# Deep Learning (CS 470, CS 570)

**Module 2, Lecture 1: Perceptron Learning Algorithm**

# Toy Example: Animal Classification

We want to write an machine learning algorithm to classify each image into two animals classes; elephant and horse. This is a toy example for ease of understanding but we will soon see some real-world image classification examples.

**Dataset:** Contains image samples from each classes we want to detect. A part of these images are used for training and rest for testing.
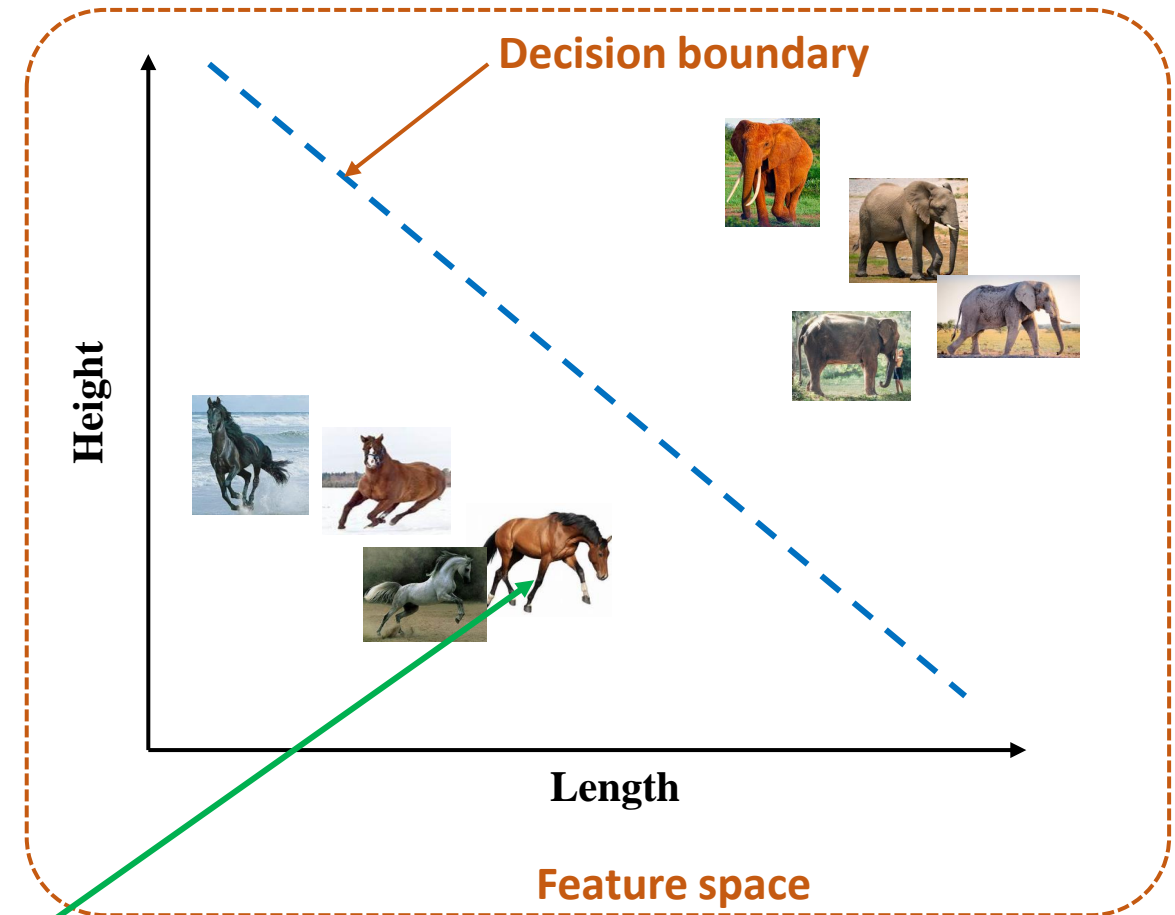


For any image we can use some image processing techniques to extract image features that encode crucial information useful for classification. We will later learn how to extract image features, but for this toy example we will assume that we know how to compute two features such as length and height of the animal present in the image.

**Decision boundary**

**Height**
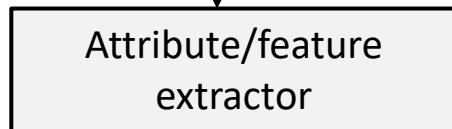
**Length**

**Feature space**

Becomes a point in a coordinate space where each axis of the coordinate space represents a feature. We refer it as data point or sample point. If we extract a $d$ dimensional feature vector, then the dimension of the coordinate space is also $d$ .

**[8ft 5ft]**

Any image

Represented by the length and height of the animal

# Toy Example: Animal Classification

Given enough training image samples, a machine learning algorithm learns to draw a boundary in the coordinate space (known as feature space) such that all horse images are in one side of the boundary (decision boundary) and elephant images are on the other side.

Once the decision boundary is learnt, the algorithm can predict the class of a new unknown image depending on which side of the decision boundary it belongs to. A typical ML algorithm framework is give below.

**Decision boundary**

**Height**

**Length**

**Feature space**

**Training data**

Attribute/feature extractor → Height, length → ML Algo → Predicted labels: elephant/horse

Labels: elephant/horse

# Taxonomy

**Supervised learning:**

**Input:** $x = (x_1, x_2, \ldots, x_d)$    Vector of dimension $d$

**Output:** y    Scalar
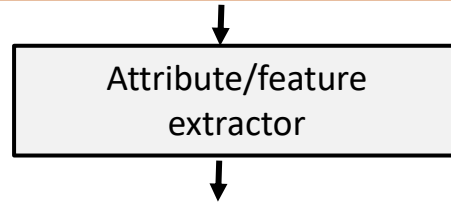
Continuous -> regression    Categorical -> classification

**Target function:** $f : x \rightarrow y$    Unknown

**Data:** $(x^1, y^1), (x^2, y^2), \ldots, (x^N, y^N)$    $N$ is the number of data points

**Hypothesis Set:** $H = \{h\}$    $H$ is a set that defines what are the possible decision boundaries a machine learning algorithm can generate. $h$ is a member of $H$.

**Learned Hypothesis:** $g: x \rightarrow y, \quad g \in H$    $g$ is an approximation of $f$

# Classification Framework



**Training data**

Attribute/feature extractor

**Data:** $(x^1, y^1), (x^2, y^2), \ldots, (x^N, y^N)$
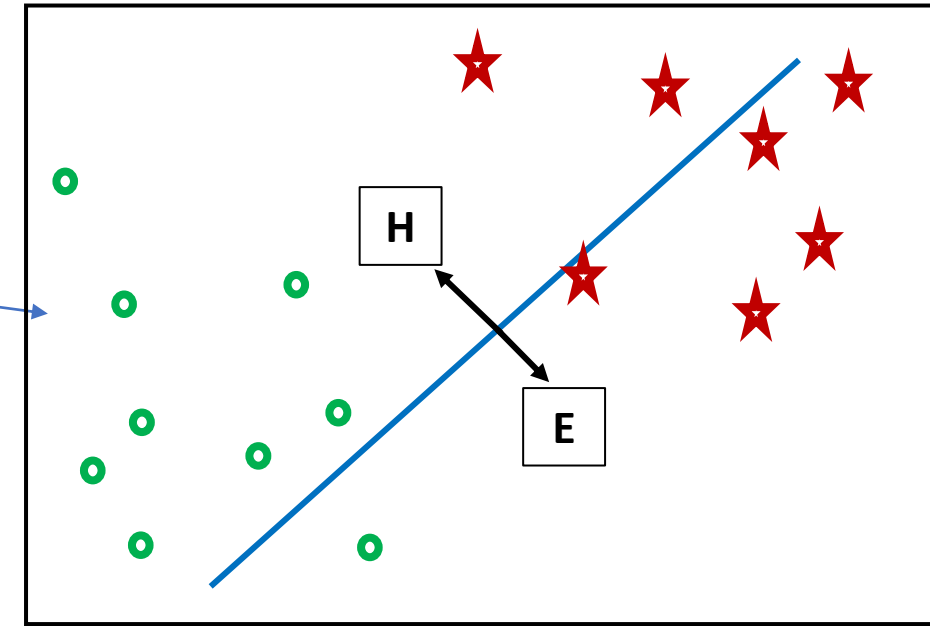
Plot data

**Hypothesis:**

Select a random line

Right side of the line is Elephant

Left side of the line is Horse

**Model training/learning/optimization:**
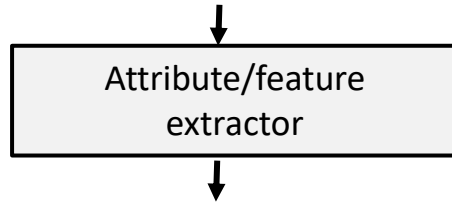
Change the line orientation

**H**

**E**

**Feature space**

○ Horse

★ Elephant

The above steps describes how a typical ML algorithm learns to classify two different group (here horse and elephant). A hypothesis h is a set of rules that defines the two classes in the feature space. The ML algorithm starts with a random hypothesis that consist of a decision boundary (here a line in arbitrary direction), and a few rules. During training, the algorithm learns from data examples how to draw a correct decision boundary (i.e. orientation and position of the line that divides the two groups of data points).

# Classification Framework

**Training data**

Attribute/feature extractor

**Data:** $(\boldsymbol{x}^1, \boldsymbol{y}^1), (\boldsymbol{x}^2, \boldsymbol{y}^2), \dots ,(\boldsymbol{x}^N, \boldsymbol{y}^N)$
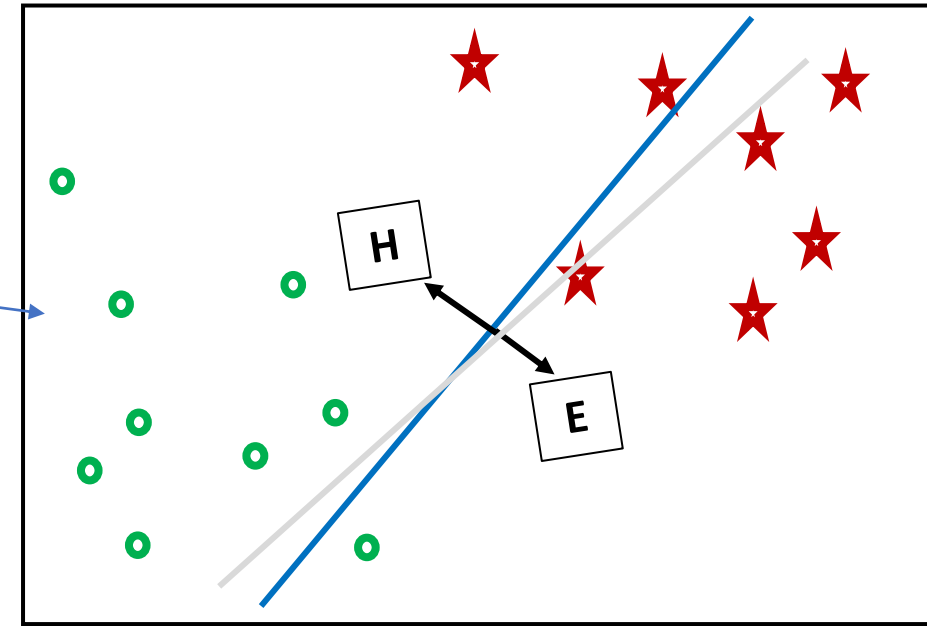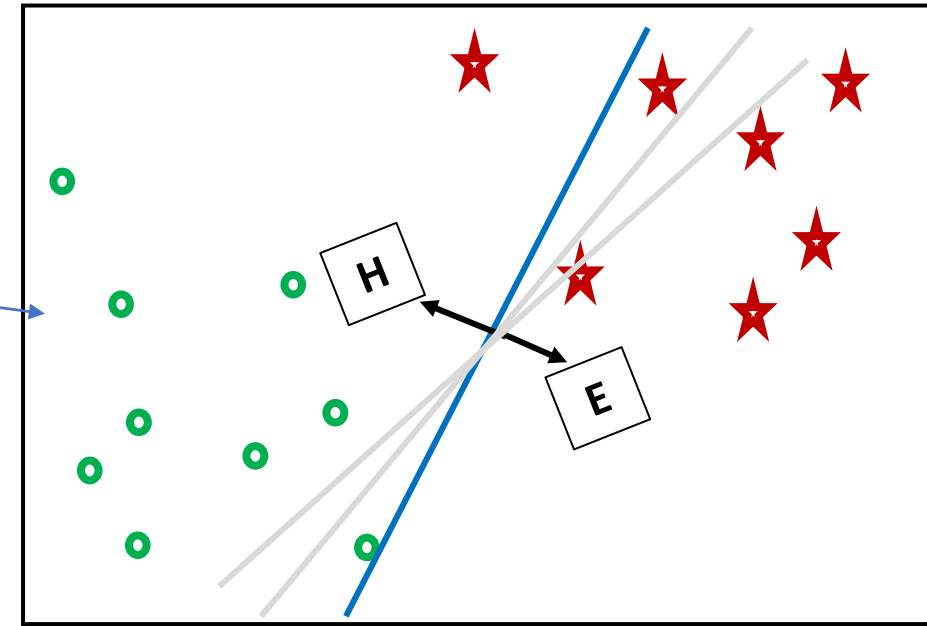
Plot data

**Hypothesis:**

Select a random line

Right side of the line is Elephant

Left side of the line is Horse

**Model training/optimization:**

Change the line orientation
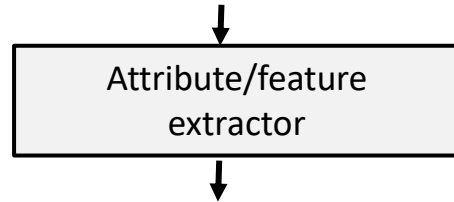
**Feature space**

○ Horse

★ Elephant

H

E

To learn how to draw the correct decision boundary, the ML algorithm runs a iterative process that slightly changes the line orientation (and position but the example doesn't show that) in every step following certain rules.

# Classification Framework



**Training data**

Attribute/feature extractor

**Data:** $(x^1, y^1), (x^2, y^2), \ldots, (x^N, y^N)$

Plot data

**Hypothesis:**

Select a random line

Right side of the line is Elephant

Left side of the line is Horse

**Feature space**

○ Horse

★ Elephant

**Model training/optimization:**

Change the line orientation

To learn how to draw the correct decision boundary, the ML algorithm runs a iterative process that slightly changes the line orientation (and position but the example doesn't show that) in every step following certain rules.

# Classification Framework



**Training data**

Attribute/feature extractor

**Data:** $(x^1, y^1), (x^2, y^2), \ldots, (x^N, y^N)$
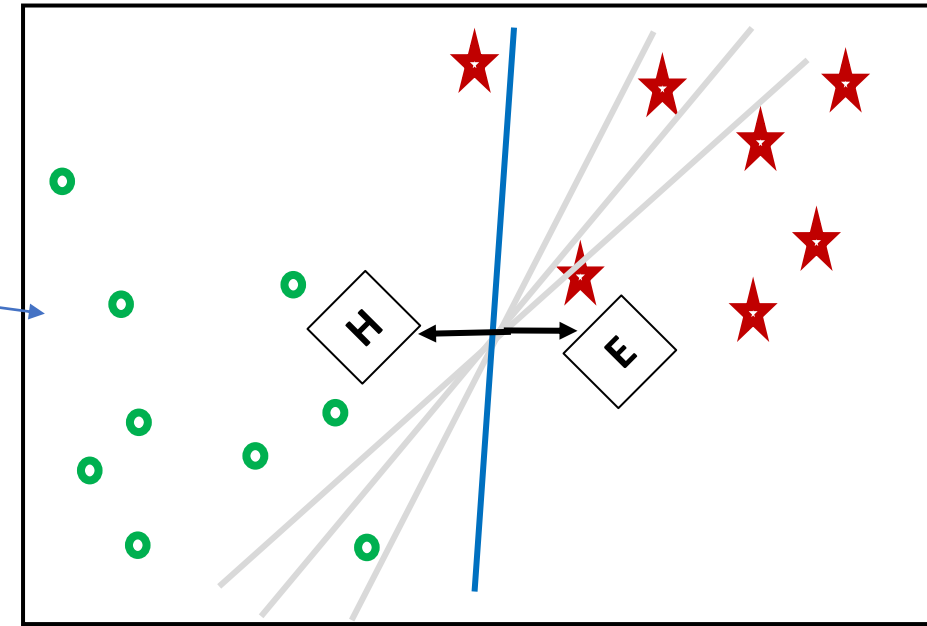
Plot data

**Hypothesis:**

Select a random line

Right side of the line is Elephant

Left side of the line is Horse

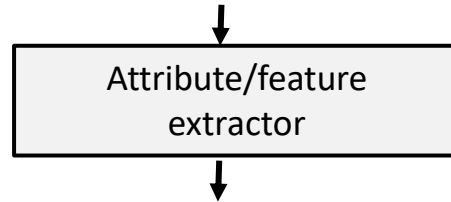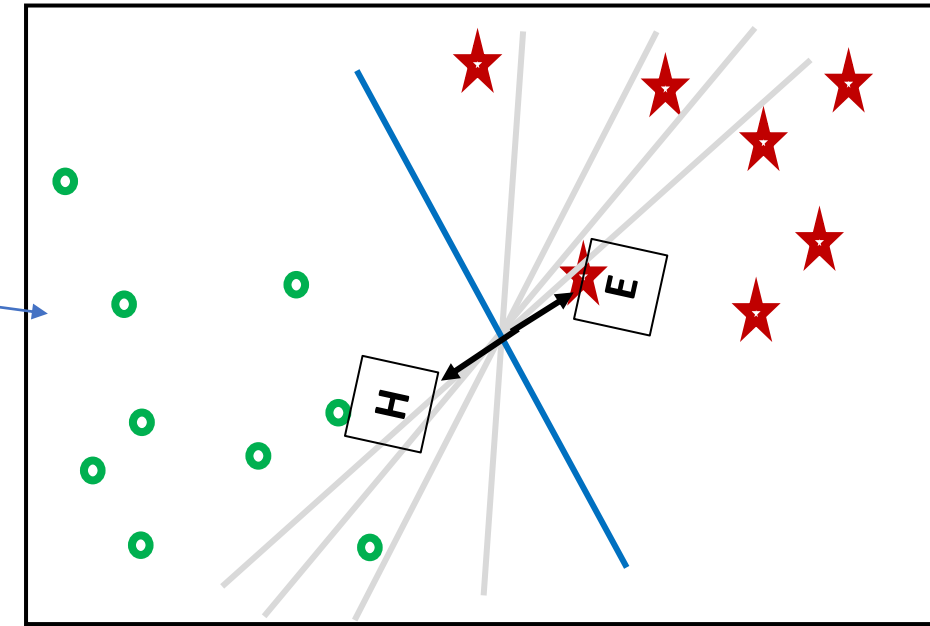**Model training/optimization:**

Change the line orientation

**Feature space**

○ Horse

★ Elephant

To learn how to draw the correct decision boundary, the ML algorithm runs a iterative process that slightly changes the line orientation (and position but the example doesn't show that) in every step following certain rules.

# Classification Framework

Attribute/feature
extractor

**Data:** $(\boldsymbol{x}^1, y^1), (\boldsymbol{x}^2, y^2), \dots, (\boldsymbol{x}^N, \boldsymbol{y}^N)$

Plot data

**Hypothesis:**

Select a random line

Right side of the line is Elephant

Left side of the line is Horse

**Model training/optimization:**

Change the line orientation

Select the line orientation that best classifies the training data

**Hypothesis set:**

Hypotheses associated with all line orientations



**Feature space**

○ Horse

★ Elephant

Finally, the algorithm finds the orientation (and position) that best classifies the two groups of training examples. This is when training stops.

# Perceptron Learning Algorithm: Hypothesis Set

**Input:** $\mathbf{x} = (x_1, x_2, \ldots, x_d)$

$x_i, i \in \{1, \ldots, d\}$ is a feature such as length and height

**Hypothesis:**

Image of an elephant if $\sum_{i=1}^{d} w_i x_i >$ threshold      Right side of line

Image of a horse if $\sum_{i=1}^{d} w_i x_i <$ threshold      Left side of line

**OR,**

$$h(x) = sign(\sum_{i=1}^{d} w_i x_i - threshold)$$      '+'ve : right side, '-'ve left side

**Hypothesis Set:**

$$H = \{ h(x) : \text{for all possible values of } w \text{ and threshold} \}$$

Here, we will learn the PLA algorithm that follows the learning technique that we just seen in the previous example. We will use a little math to understand the algorithm in details.

# Perceptron Learning Algorithm: Hypothesis Set

$$\sum_{i=1}^{d} w_i x_i - threshold$$

Is this equation familiar?

$$ax + by + c = 0$$  Equation of a line

$$ax + by + cz + d = 0$$  Equation of a plane

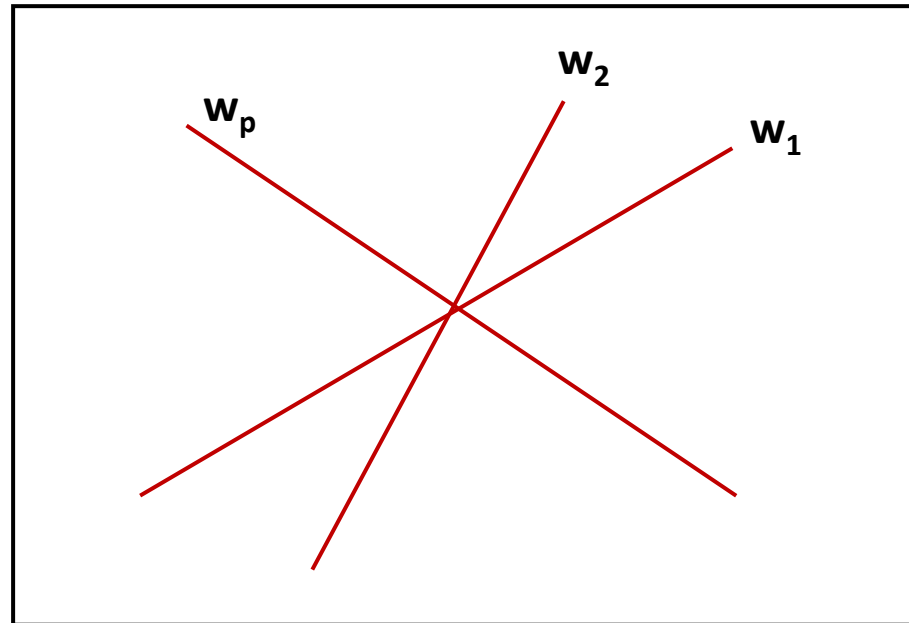$$\text{w}_1 x_1 + \text{w}_2 x_2 + \cdots + \text{w}_d x_d + threshold = 0$$  Equation of a hyperplane in a d dimensional space

**Hypothesis set:** Set of all possible hyper-planes

**Hypothesis:**

$$h(x) = sign(\sum_{i=0}^{d} w_i x_i)$$ , $x_0$ is always 1

$$h(x) = sign(\mathbf{w}^T \mathbf{x})$$  Vector form

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \ldots \\ w_d \end{bmatrix} \qquad \mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ \ldots \\ x_d \end{bmatrix}$$

# Perceptron Learning Algorithm: Learning

Random initialization of $w$

$w$ is perpendicular to the hyperplane $w^Tx$     WHY?

At any step, t-1, of the algorithm, pick a misclassified point

$$h(x) = sign(w_{t\_1}{}^T x_n) \neq sign(y_n)$$

Update the weight vector as:

$$w_t \leftarrow w_{t-1} + y_n x_n \quad \text{Update rule}$$

# Perceptron Learning Algorithm: Learning

Random initialization of $w$

$w$ is perpendicular to the hyperplane $w^Tx$     WHY?

At any step, t-1, of the algorithm, pick a misclassified point

$$h(x) = sign(w_{t\_1}^T x_n) \neq sign(y_n)$$

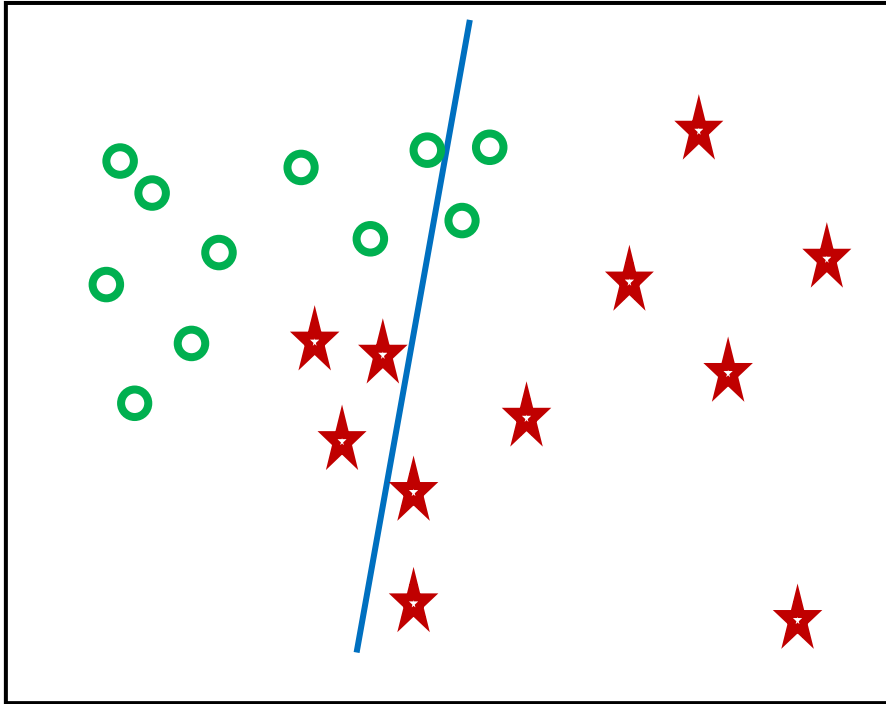Update the weight vector as:

$$w_t \leftarrow w_{t-1} + y_n x_n \quad \text{Update rule}$$



Continue picking points and updating $w$ as long there are any misclassified points.

When data is linearly separable, algorithm will find a $w_{final}$ that will classify all training points correctly.

$$g(x) = sign(w_{final}^T x)$$

# When Data is not Linearly Separable
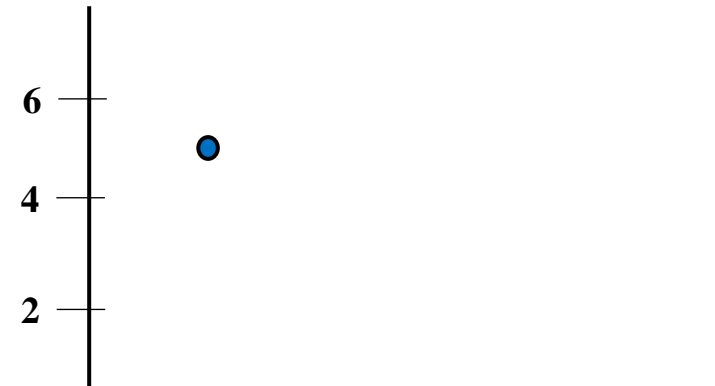


Random initialization of $w$

Estimate classification error. Save error and $w$

For a fixed number of iterations
    Update $w$
    Re-estimate error
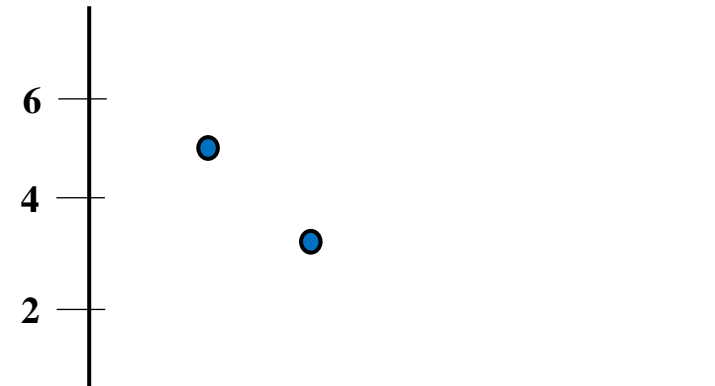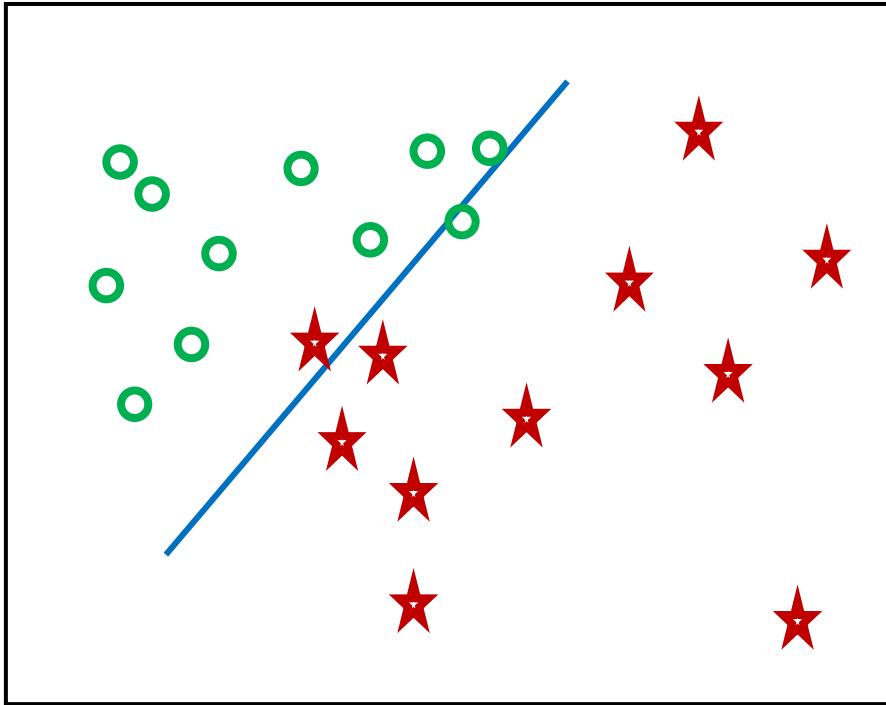    Save the error rate and $w$

At any stage of the algorithm, the classification error is the percentage of misclassified sample points.

# When Data is not Linearly Separable



Random initialization of $w$

Estimate classification error. Save error and $w$

For a fixed number of iterations
 Update $w$
 Re-estimate error
 Save the error rate and $w$

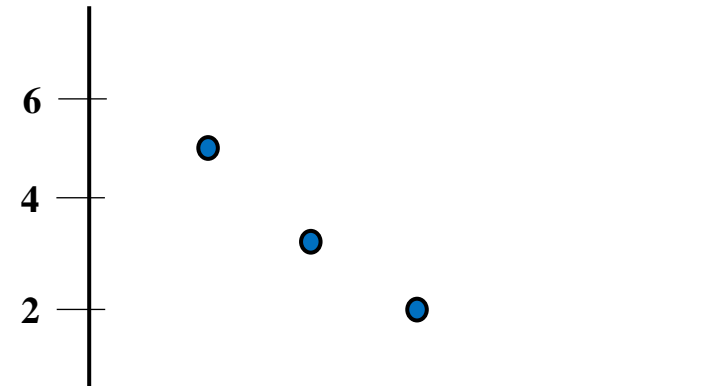# When Data is not Linearly Separable



Random initialization of $w$

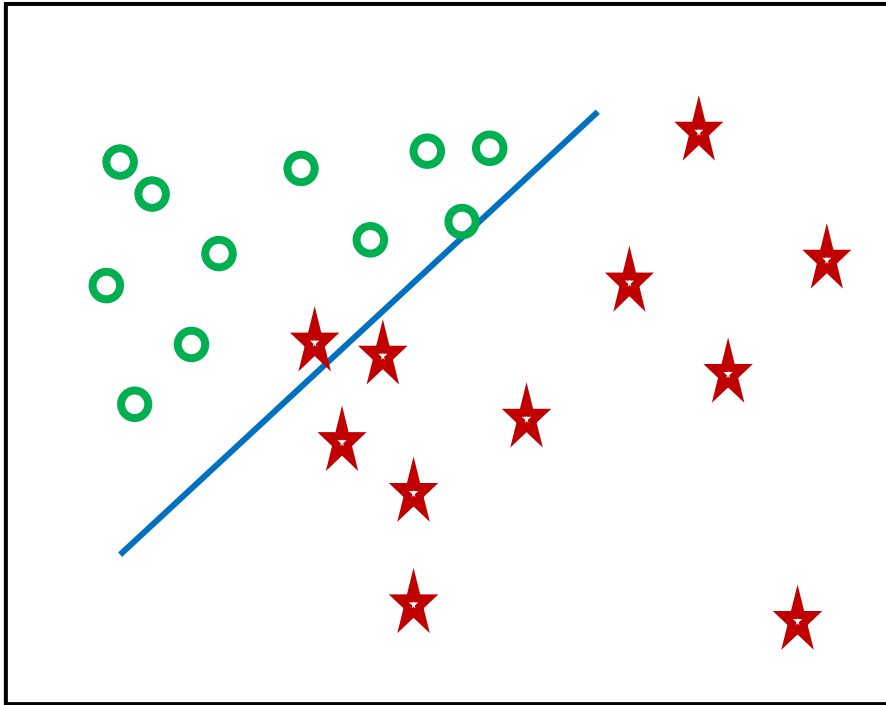Estimate classification error. Save error and $w$

For a fixed number of iterations
    Update $w$
    Re-estimate error
    Save the error rate and $w$

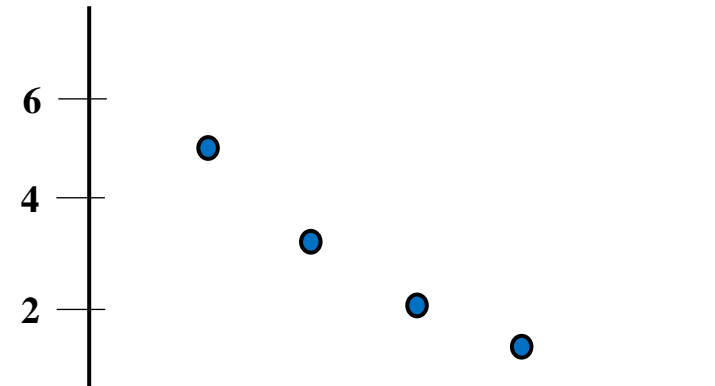# When Data is not Linearly Separable



Random initialization of $w$

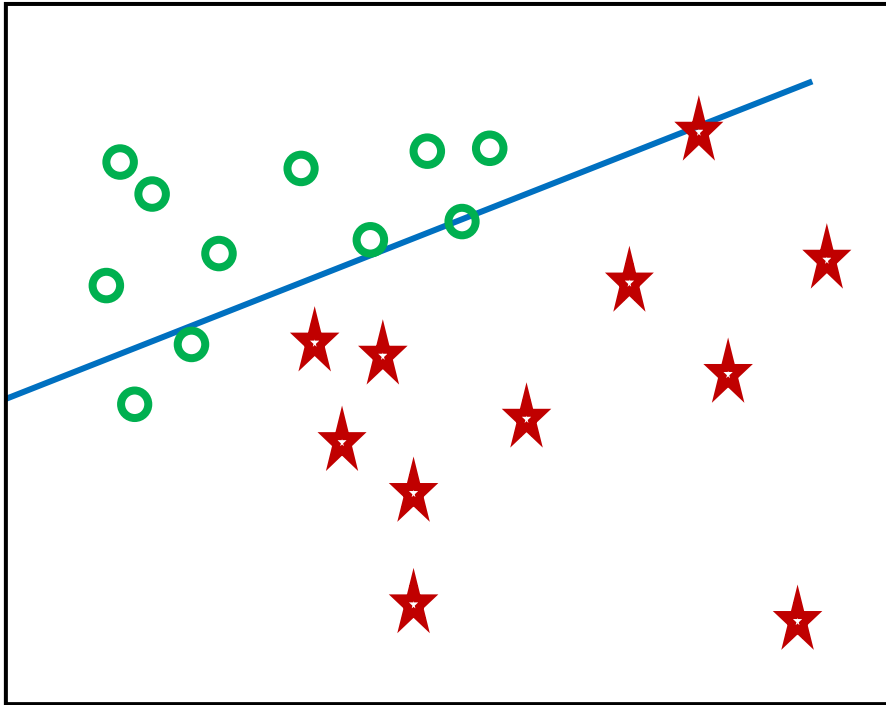Estimate classification error. Save error and $w$

For a fixed number of iterations
      Update $w$
      Re-estimate error
      Save the error rate and $w$

# When Data is not Linearly Separable

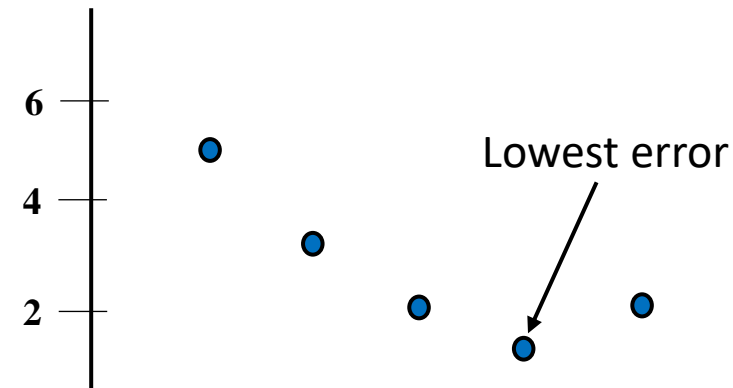Random initialization of $w$

Estimate classification error. Save error and $w$

For a fixed number of iterations
    Update $w$
    Re-estimate error
    Save the error rate and $w$

At the end, report lowest error and corresponding $w$

# Additional Reading

PLA [description and code](#)