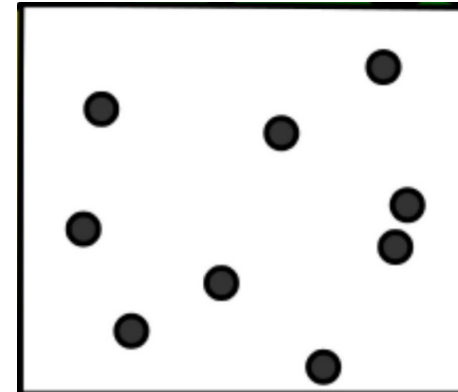
- Idea: set the hyperparameters to optimize some performance metric of the model
- Issue: if we have many hyperparameters that can all take on lots of different values, we might not be able to test all possible combinations
- Commonly used methods:
 - **Grid search**
 - **Random search**
 - Bayesian optimization (used by Google DeepMind to optimize the hyperparameters of AlphaGo: <https://arxiv.org/pdf/1812.06855v1.pdf>)
 - Evolutionary algorithms
 - Graduate-student descent

Grid Search vs. Random Search (Bergstra and Bengio, 2012)

Grid Layout



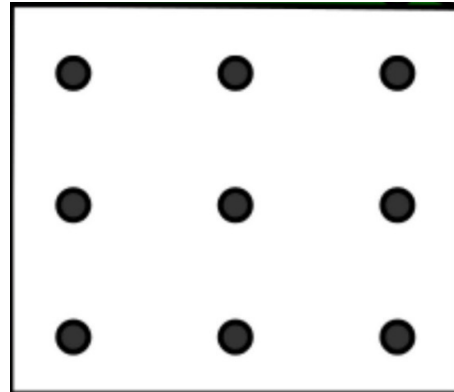
Random Layout



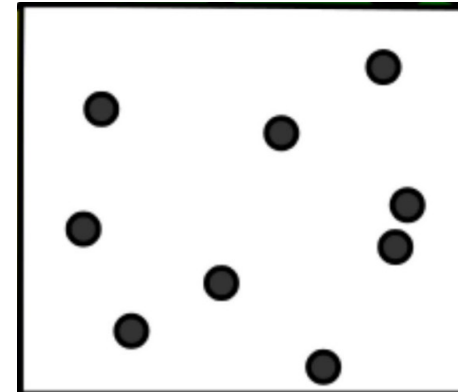
Poll Question 1:

Which hyperparameter optimization method do you think will perform better?

Grid Layout

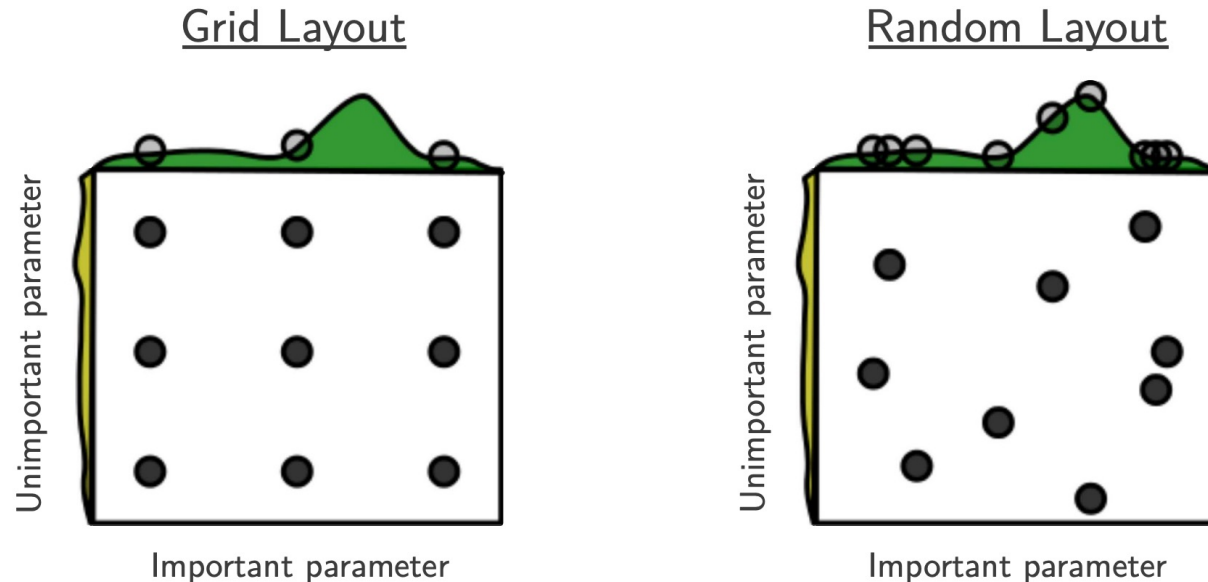


Random Layout



- A. Graduate student descent (**TOXIC**)
- B. Grid search
- C. Random search

Grid Search vs. Random Search (Bergstra and Bengio, 2012)



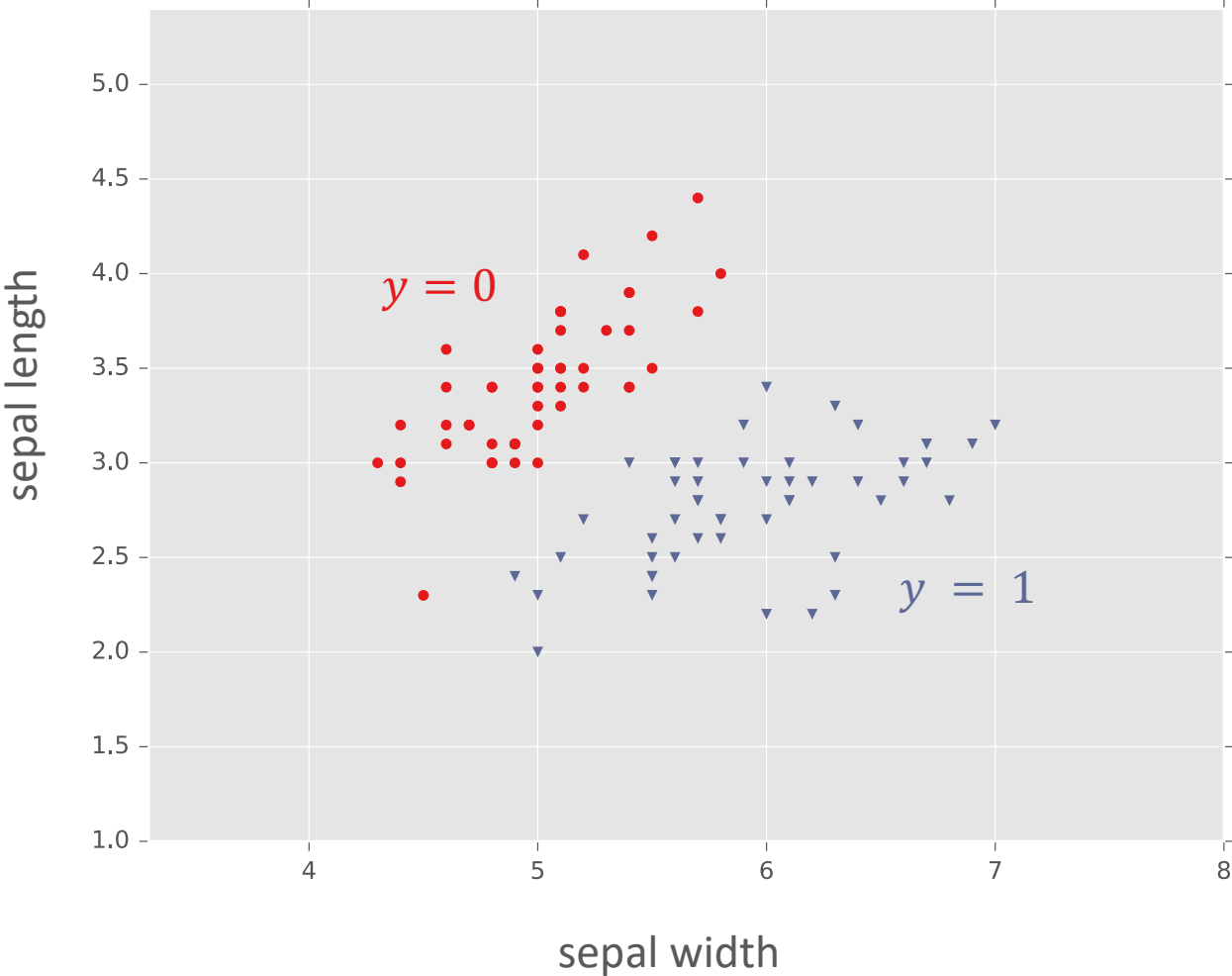
Grid and random search of nine trials for optimizing a function $f(x, y) = g(x) + h(y) \approx g(x)$ with *low effective dimensionality*. Above each square $g(x)$ is shown in green, and left of each square $h(y)$ is shown in yellow. With grid search, nine trials only test $g(x)$ in three distinct places. With random search, all nine trials explore distinct values of g . This failure of grid search is the rule rather than the exception in high dimensional hyper-parameter optimization.

Model Selection Learning Objectives

You should be able to...

- Plan an experiment that uses training, validation, and test datasets to predict the performance of a classifier on unseen data (without cheating)
- Explain the difference between (1) training error, (2) validation error, (3) cross-validation error, (4) test error, and (5) true error
- For a given learning technique, identify the model, learning algorithm, parameters, and hyperparameters
- Select an appropriate algorithm for optimizing (aka. learning) hyperparameters

Recall: Fisher Iris Dataset



Linear Algebra Review

- Notation: in this class vectors will be assumed to be column vectors by default, i.e.,

$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_D \end{bmatrix} \text{ and } \mathbf{a}^T = [a_1 \quad a_2 \quad \cdots \quad a_D]$$

- The dot product between two D -dimensional vectors is

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

- The L_2 -norm of $\mathbf{a} = \|\mathbf{a}\|_2 = \sqrt{\mathbf{a}^T \mathbf{a}}$
- Two vectors are *orthogonal* iff

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

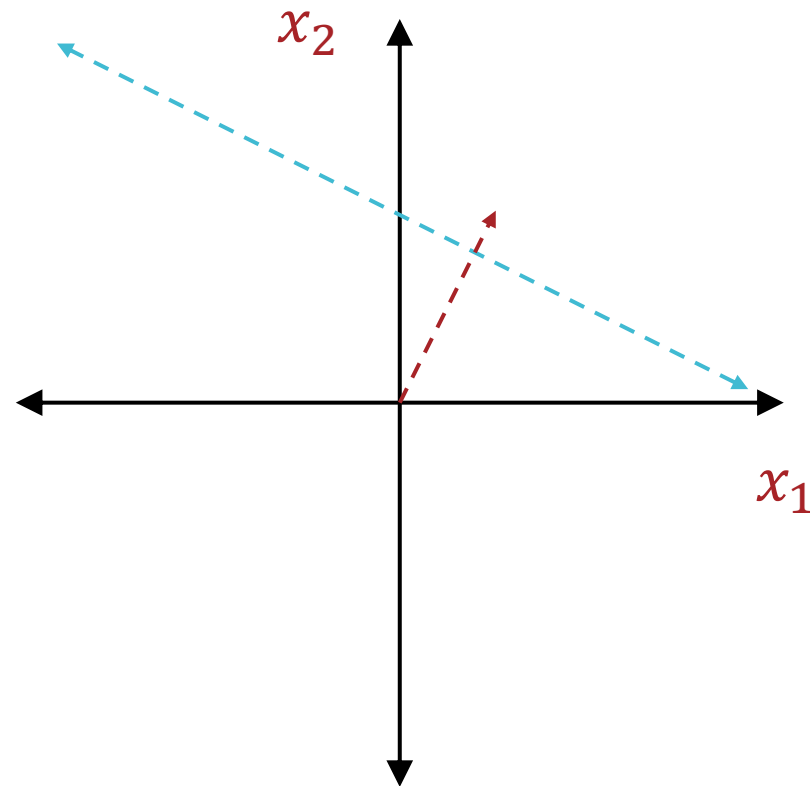
Geometry Warm-up

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

$$w_1x_1 + w_2x_2 + b > 0$$

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

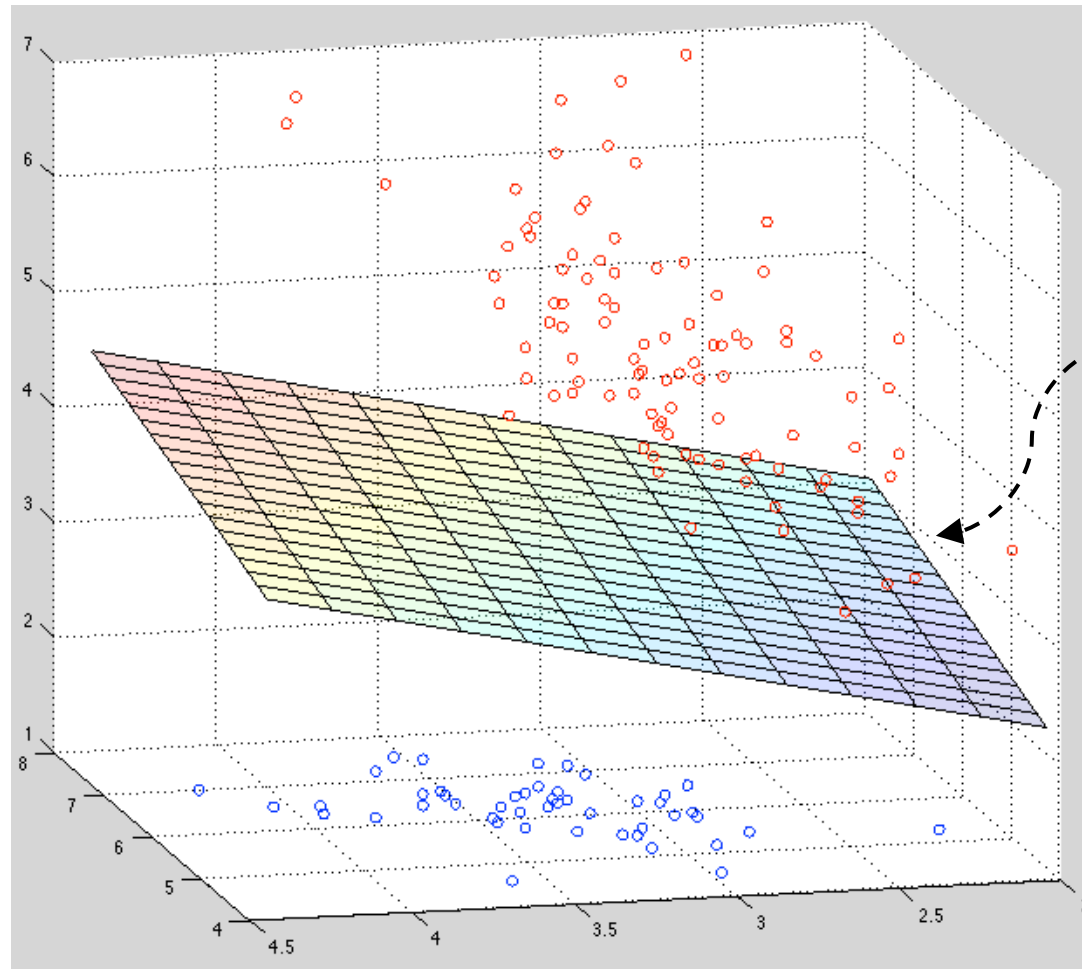
2. Then draw the vector $w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$



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, $\mathbf{w}^T \mathbf{x} + b = 0$ defines a *hyperplane*
 - The vector \mathbf{w} is always orthogonal to this hyperplane and always points in the direction where $\mathbf{w}^T \mathbf{x} + b > 0$!
- A hyperplane creates two *halfspaces*:
 - $\mathcal{S}_+ = \{\mathbf{x}: \mathbf{w}^T \mathbf{x} + b > 0\}$ or all \mathbf{x} s.t. $\mathbf{w}^T \mathbf{x} + b$ is positive
 - $\mathcal{S}_- = \{\mathbf{x}: \mathbf{w}^T \mathbf{x} + b < 0\}$ or all \mathbf{x} s.t. $\mathbf{w}^T \mathbf{x} + b$ is negative

Linear Decision Boundaries: Example



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

Key question: how do we learn the *parameters*, \mathbf{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 data points arrive gradually over time and we learn continuously
- Examples of online learning:
 - Predicting stock prices
 - Recommender systems
 - Medical diagnosis
 - Robotics

Online Learning: Setup

- For $t = 1, 2, 3, \dots$
 - Receive an unlabeled data point, $\mathbf{x}^{(t)}$
 - Predict its label, $\hat{y} = h_{\mathbf{w},b}(\mathbf{x}^{(t)})$
 - Observe its true label, $y^{(t)}$
 - Pay a penalty if we made a mistake, $\hat{y} \neq y^{(t)}$
 - Update the parameters, \mathbf{w} and b
- Goal: minimize the number of mistakes made

(Online) Perceptron Learning Algorithm

- Initialize the weight vector and intercept to all zeros:

$$\mathbf{w} = [0 \quad 0 \quad \dots \quad 0] \text{ and } b = 0$$

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

- Receive an unlabeled data point, $\mathbf{x}^{(t)}$

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

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

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

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

- $b \leftarrow b + 1$

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

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

- $b \leftarrow b - 1$

(Online) Perceptron Learning Algorithm

- Initialize the weight vector and intercept to all zeros:

$$\mathbf{w} = [0 \quad 0 \quad \dots \quad 0] \text{ and } b = 0$$

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

- Receive an unlabeled data point, $\mathbf{x}^{(t)}$

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

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

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

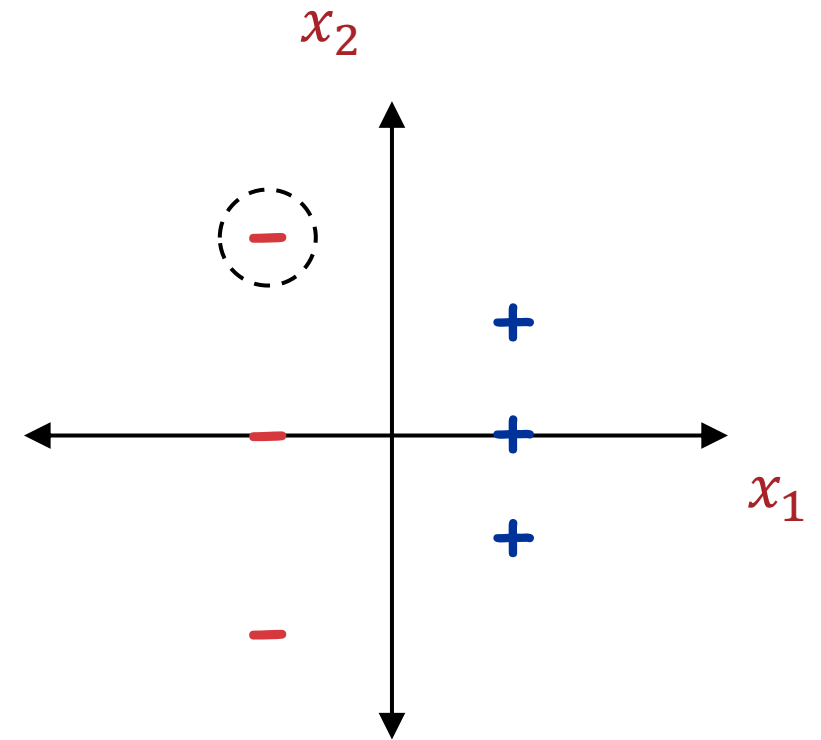
- $\mathbf{w} \leftarrow \mathbf{w} + y^{(t)} \mathbf{x}^{(t)}$

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

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

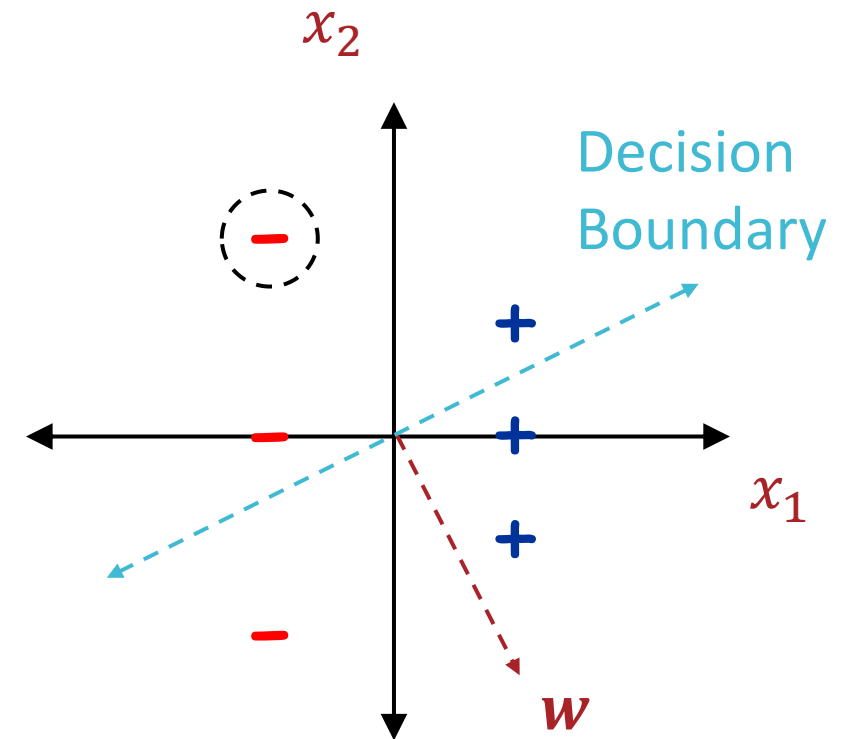
x_1	x_2	\hat{y}	y	Mistake?
-1	2	+	-	Yes

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



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

x_1	x_2	\hat{y}	y	Mistake?
-1	2	+	-	Yes



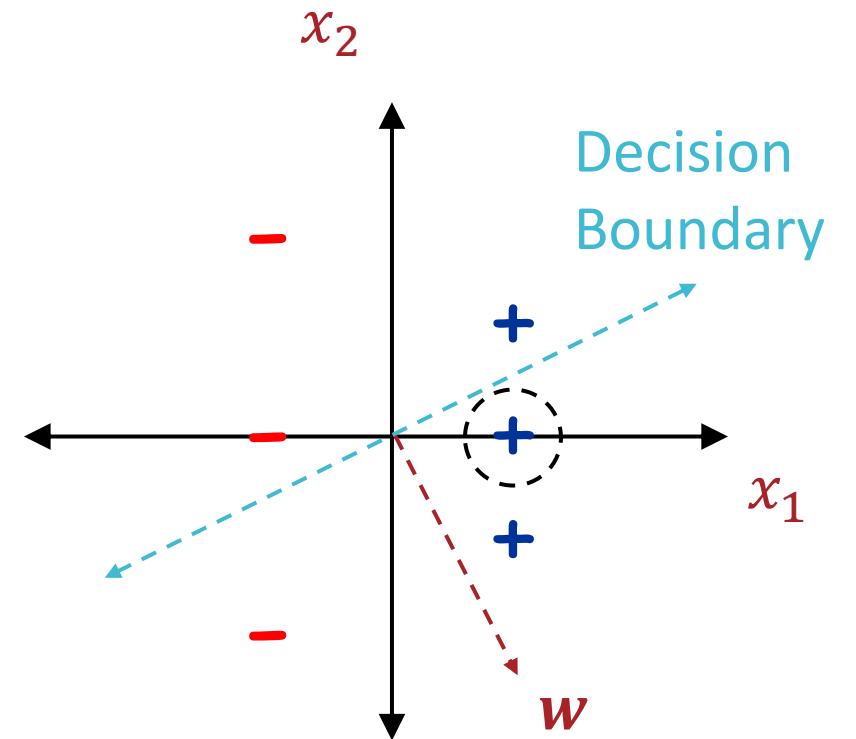
$$\mathbf{w} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

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

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

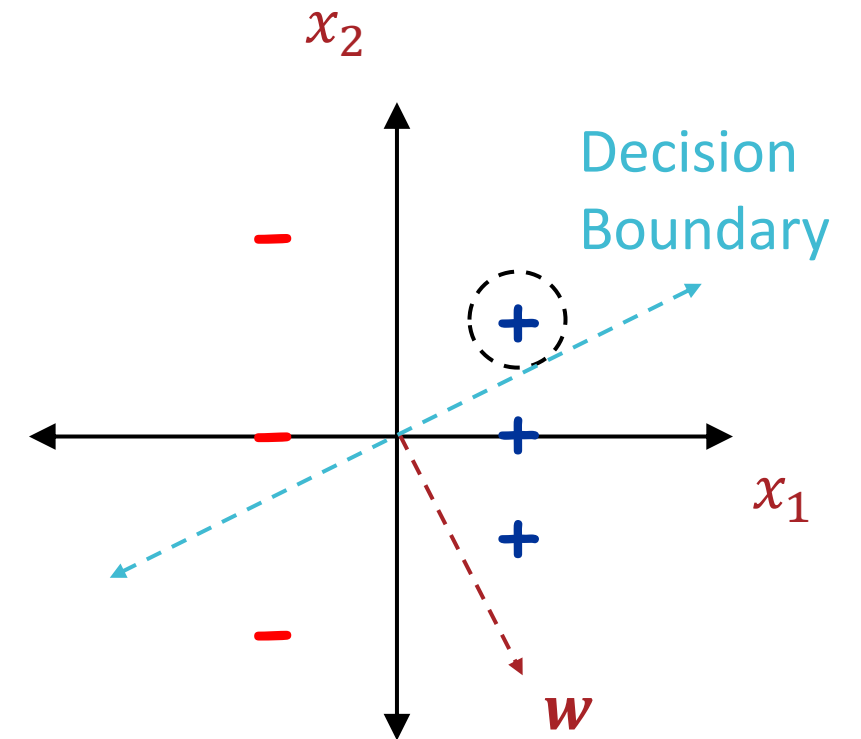
x_1	x_2	\hat{y}	y	Mistake?
-1	2	+	-	Yes
1	0	+	+	No

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



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

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



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

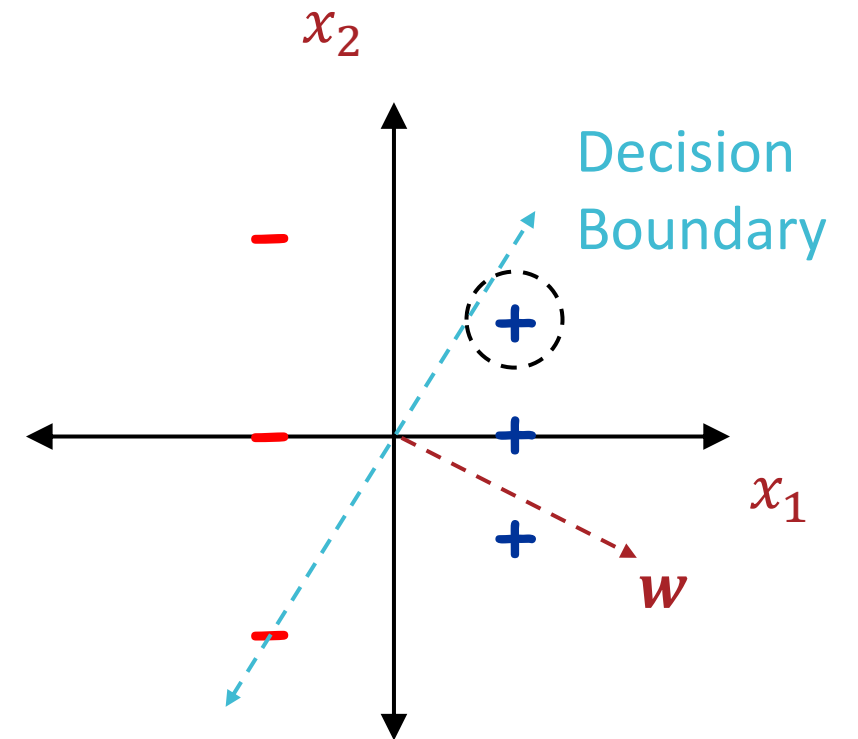
$$w \leftarrow w + y^{(3)}x^{(3)} = \begin{bmatrix} 1 \\ -2 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$

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

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

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

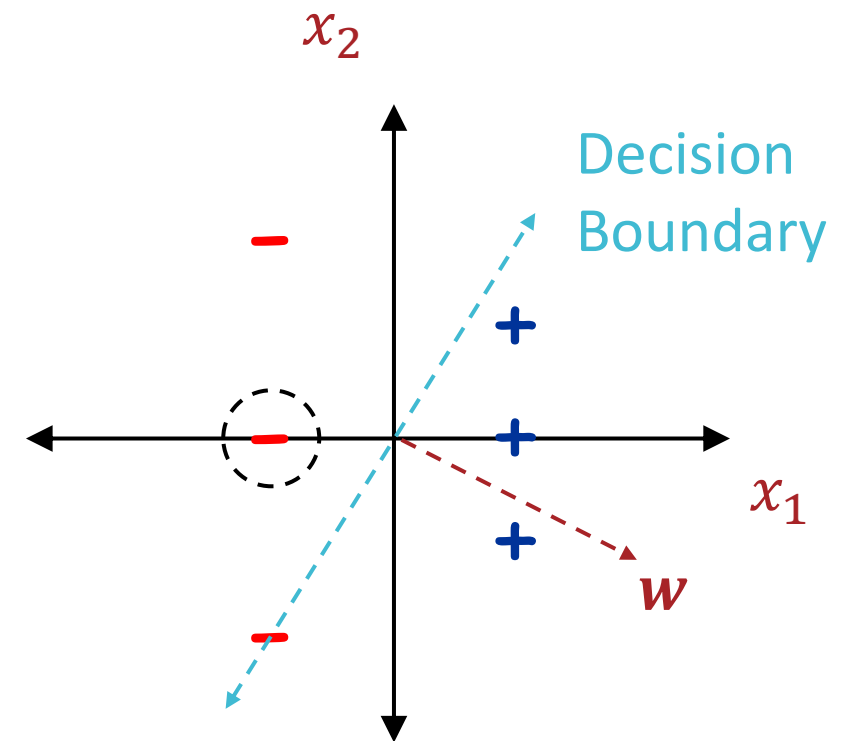
$$\mathbf{w} \leftarrow \mathbf{w} + y^{(3)} \mathbf{x}^{(3)} = \begin{bmatrix} 1 \\ -2 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$



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

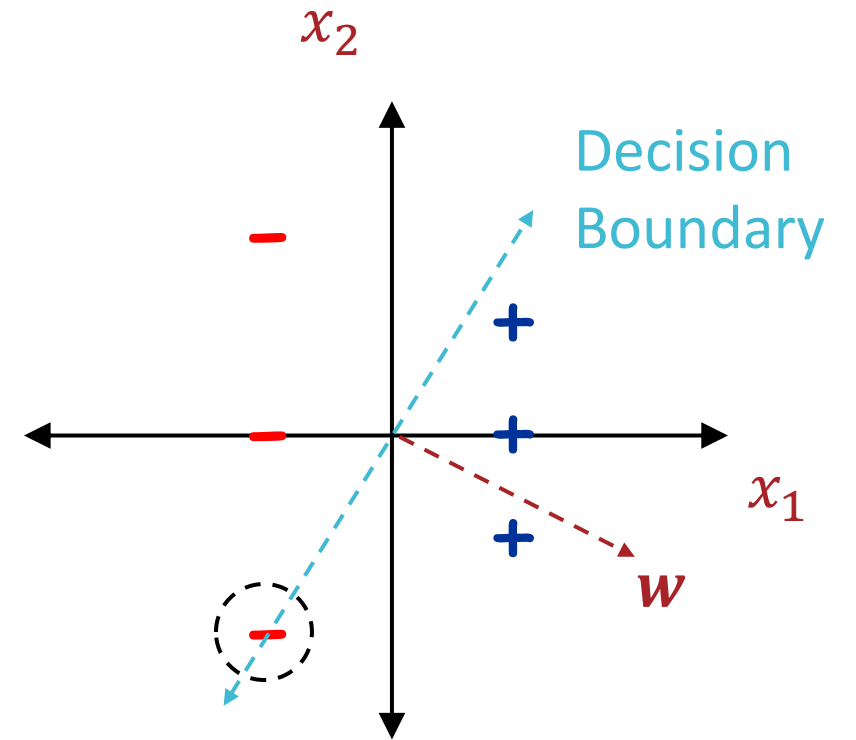
x_1	x_2	\hat{y}	y	Mistake?
-1	2	+	-	Yes
1	0	+	+	No
1	1	-	+	Yes
-1	0	-	-	No

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



(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

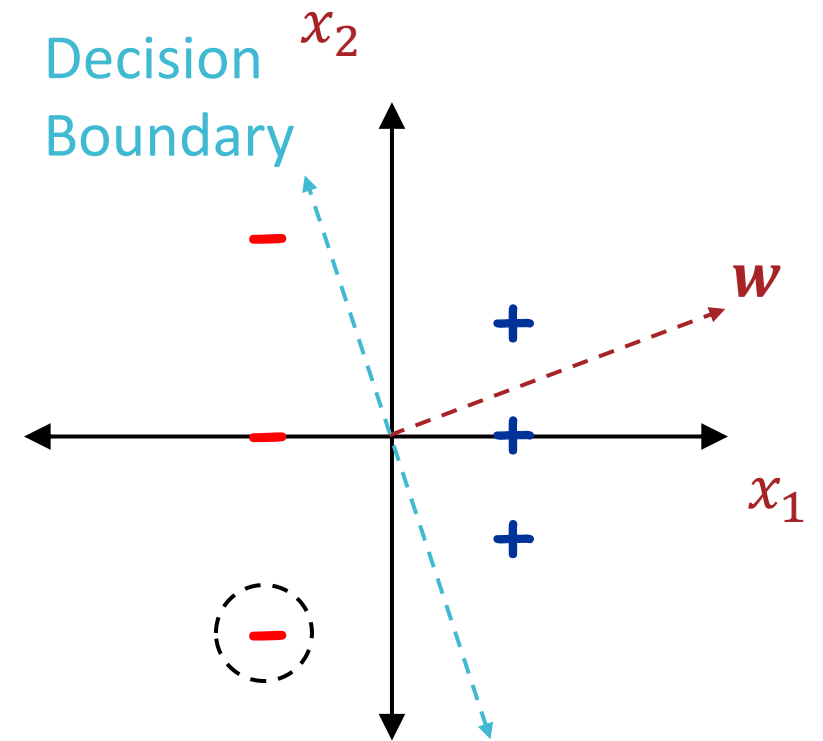


$$w = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$

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

(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



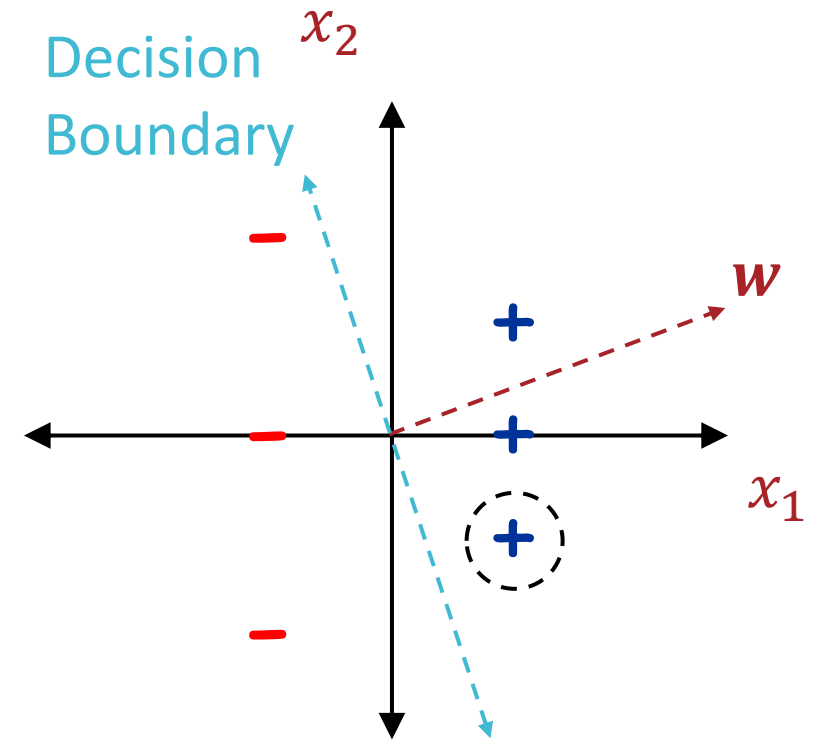
$$\mathbf{w} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$

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

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

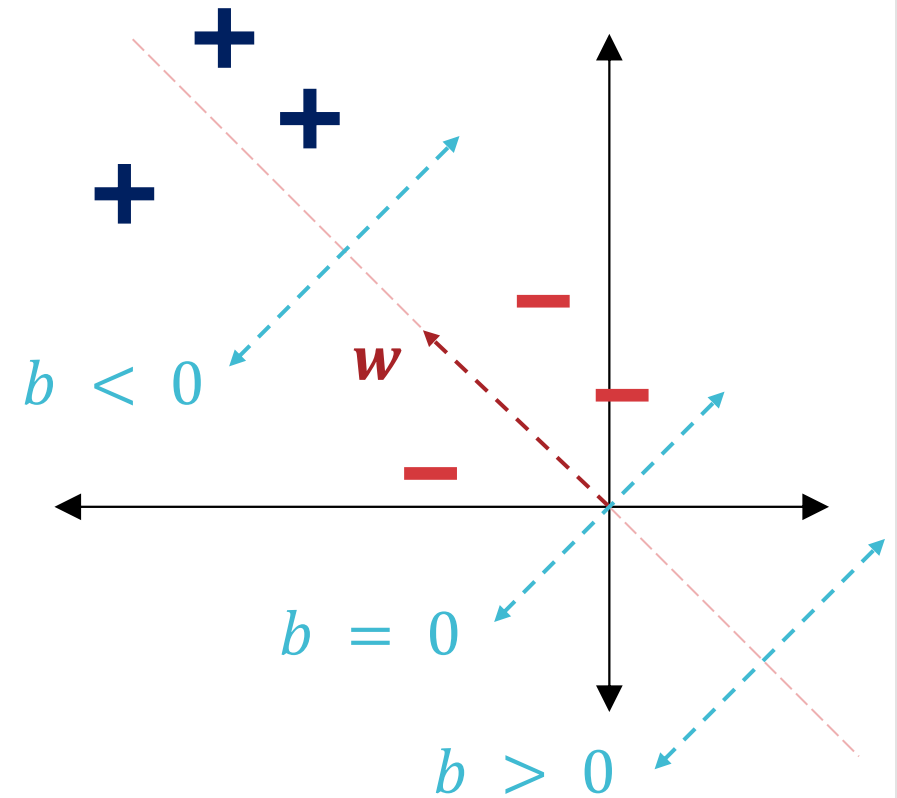
x_1	x_2	\hat{y}	y	Mistake?
-1	2	+	-	Yes
1	0	+	+	No
1	1	-	+	Yes
-1	0	-	-	No
-1	-2	+	-	Yes
1	-1	+	+	No

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



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



Poll Question 2:

- **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 the ~~training~~ validation data.
 - A. True
 - B. True and False (**TOXIC**)
 - C. False

Notational Hack

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

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

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

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

(Online) Perceptron Learning Algorithm

- Initialize the weight vector and intercept to all zeros:

$$\mathbf{w} = [0 \quad 0 \quad \dots \quad 0] \text{ and } b = 0$$

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

- Receive an unlabeled data point, $\mathbf{x}^{(t)}$

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

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

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

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

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

(Online) Perceptron Learning Algorithm

- Initialize the parameters to all zeros:

$$\boldsymbol{\theta} = [0 \quad 0 \quad \dots \quad 0]$$

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

- Receive an unlabeled data point, $\mathbf{x}^{(t)}$

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

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

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

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

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

Automatically handles
updating the intercept

(Online) Perceptron Learning Algorithm: Inductive Bias

- The decision boundary is linear and *recent mistakes are more important than older ones* (and should be corrected immediately)

(Online) Perceptron Learning Algorithm

- Initialize the parameters to all zeros:

$$\boldsymbol{\theta} = [0 \quad 0 \quad \dots \quad 0]$$

- For $t = 1, 2, 3, \dots$
 - Receive an unlabeled data point, $\mathbf{x}^{(t)}$
 - Predict its label, $\hat{y} = \text{sign}(\boldsymbol{\theta}^T \mathbf{x}'^{(t)})$
 - Observe its true label, $y^{(t)}$
 - If we misclassified a point ($y^{(t)} \neq \hat{y}$):
 - $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + y^{(t)} \mathbf{x}'^{(t)}$

(Batch) Perceptron Learning Algorithm

- Input: $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$
- Initialize the parameters to all zeros:

$$\boldsymbol{\theta} = [0 \quad 0 \quad \dots \quad 0]$$

- While NOT CONVERGED
 - For $t \in \{1, \dots, N\}$
 - Predict the label of $\mathbf{x}'^{(t)}$, $\hat{y} = \text{sign}(\boldsymbol{\theta}^T \mathbf{x}'^{(t)})$
 - Observe its true label, $y^{(t)}$
 - If we misclassified $\mathbf{x}'^{(t)}$ ($y^{(t)} \neq \hat{y}$):
 - $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + y^{(t)} \mathbf{x}'^{(t)}$

~~Poll Question 3:~~ (SKIPPED)

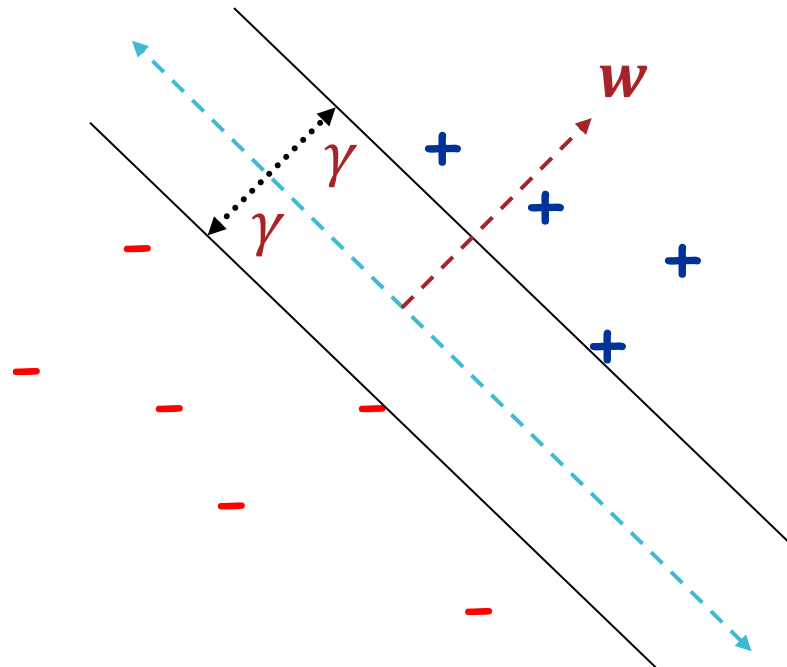
- **True or False:** The parameter vector \mathbf{w} learned by the batch Perceptron Learning Algorithm can be **written as a linear combination** of the examples, i.e.,

$$\mathbf{w} = c_1 \mathbf{x}^{(1)} + c_2 \mathbf{x}^{(2)} + \dots + c_N \mathbf{x}^{(M)}$$

- A. True and False (**TOXIC**)
- B. True
- C. False

Perceptron Mistake Bound

- Definitions:
 - A dataset \mathcal{D} is *linearly separable* if \exists a linear decision boundary that perfectly classifies the data points in \mathcal{D}
 - The margin, γ , of a dataset \mathcal{D} is the greatest possible distance between a linear separator and the closest data point in \mathcal{D} to that linear separator



Perceptron Mistake Bound

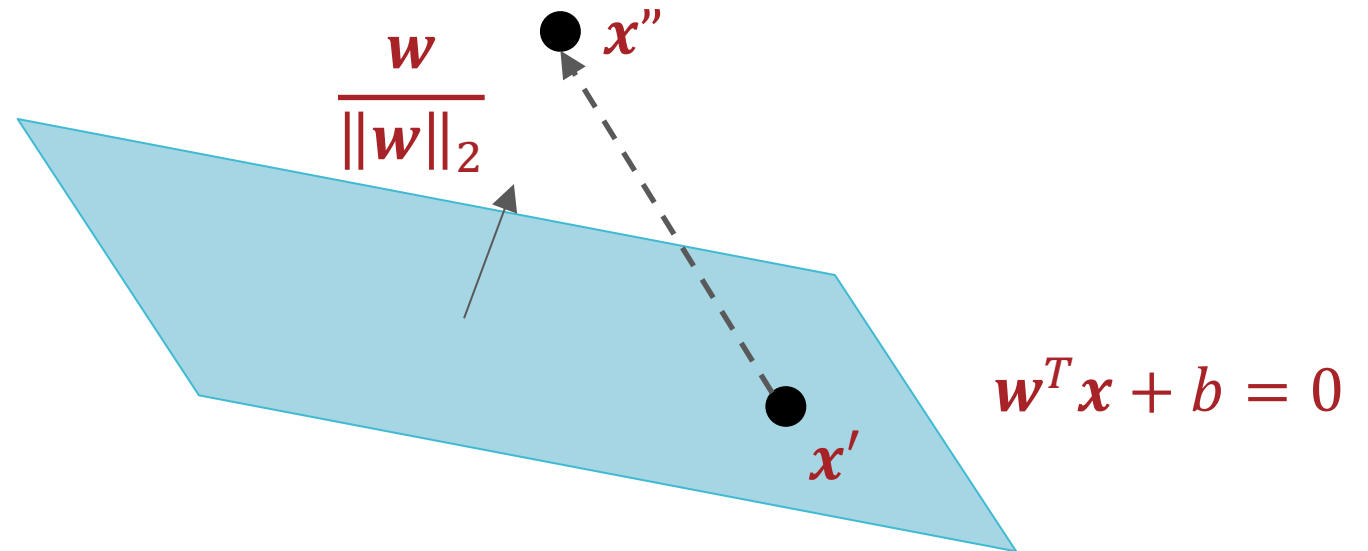
- Theorem: if the data points seen by the Perceptron Learning Algorithm (online and batch)
 1. lie in a ball of radius R (centered around the origin)
 2. have a margin of γ

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!

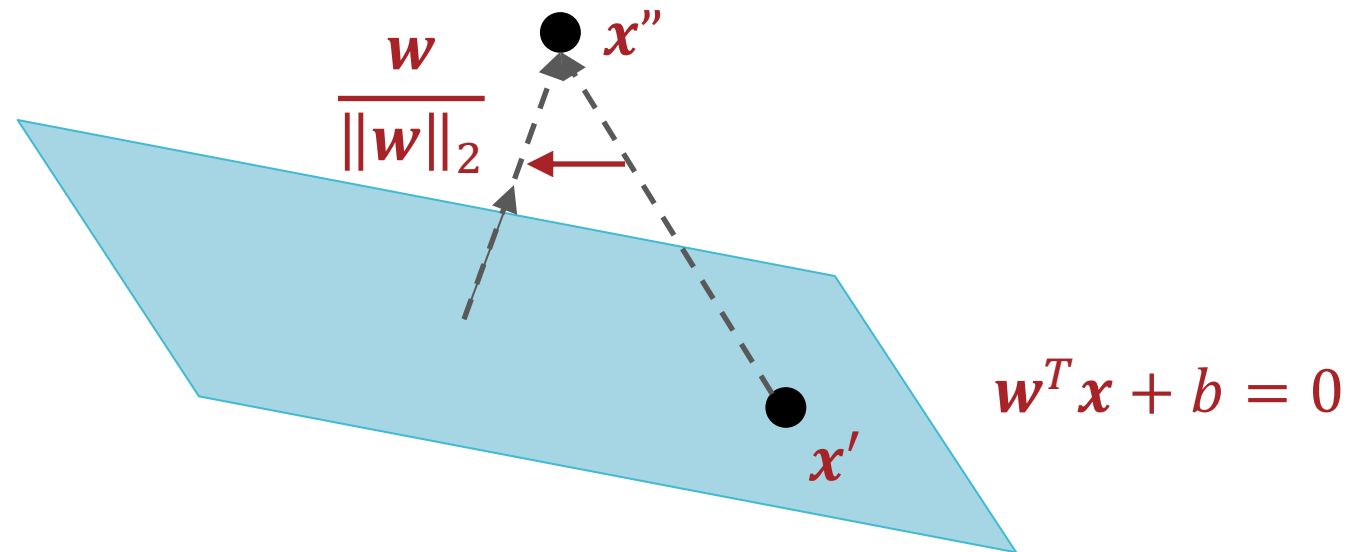
Computing the Margin

- Let \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 \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



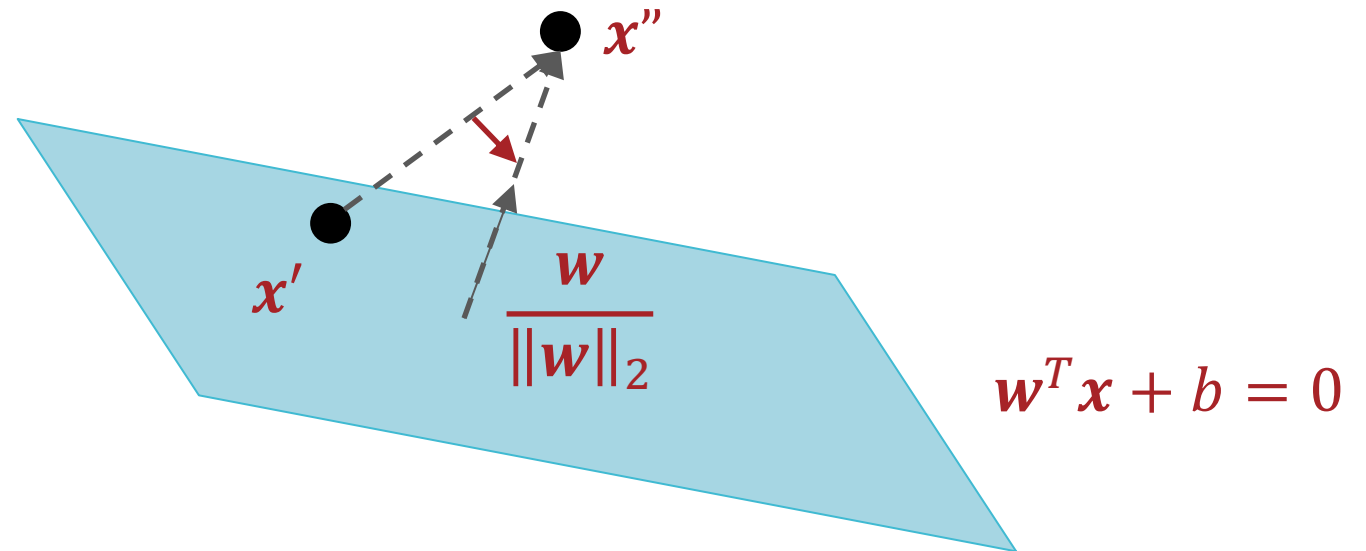
Computing the Margin

- 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 \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



Computing the Margin

- 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 \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



Computing the Margin

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

$$\left| \frac{\mathbf{w}^T (\mathbf{x}'' - \mathbf{x}')}{\|\mathbf{w}\|_2} \right| = \frac{|\mathbf{w}^T \mathbf{x}'' - \mathbf{w}^T \mathbf{x}'|}{\|\mathbf{w}\|_2} = \frac{|\mathbf{w}^T \mathbf{x}'' + b|}{\|\mathbf{w}\|_2}$$

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)