

Binary Classification

- We will focus on the binary classification problem to build up our intuition about machine learning.
- Binary classification is closely related to binary hypothesis testing, except that we are not provided with conditional probability models for the observation.
- Instead, we have a dataset consisting of labeled examples.
- While we can use these examples to estimate the probability models, our goal is to design a **decision rule** to classify new examples.
 - Ex: Design a classifier for cat and dog images.
 - Ex: Classify MRI tumor images as malignant or benign.
 - Ex: Automatically sort harvested fruit into quality bins.

• Binary Classification Framework:

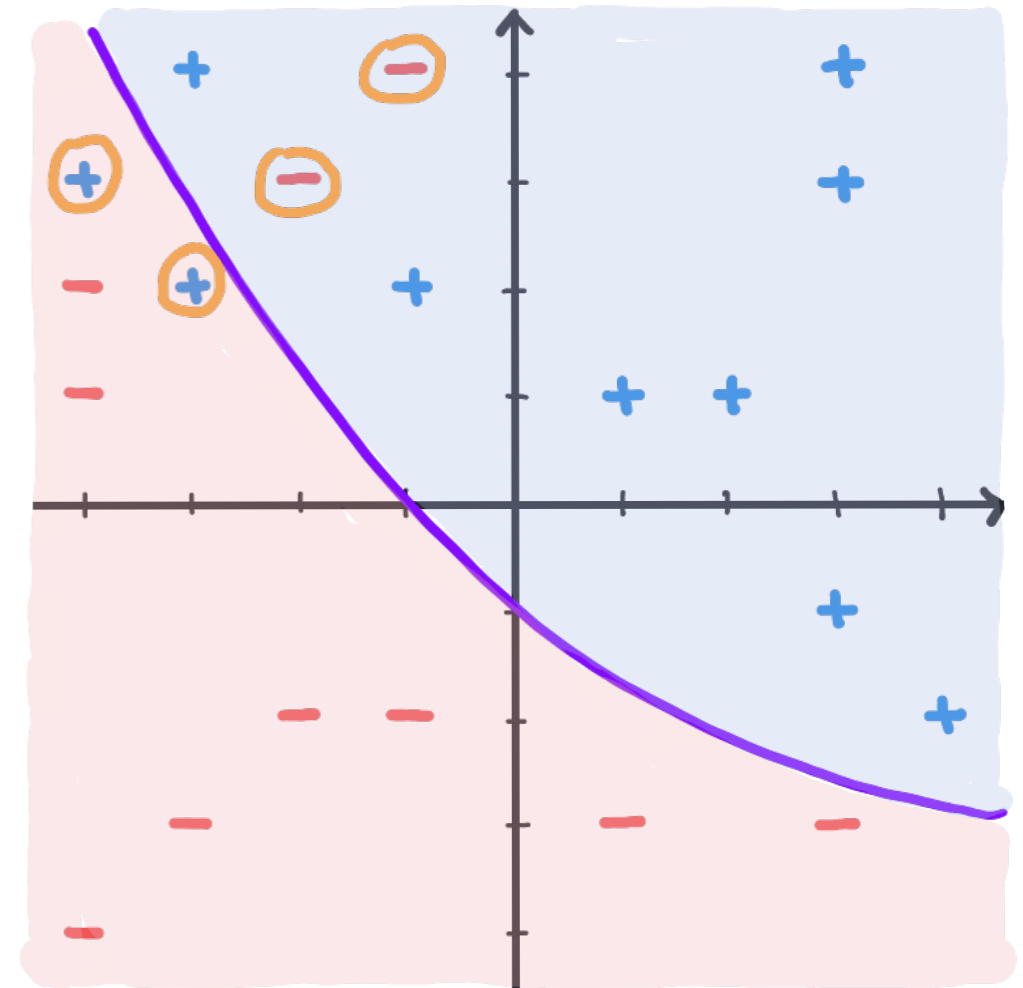
→ We have a **dataset** consisting of a collection of labeled examples of the form (\underline{x}_i, y_i) for $i=1, \dots, n$ where

* $\underline{x}_i = \begin{bmatrix} x_{i,1} \\ \vdots \\ x_{i,d} \end{bmatrix}$ is the d -dimensional feature vector

* $y_i \in \{+1, -1\}$ is the label

→ Our task is to design a **classifier** (or decision rule) $D(\underline{x})$ that outputs either $+1$ or -1 for any possible input vector.

→ We measure performance using the **error rate**, which is the fraction of misclassified points.



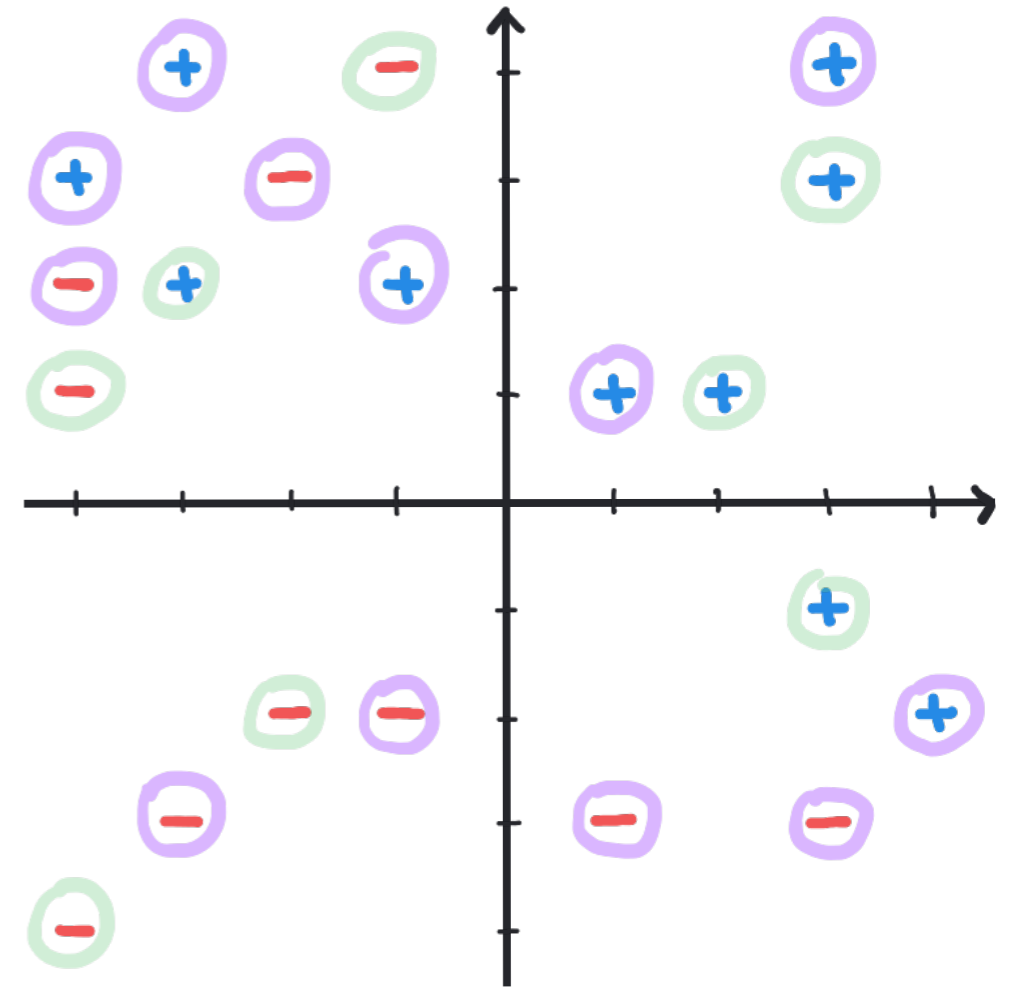
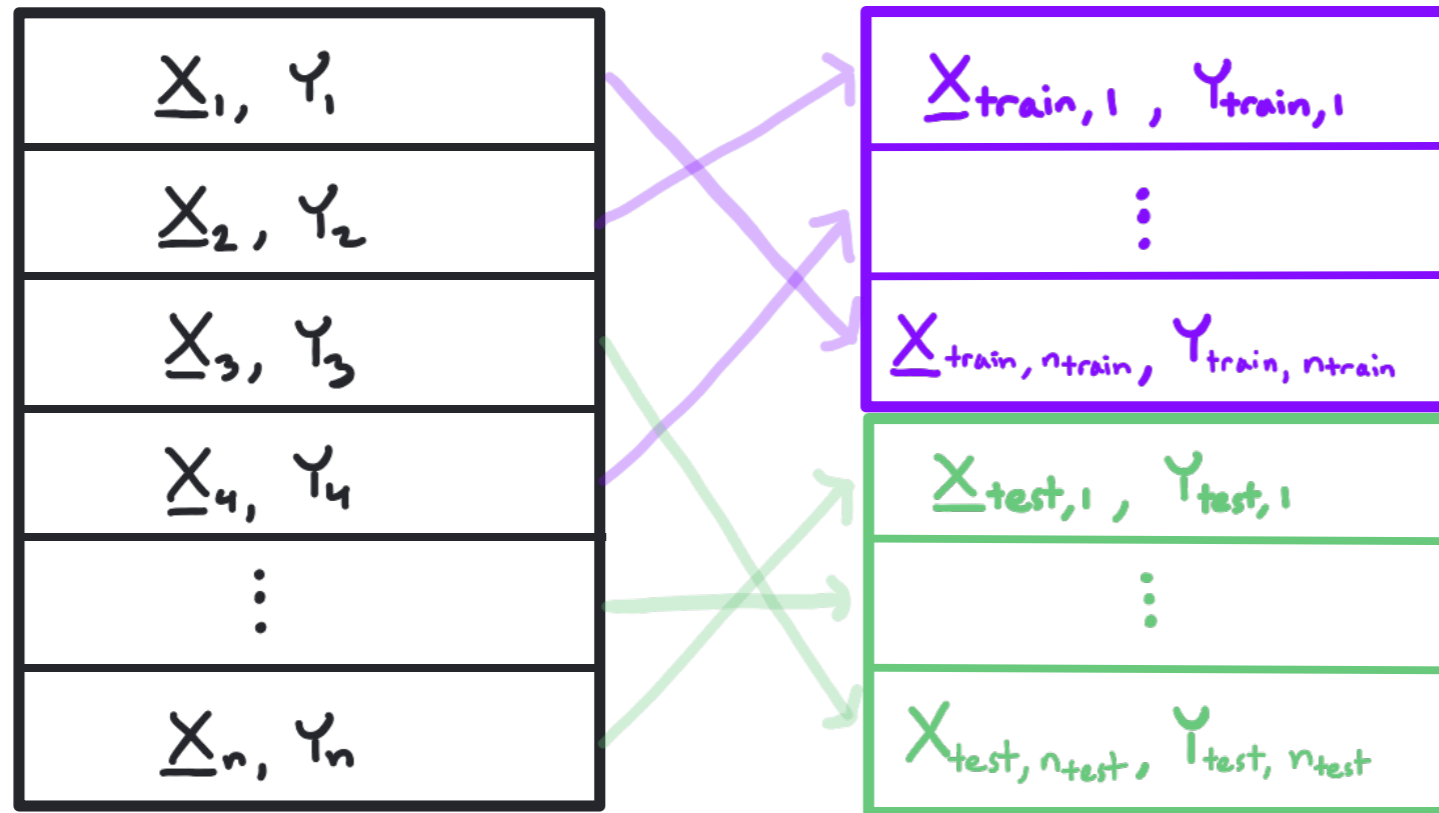
2-dimensional example dataset

examples $n = 20$

misclassified = 4

Error Rate = $\frac{4}{20} = 20\%$

- Recall that we need to split our dataset into **training data** and **test data**.



- The **training data** $(\underline{X}_{\text{train},i}, Y_{\text{train},i})$, $i = 1, \dots, n_{\text{train}}$ is used to design a good decision rule and tune its parameters.
- The **test data** $(\underline{X}_{\text{test},i}, Y_{\text{test},i})$, $i = 1, \dots, n_{\text{test}}$ is used to estimate the error rate and **should not be used for parameter tuning**.
- Note that $n_{\text{train}} + n_{\text{test}} = n$.

- It will be helpful to define a function to count errors:

$$g_{\text{error}}(D(\underline{x}_i), y_i) = \begin{cases} 1, & D(\underline{x}_i) \neq y_i & \underline{x}_i: \text{misclassified} \\ 0, & D(\underline{x}_i) = y_i & \underline{x}_i: \text{correctly classified} \end{cases}$$

- The **training error rate** is the fraction of misclassified **training examples**:

$$\begin{aligned} \text{Err}_{\text{train}} &= \frac{\# \text{ misclassified training examples}}{\# \text{ total training examples}} \\ &= \frac{1}{n_{\text{train}}} \sum_{i=1}^{n_{\text{train}}} g_{\text{error}}(D(\underline{x}_{\text{train},i}), Y_{\text{train},i}) \end{aligned}$$

- The **test error rate** is the fraction of misclassified **test examples**:

$$\begin{aligned} \text{Err}_{\text{test}} &= \frac{\# \text{ misclassified test examples}}{\# \text{ total test examples}} \\ &= \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} g_{\text{error}}(D(\underline{x}_{\text{test},i}), Y_{\text{test},i}) \end{aligned}$$

- Remember we can get to **0% training error rate** by memorizing dataset but this may not lead to a low **test error rate**.

• Nearest Neighbor Classifier:

→ Basic Idea: Find the training example that is closest to the input vector, and output its label as the guess.

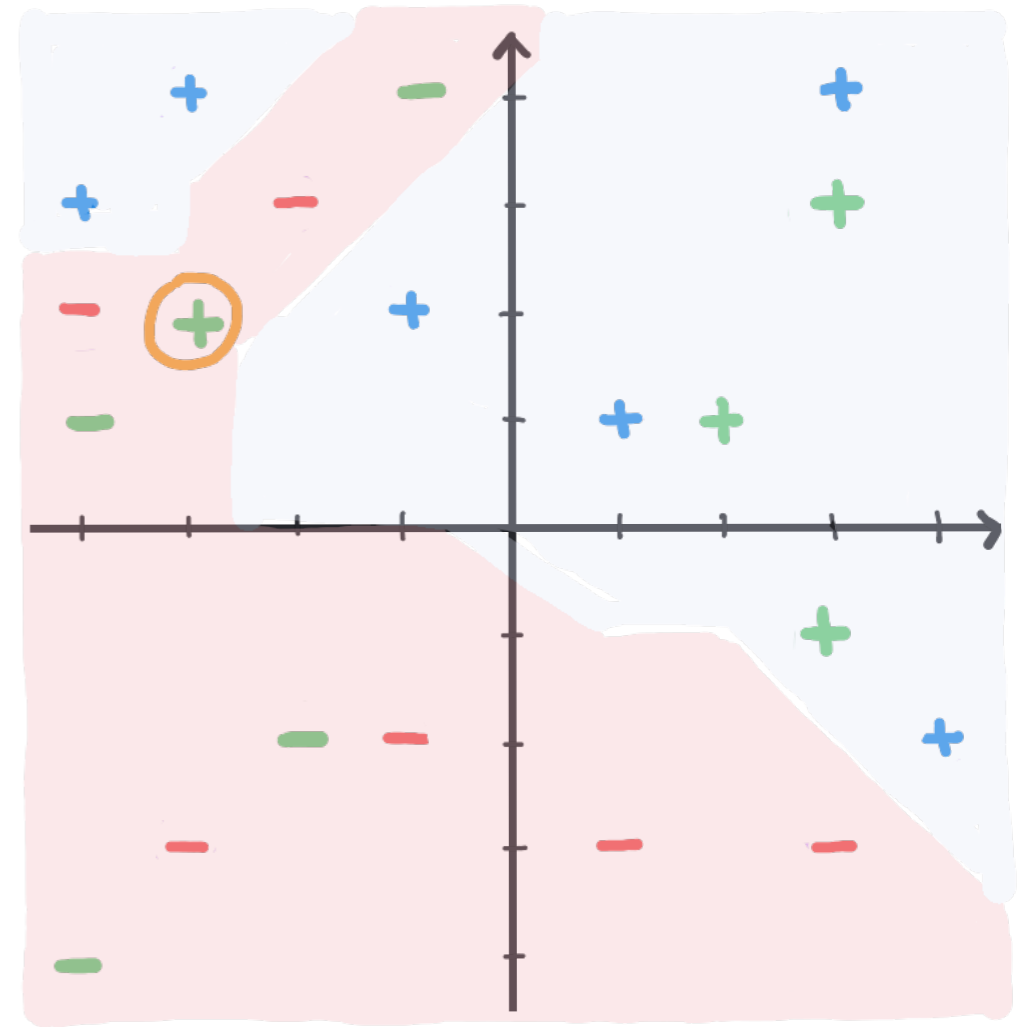
$$D_{NN}(\underline{x}) = Y_{\text{train}, i_{\text{closest}}}$$

where i_{closest} is the index of the closest training example

$$\| \underline{x} - X_{\text{train}, i_{\text{closest}}} \| = \min_{i=1, \dots, n_{\text{train}}} \| \underline{x} - X_{\text{train}, i} \|$$

$$\| \underline{x} \| = \sqrt{\sum_{j=1}^d x_j^2} \text{ is the length of } \underline{x}$$

→ We can also consider a k-Nearest Neighbor Classifier that finds the k training examples closest to the input vector and outputs the majority vote amongst their labels as its guess.



$$n_{\text{train}} = 12$$

$$\# \text{ Training Errors} = 0$$

$$\text{Training Error Rate} = 0\%$$

$$n_{\text{test}} = 8$$

$$\# \text{ Test Errors} = 1$$

$$\text{Test Error Rate} = \frac{1}{8} = 12.5\%$$

• Closest Average Classifier:

→ Basic Idea: Compute the conditional sample mean vector of the training data given $Y = +1$ and $Y = -1$.

Find the closest mean vector to the input vector and output its label.

$n_{\text{train},+}$ = # training examples with $Y_i = +1$

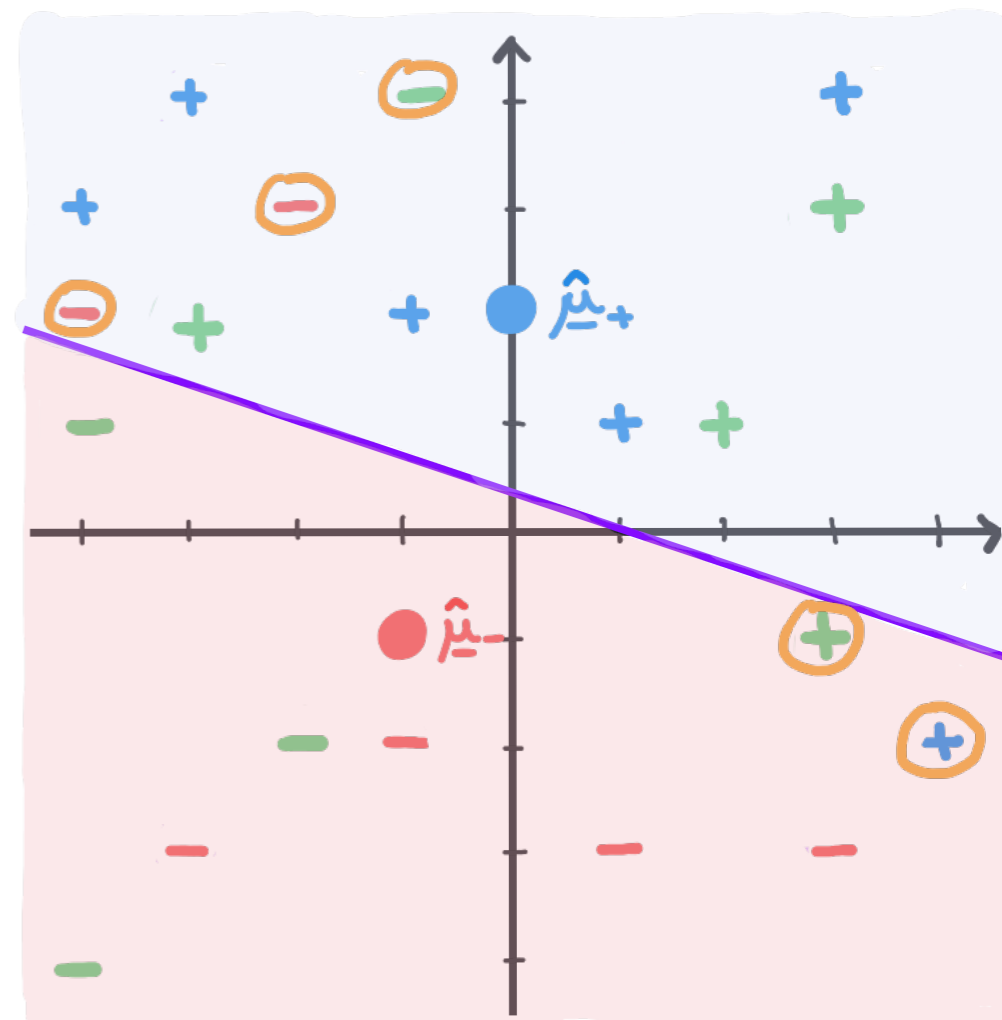
$n_{\text{train},-}$ = # training examples with $Y_i = -1$

$$L_+ = \{ i \in \{1, \dots, n_{\text{train}}\} : Y_{\text{train},i} = +1 \}$$

$$L_- = \{ i \in \{1, \dots, n_{\text{train}}\} : Y_{\text{train},i} = -1 \}$$

$$\hat{\mu}_+ = \frac{1}{n_{\text{train},+}} \sum_{i \in L_+} \underline{x}_{\text{train},i} \quad \hat{\mu}_- = \frac{1}{n_{\text{train},-}} \sum_{i \in L_-} \underline{x}_{\text{train},i}$$

$$D_{\text{avg}}(\underline{x}) = \begin{cases} +1 & \|\underline{x} - \hat{\mu}_+\| \leq \|\underline{x} - \hat{\mu}_-\| \\ -1 & \|\underline{x} - \hat{\mu}_-\| < \|\underline{x} - \hat{\mu}_+\| \end{cases}$$



$n_{\text{train}} = 12$

Training Errors = 3

Training Error Rate = $\frac{3}{12} = 25\%$

$n_{\text{test}} = 8$

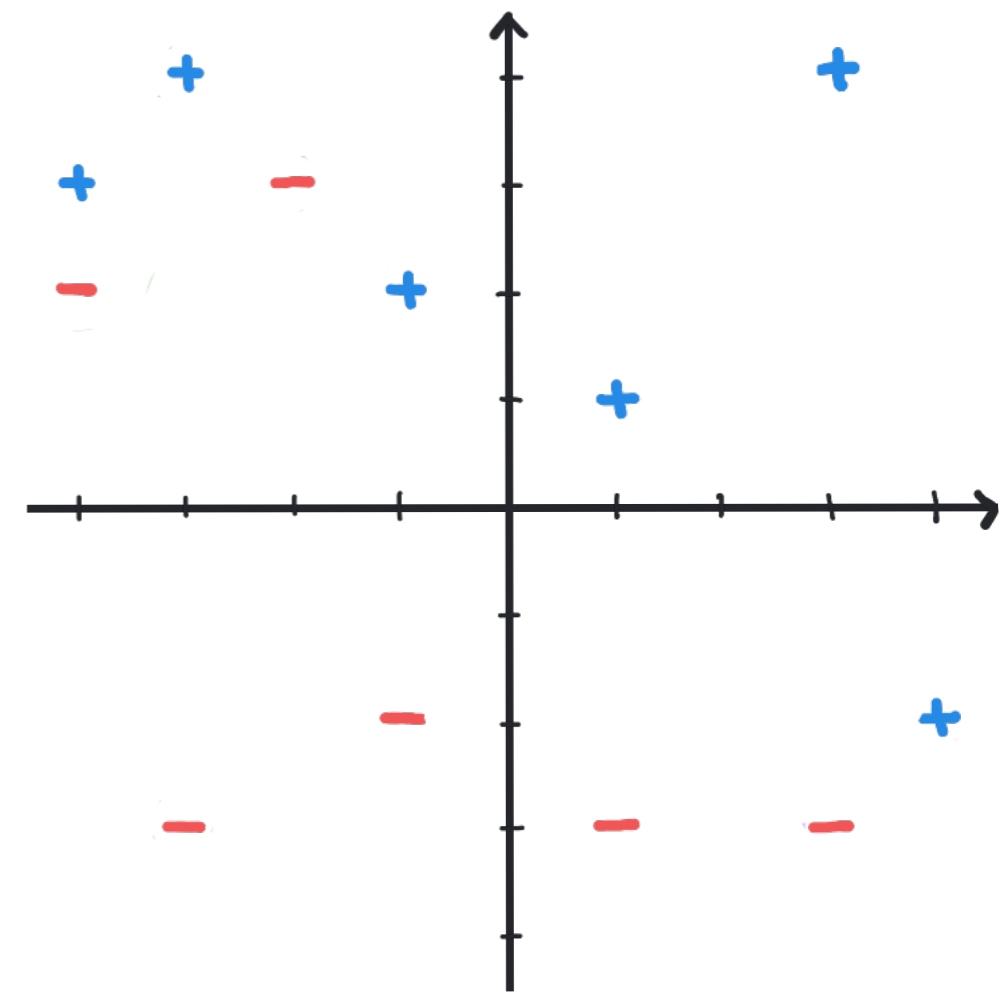
Test Errors = 2

Test Error Rate = $\frac{2}{8} = 25\%$

• Linear Discriminant Analysis:

→ Basic Idea: Assume that, given the label, the input vector is Gaussian with mean vector μ_+ or μ_- and the same covariance matrix Σ . Estimate these parameters from the data and apply the ML rule.

→ Define $n_{train,+}$, $n_{train,-}$, L_+ , L_- , $\hat{\mu}_+$, $\hat{\mu}_-$ as on previous slide.



$$\hat{\Sigma}_+ = \frac{1}{n_{train,+} - 1} \sum_{i \in L_+} (\underline{x}_{train,i} - \hat{\mu}_+) (\underline{x}_{train,i} - \hat{\mu}_+)^T$$

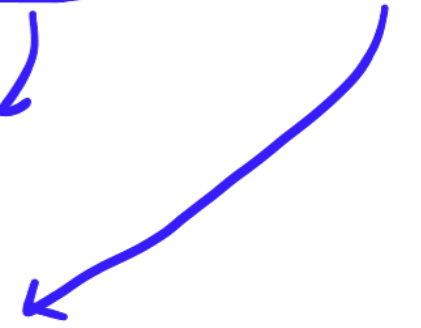
$$\hat{\Sigma}_- = \frac{1}{n_{train,-} - 1} \sum_{i \in L_-} (\underline{x}_{train,i} - \hat{\mu}_-) (\underline{x}_{train,i} - \hat{\mu}_-)^T$$

$$\hat{\Sigma} = \frac{1}{n_{train} - 2} \left((n_{train,+} - 1) \hat{\Sigma}_+ + (n_{train,-} - 1) \hat{\Sigma}_- \right)$$

Conditional PDFs of \underline{x} given $Y = +1$ and -1

$$\hat{f}_{\underline{x}|Y}(\underline{x} | Y = +1) = \frac{1}{\sqrt{(2\pi)^d \det(\hat{\Sigma})}} \exp\left(-\frac{1}{2} (\underline{x} - \hat{\mu}_+)^T \hat{\Sigma}^{-1} (\underline{x} - \hat{\mu}_+)\right)$$

$$\hat{f}_{\underline{x}|Y}(\underline{x} | Y = -1) = \frac{1}{\sqrt{(2\pi)^d \det(\hat{\Sigma})}} \exp\left(-\frac{1}{2} (\underline{x} - \hat{\mu}_-)^T \hat{\Sigma}^{-1} (\underline{x} - \hat{\mu}_-)\right)$$



• Linear Discriminant Analysis:

$$\hat{f}_{x|y}(x|y=+1) = \frac{1}{\sqrt{(2\pi)^d \det(\hat{\Sigma})}} \exp\left(-\frac{1}{2} (x - \hat{\mu}_+)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_+)\right)$$

$$\hat{f}_{x|y}(x|y=-1) = \frac{1}{\sqrt{(2\pi)^d \det(\hat{\Sigma})}} \exp\left(-\frac{1}{2} (x - \hat{\mu}_-)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_-)\right)$$

ML Rule: Guess +1 if $\hat{f}_{x|y}(x|y=+1) \geq \hat{f}_{x|y}(x|y=-1)$ or, equivalently,

$$\begin{aligned} (x - \hat{\mu}_+)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_+) &\leq (x - \hat{\mu}_-)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_-) \\ \cancel{x^T \hat{\Sigma}^{-1} x} - 2 \hat{\mu}_+^T \hat{\Sigma}^{-1} x + \hat{\mu}_+^T \hat{\Sigma}^{-1} \hat{\mu}_+ &\leq \cancel{x^T \hat{\Sigma}^{-1} x} - 2 \hat{\mu}_-^T \hat{\Sigma}^{-1} x + \hat{\mu}_-^T \hat{\Sigma}^{-1} \hat{\mu}_- \\ \hat{\mu}_+^T \hat{\Sigma}^{-1} \hat{\mu}_+ - \hat{\mu}_-^T \hat{\Sigma}^{-1} \hat{\mu}_- &\leq 2 (\hat{\mu}_+ - \hat{\mu}_-)^T \hat{\Sigma}^{-1} x \end{aligned}$$

$c \leq \underline{b}^T x$ Linear Classifier!

$$D_{LDA}(x) = \begin{cases} +1 & 2(\hat{\mu}_+ - \hat{\mu}_-)^T \hat{\Sigma}^{-1} x \geq \hat{\mu}_+^T \hat{\Sigma}^{-1} \hat{\mu}_+ - \hat{\mu}_-^T \hat{\Sigma}^{-1} \hat{\mu}_- \\ -1 & \text{otherwise} \end{cases}$$

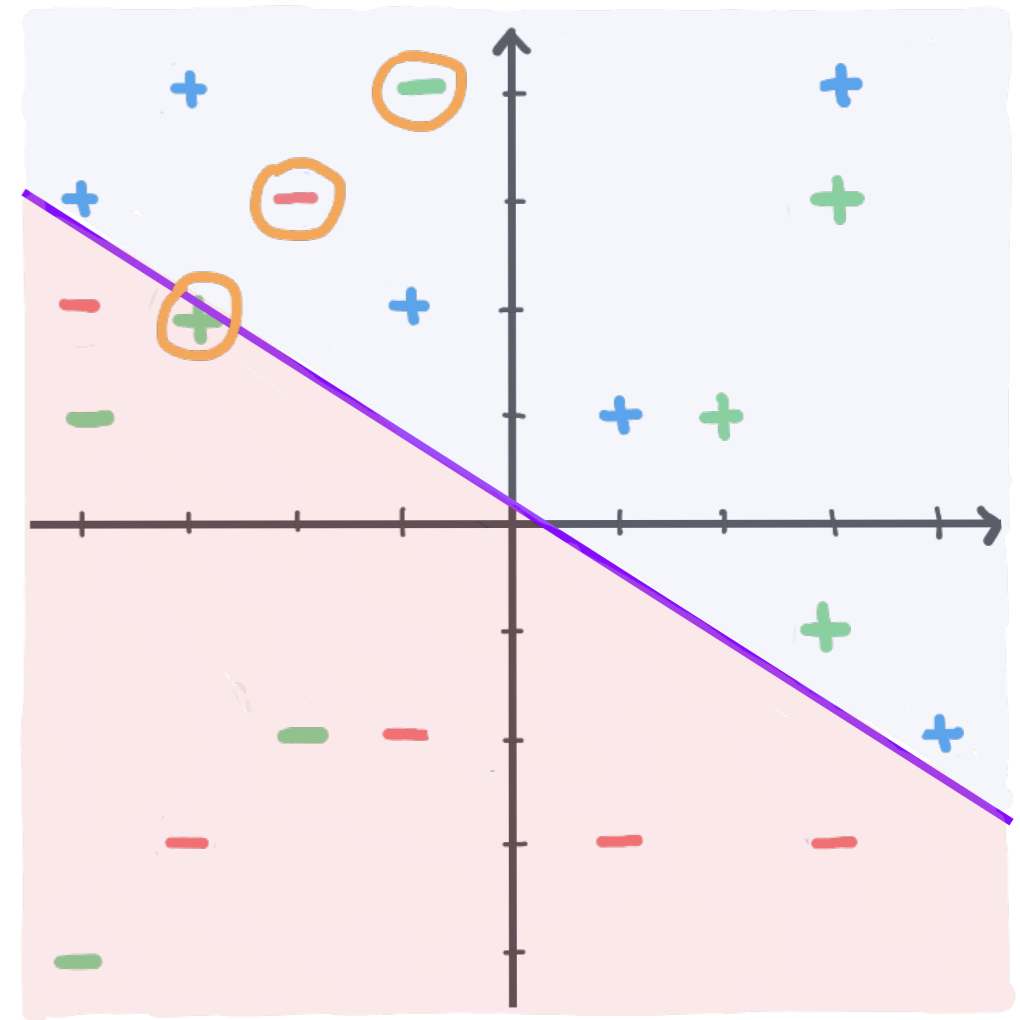
→ Closest average classifier can be viewed as a special case with the assumption $\Sigma = \mathbf{I}$, the identity matrix.

• Linear Discriminant Analysis:

→ Basic Idea: Assume that, given the label, the input vector is Gaussian with mean vector μ_+ or μ_- and the same covariance matrix Σ . Estimate these parameters from the data and apply the ML rule.

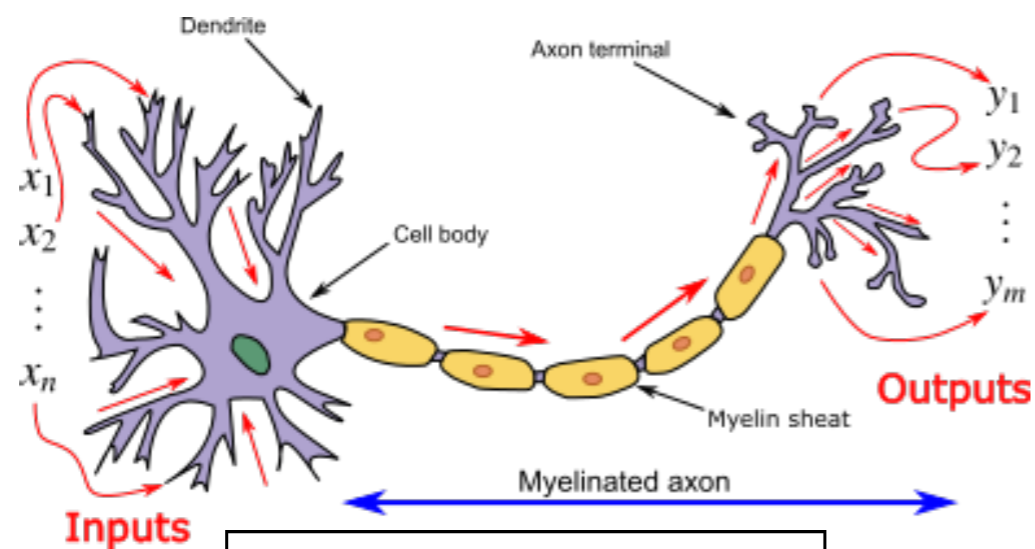
$$D_{LDA}(x) = \begin{cases} +1 & 2(\hat{\mu}_+ - \hat{\mu}_-)^T \hat{\Sigma}^{-1} x \\ & \geq \hat{\mu}_+^T \hat{\Sigma}^{-1} \hat{\mu}_+ - \hat{\mu}_-^T \hat{\Sigma}^{-1} \hat{\mu}_- \\ -1 & \text{otherwise} \end{cases}$$

→ This is a linear classifier.
 → We could also assume the covariance matrices are different.
 This is called Quadratic Discriminant Analysis.

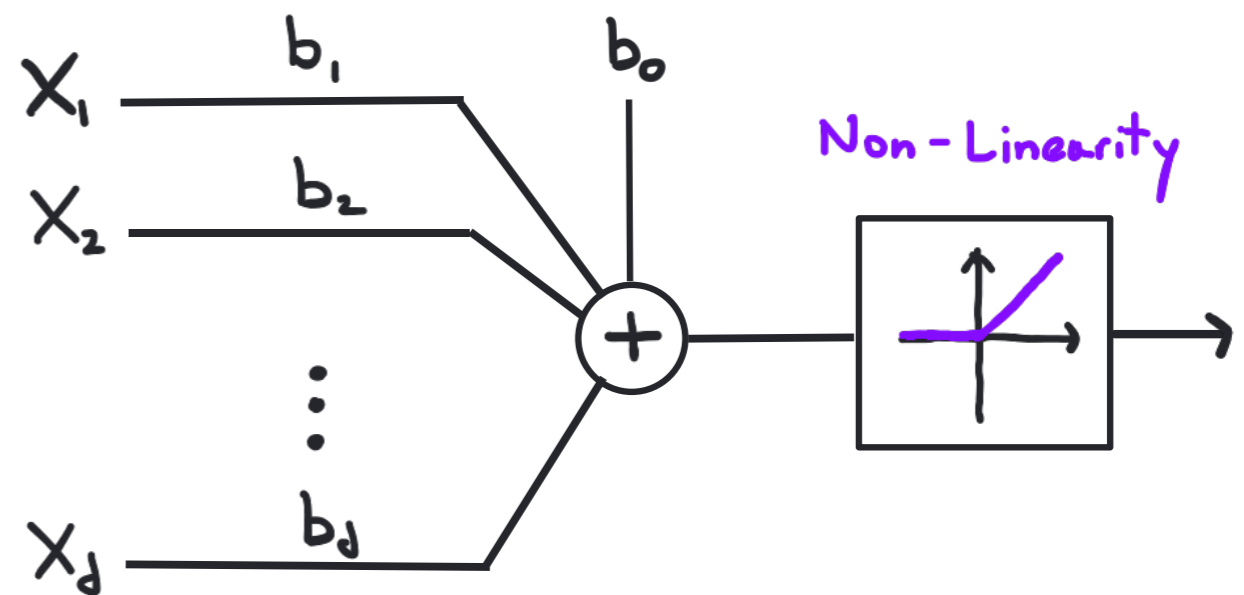


$n_{\text{train}} = 12$
 # Training Errors = 1
 Training Error Rate = $\frac{1}{12}$
 = 8.3%
 $n_{\text{test}} = 8$
 # Test Errors = 2
 Test Error Rate = $\frac{2}{8}$
 = 25%

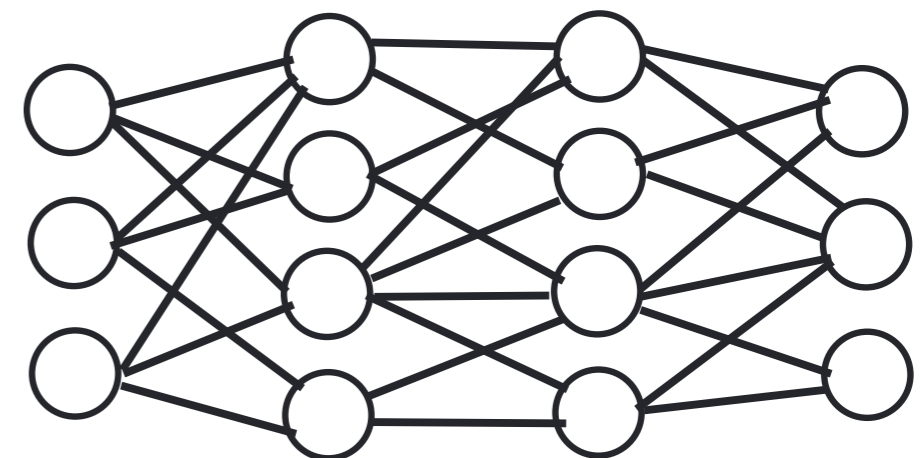
- Modern approaches to binary classification rely heavily on large-scale optimization techniques to carefully select the parameters of decision rules drawn from much richer families of functions.
- For instance, neural networks use gradient descent to set the weights of the individual units, which are loosely based on biological neurons.



commons.wikimedia.org/wiki/File:Neuron3.png
 Creative Commons License
 Author: Prof. Loc Vu-Quoc

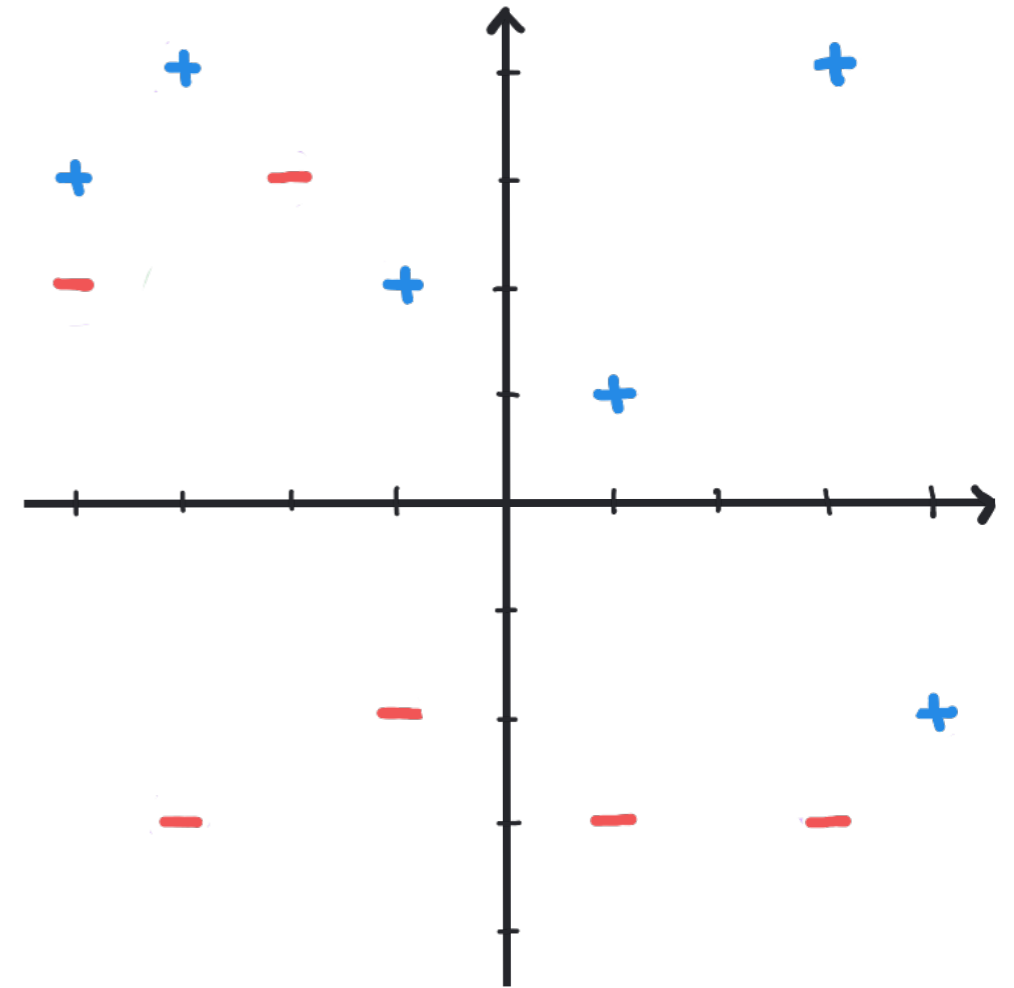
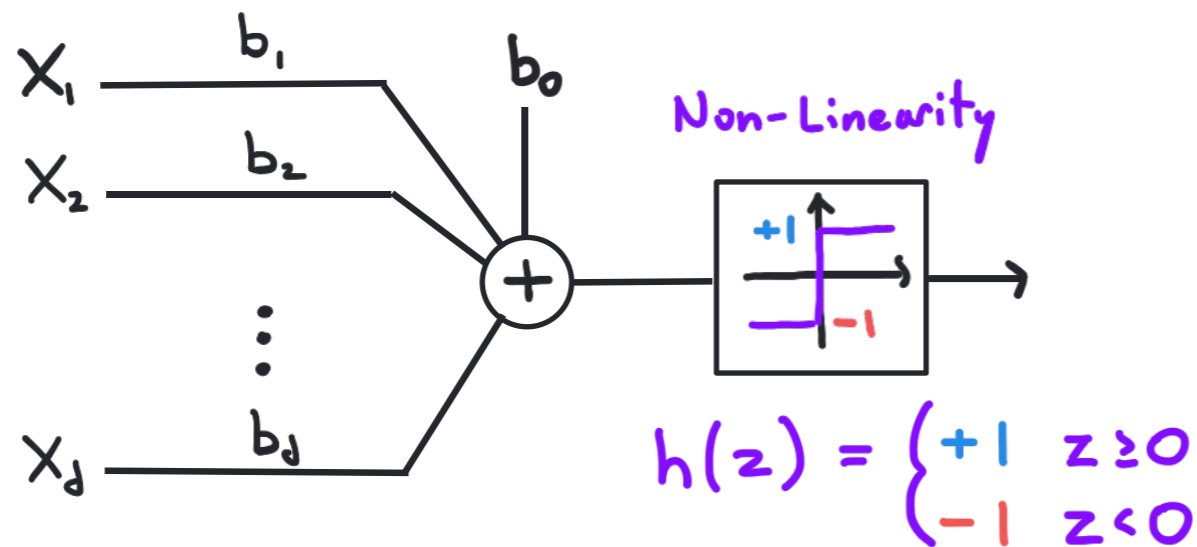


- Deep learning uses many layers of these units to tackle massive datasets. (Beyond our scope.)



• Perceptron Classifier:

→ Basic Idea: Although gradient descent is beyond our scope, we can implement a simple variant of the perceptron, the first artificial neural network.



→ Train Weights: $\underline{b} = \begin{bmatrix} b_0 \\ \vdots \\ b_d \end{bmatrix}$, initialize $\underline{b} = \underline{0}$

for $i = 1$ to n_{train}

$$Y_{\text{temp}} = h\left(\underline{b}^T \begin{bmatrix} 1 \\ \underline{x}_{\text{train},i} \end{bmatrix}\right)$$

Guess label. Note that $\underline{b}^T \begin{bmatrix} 1 \\ \underline{x}_{\text{train},i} \end{bmatrix} = b_0 + \sum_{j=1}^d b_j x_{\text{train},i,j}$

$$\underline{b} = \underline{b} + r (Y_{\text{train},i} - Y_{\text{temp}}) \begin{bmatrix} 1 \\ \underline{x}_{\text{train},i} \end{bmatrix}$$

Update weights if guess is wrong.

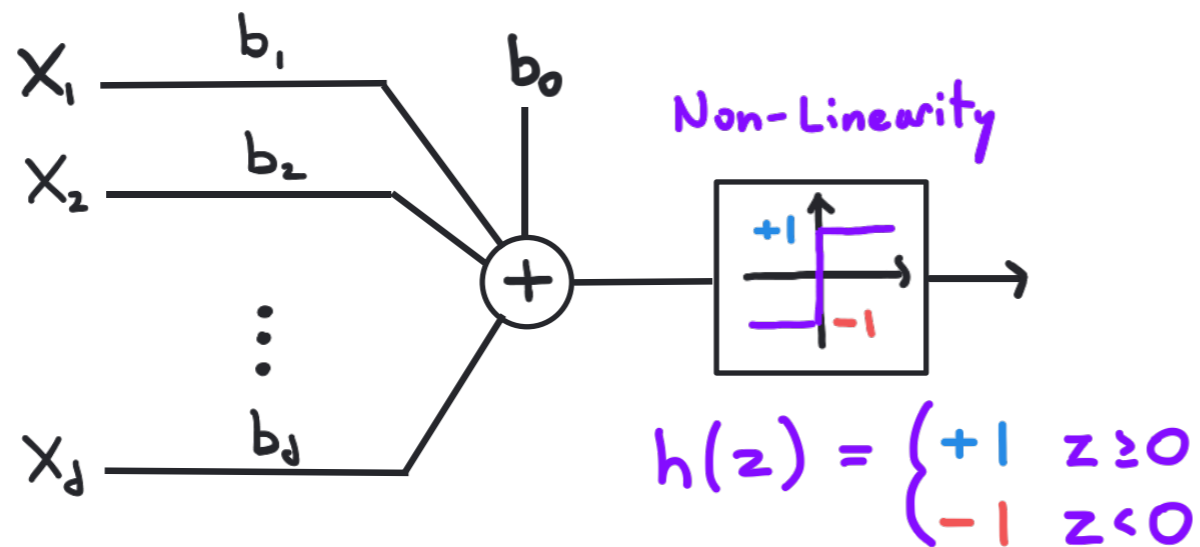
end

$0 < r < 1$ is the learning rate.

→ Use trained weights as a linear classifier.

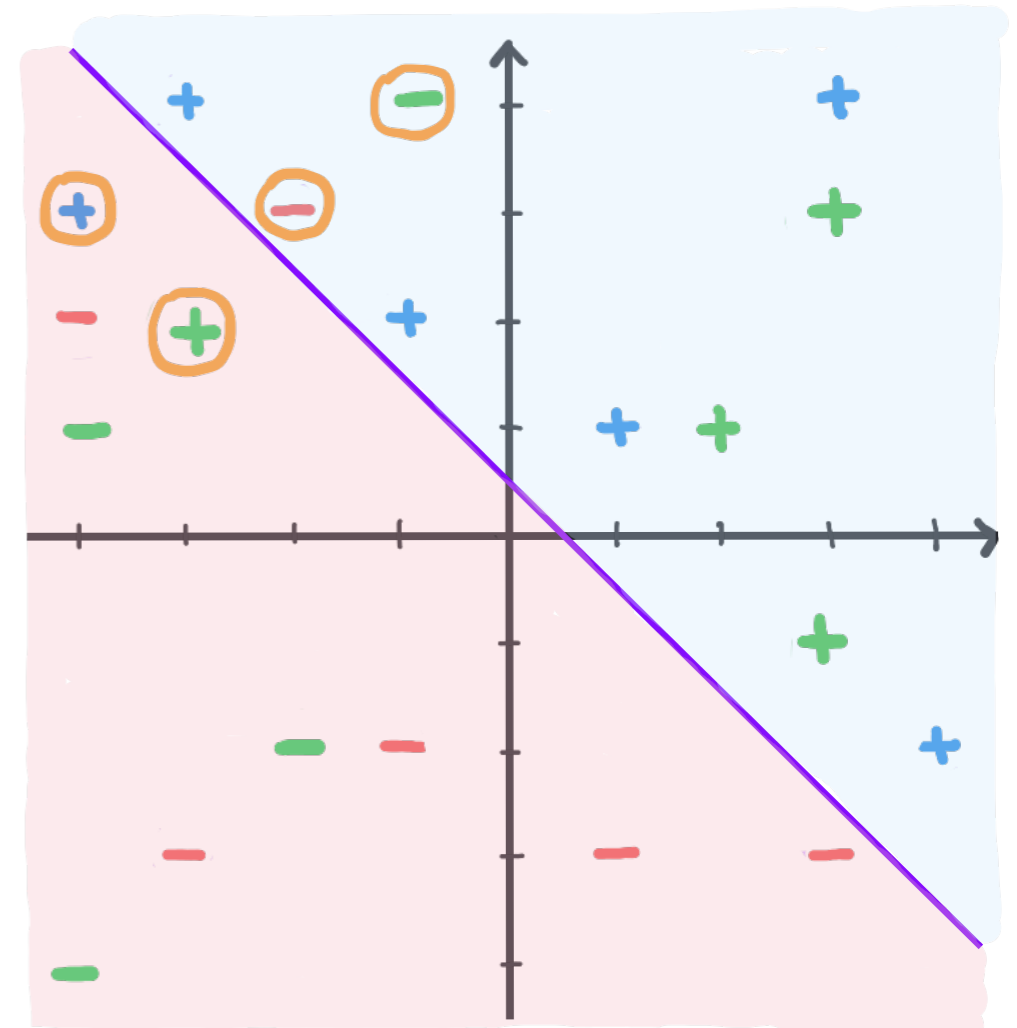
• Perceptron Classifier:

→ Basic Idea: Although gradient descent is beyond our scope, we can implement a simple variant of the perceptron, the first artificial neural network.



$$D_{\text{perceptron}}(\underline{x}) = \begin{cases} +1 & \underline{b}^T \begin{bmatrix} 1 \\ \underline{x} \end{bmatrix} \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

where $\underline{b} = \begin{bmatrix} b_0 \\ \vdots \\ b_d \end{bmatrix}$ are the trained weights.



$n_{\text{train}} = 12$

Training Errors = 2

Training Error Rate = $\frac{2}{12} = 16.7\%$

$n_{\text{test}} = 8$

Test Errors = 2

Test Error Rate = $\frac{2}{8} = 25\%$