

# Model Selection, Perceptron

CSC 461: Machine Learning

Fall 2022

Prof. Marco Alvarez  
University of Rhode Island

## Evaluation (model selection)

### Actuals/Predictions (example)

ID	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	1	1	1	0	0	0	0
Predicted	0	0	1	1	1	1	1	1	1	0	0	0
	FN	FN	TP	TP	TP	TP	TP	TP	FP	TN	TN	TN

### Confusion matrix (2 classes)

		Predicted condition	
		Positive (PP)	Negative (PN)
Actual condition	Total population = P + N		
	Positive (P)	True positive (TP)	False negative (FN)
	Negative (N)	False positive (FP)	True negative (TN)

[https://en.wikipedia.org/wiki/Confusion\\_matrix](https://en.wikipedia.org/wiki/Confusion_matrix)

## Confusion matrix (example)

		Predicted condition	
		Cancer	Non-cancer
Actual condition	Total 8 + 4 = 12	7	5
	Cancer 8	6	2
	Non-cancer 4	1	3

[https://en.wikipedia.org/wiki/Confusion\\_matrix](https://en.wikipedia.org/wiki/Confusion_matrix)

## Evaluation metrics (2 classes)

### accuracy (ACC)

$$ACC = \frac{TP + TN}{P + N} = \frac{TP + TN}{TP + TN + FP + FN}$$

### F1 score

is the harmonic mean of precision and sensitivity

$$F_1 = 2 \cdot \frac{PPV \cdot TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$$

### Matthews correlation coefficient (MCC)

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

[https://en.wikipedia.org/wiki/Confusion\\_matrix](https://en.wikipedia.org/wiki/Confusion_matrix)

## Evaluation metrics (2 classes)

### sensitivity, recall, hit rate, or true positive rate (TPR)

$$TPR = \frac{TP}{P} = \frac{TP}{TP + FN} = 1 - FNR$$

### specificity, selectivity or true negative rate (TNR)

$$TNR = \frac{TN}{N} = \frac{TN}{TN + FP} = 1 - FPR$$

### precision or positive predictive value (PPV)

$$PPV = \frac{TP}{TP + FP} = 1 - FDR$$

### negative predictive value (NPV)

$$NPV = \frac{TN}{TN + FN} = 1 - FOR$$

### miss rate or false negative rate (FNR)

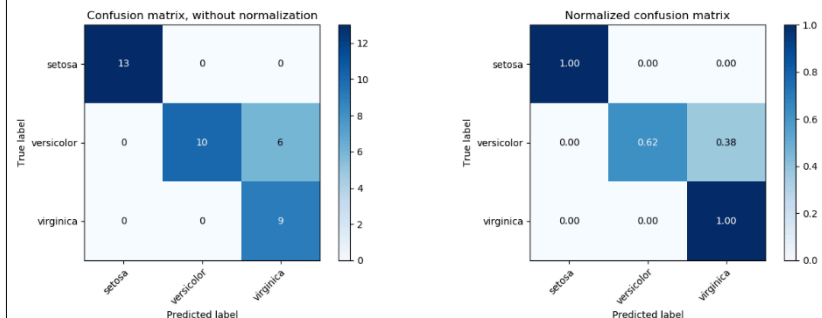
$$FNR = \frac{FN}{P} = \frac{FN}{FN + TP} = 1 - TPR$$

### fall-out or false positive rate (FPR)

$$FPR = \frac{FP}{N} = \frac{FP}{FP + TN} = 1 - TNR$$

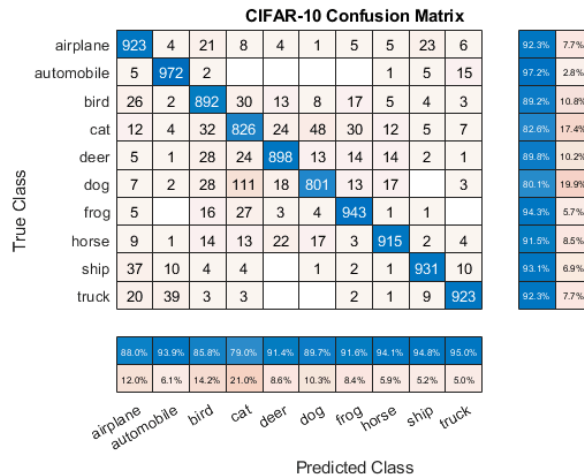
[https://en.wikipedia.org/wiki/Confusion\\_matrix](https://en.wikipedia.org/wiki/Confusion_matrix)

## Confusion matrix (>2 classes)



[https://scikit-learn.org/stable/auto\\_examples/model\\_selection/plot\\_confusion\\_matrix.html](https://scikit-learn.org/stable/auto_examples/model_selection/plot_confusion_matrix.html)

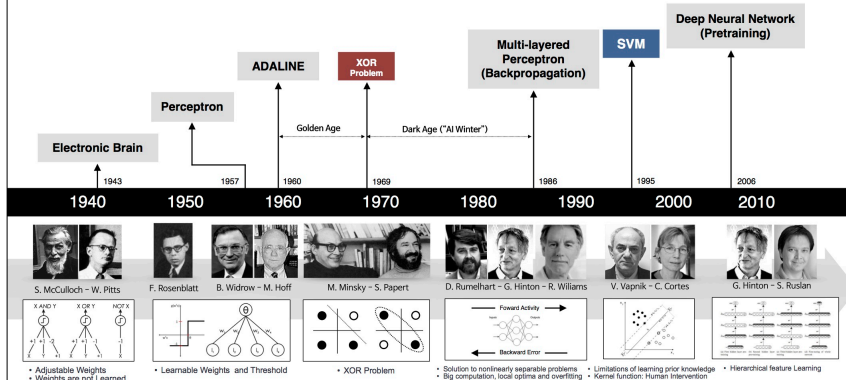
## Confusion matrix (>2 classes)



<https://www.mathworks.com/help/deeplearning/ref/confusionchart.html>

## The Perceptron

## Neural Networks



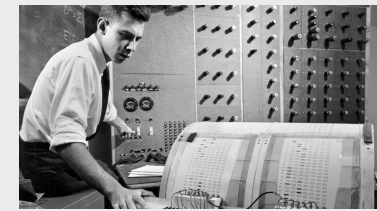
[http://beamandrew.github.io/deeplearning/2017/02/23/deep\\_learning\\_101\\_part1.html](http://beamandrew.github.io/deeplearning/2017/02/23/deep_learning_101_part1.html)

## Rosenblatt (1958)

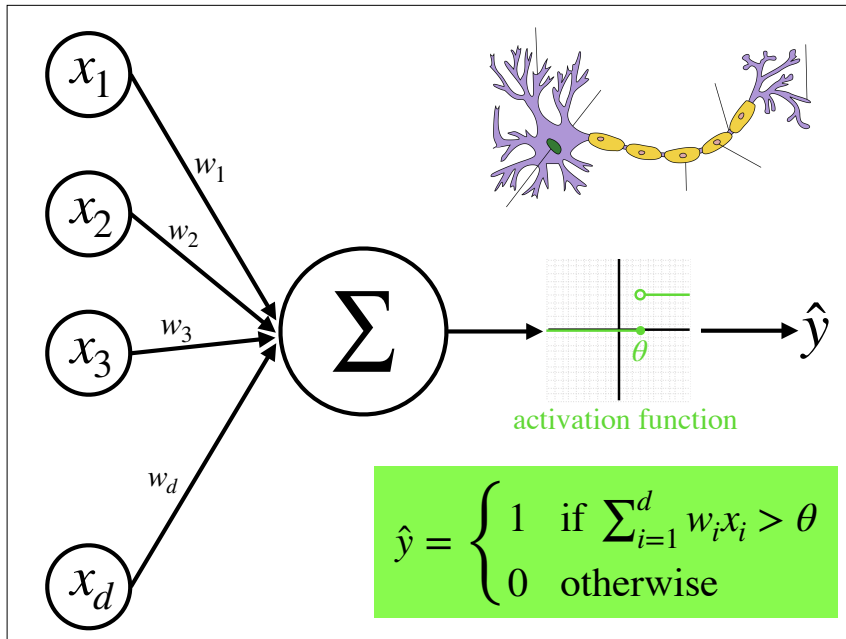
- ▶ Perceptron introduced by Frank Rosenblatt (psychologist, logician)
- ✓ based on work from McCulloch-Pitts and Hebb
- ✓ very powerful **learning** algorithm with high expectations

**NEW NAVY DEVICE LEARNS BY DOING; Psychologist Shows Embryo of Computer Designed to Read and Grow Wiser**

WASHINGTON, July 7, 1958 (UPI) -- The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence.



<https://news.cornell.edu/stories/2019/09/professors-perceptron-paved-way-ai-60-years-too-soon>



## Absorbing the threshold/bias

$$\hat{y} = \begin{cases} 1 & \text{if } \sum_{i=1}^d w_i x_i > \theta \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{y} = \begin{cases} 1 & \text{if } \sum_{i=1}^d w_i x_i - \theta > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{y} = \begin{cases} 1 & \text{if } \sum_{i=0}^d w_i x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$x_0 = +1, \quad w_0 = -\theta$$

## Another look

For convenience we will use +1 and -1 instead of 1 and 0

$$h_{\mathbf{w}}(\mathbf{x}) = \sigma \left( \sum_{i=0}^d w_i x_i \right) = \sigma(\mathbf{w}^T \mathbf{x})$$

$$\sigma(z) = \begin{cases} 1 & \text{if } z > 0 \\ -1 & \text{if } z \leq 0 \end{cases}$$

$$\mathcal{H} = \{h_{\mathbf{w}} : \mathbf{w} \in \mathbb{R}^{d+1}\}$$

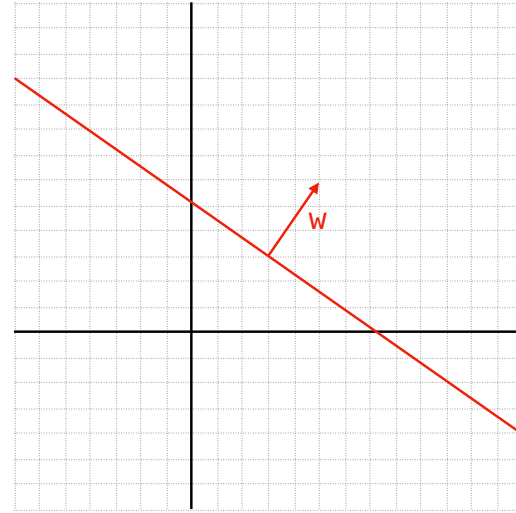
## Perceptron Algorithm

- ▶ Start with a null vector  $\mathbf{w}$
- ▶ Repeat for T epochs
  - ✓ shuffle the data instances
  - ✓ for all examples in training data
    - ✓ if misclassified
      - update the weight vector by adding  $\mathbf{x}$  to  $\mathbf{w}$  if the actual label is positive and subtracting  $\mathbf{x}$  from  $\mathbf{w}$  otherwise
- ▶ Return  $\mathbf{w}$

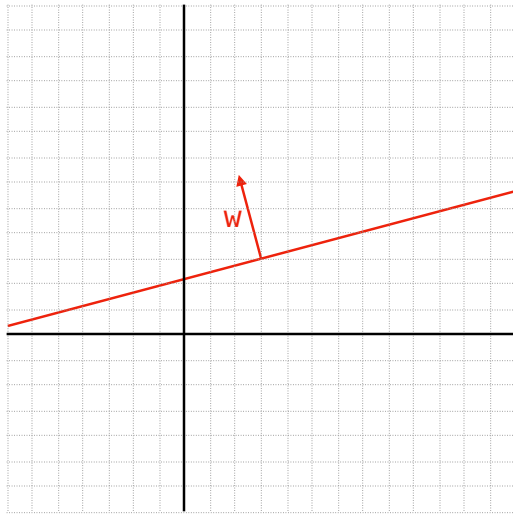
# Write the pseudocode

<https://colab.research.google.com/drive/1p7LjENLd6whkoVY2yRGg-Nm5v4XNjdDt>

## Mistake on a positive (update)



## Mistake on a negative (update)



## Intuition

- Suppose a mistake on the positive side:

$$y = +1 \quad \mathbf{w}^T \mathbf{x} \leq 0$$

- After 1 update the new weight vector will be:

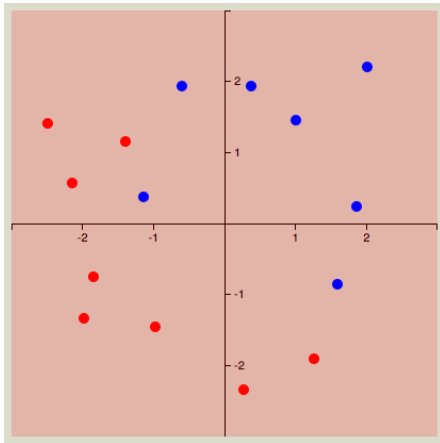
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{x}$$

- Classifying the datapoint with the new weight vector:

$$\mathbf{w}_{t+1}^T \mathbf{x} = (\mathbf{w}_t + \mathbf{x})^T \mathbf{x} = \mathbf{w}_t^T \mathbf{x} + \mathbf{x}^T \mathbf{x} \geq \mathbf{w}_t^T \mathbf{x}$$

use same idea for mistakes on the negative side

## Demo



[https://planspace.org/20150907-interactive\\_perceptron\\_training\\_toy/](https://planspace.org/20150907-interactive_perceptron_training_toy/)

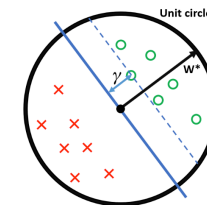
## Perceptron convergence theorem

The argument goes as follows: Suppose  $\exists \mathbf{w}^*$  such that  $y_i(\mathbf{x}_i^\top \mathbf{w}^*) > 0 \forall (\mathbf{x}_i, y_i) \in D$ .

Now, suppose that we rescale each data point and the  $\mathbf{w}^*$  such that

$$\|\mathbf{w}^*\| = 1 \quad \text{and} \quad \|\mathbf{x}_i\| \leq 1 \quad \forall \mathbf{x}_i \in D$$

Let us define the Margin  $\gamma$  of the hyperplane  $\mathbf{w}^*$  as  $\gamma = \min_{(\mathbf{x}_i, y_i) \in D} |\mathbf{x}_i^\top \mathbf{w}^*|$ .



To summarize our setup:

- All inputs  $\mathbf{x}_i$  live within the unit sphere
- There exists a separating hyperplane defined by  $\mathbf{w}^*$ , with  $\|\mathbf{w}^*\| = 1$  (i.e.  $\mathbf{w}^*$  lies exactly on the unit sphere).
- $\gamma$  is the distance from this hyperplane (blue) to the closest data point.

**Theorem:** If all of the above holds, then the Perceptron algorithm makes at most  $1/\gamma^2$  mistakes.

<http://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote03.html>

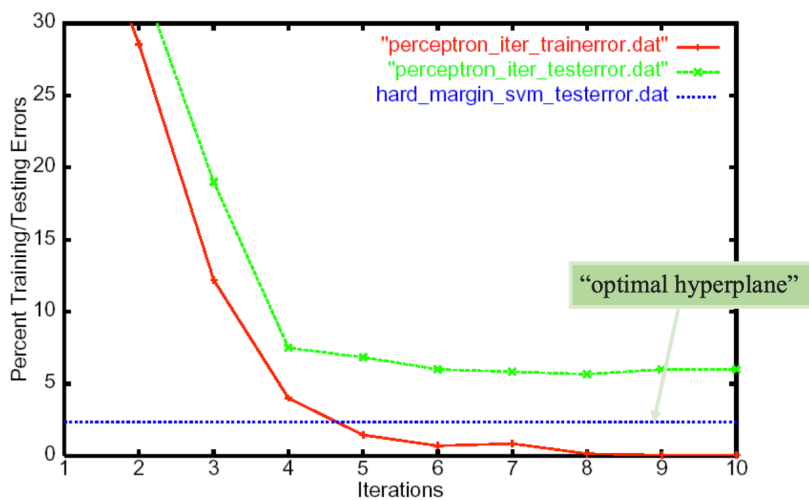
## Perceptron (remarks)

- ▶ Assumes data is linearly separable
  - ✓ **does not converge** if classes are not linearly separable
- ▶ Different correct solutions can be found
  - ✓ most are not optimal in terms of generalization
- ▶ Averaged Perceptron
  - ✓ returns a weighted average of earlier hypotheses

## Parameters vs Hyperparameters

- ▶ **Parameters**
  - ✓ weights and bias
- ▶ **Hyperparameters**
  - ✓ number of epochs (one epoch is one pass over the training data)

## Example: Reuters Text Classification

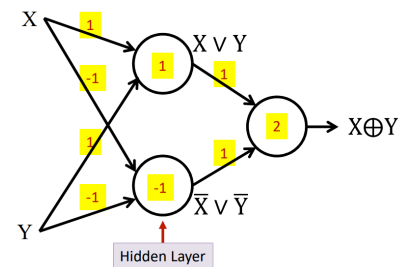
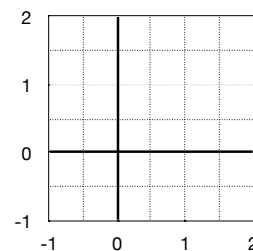


Credit: <http://www.cs.cornell.edu/courses/cs4780/2019fa/lectures/06-perceptron.pdf>

## Minsky & Papert (1969)

### Perceptrons

- ✓ influential book
- ✓ analyzed the algorithm and showed limitations (e.g. XOR)



<http://deeplearning.cs.cmu.edu/document/lecture/lecture-1.pdf>

## Example

X0	X1	X2	Y
1	0	0	-1
1	1	0	+1
1	1	1	+1
1	0	1	+1

$$\sigma(z) = \begin{cases} +1 & \text{if } z > 0 \\ -1 & \text{if } z \leq 0 \end{cases}$$

$$h_w(\mathbf{x}) = \sigma\left(\sum_{i=0}^d w_i x_i\right) = \sigma(\mathbf{w}^T \mathbf{x})$$

$$L_{0/1}(h, \mathcal{D}) = \frac{1}{n} \sum_{(x_i, y_i) \in \mathcal{D}} I(h(x_i) \neq y_i)$$

$$\mathbf{w}_a = [0, 0, 0]^T$$

$$\mathbf{w}_b = [0, 1, 0]^T$$

$$\mathbf{w}_c = [-1, 2, 2]^T$$

$$L_{0/1}(h_{\mathbf{w}_a}, \mathcal{D}) = ?$$

$$L_{0/1}(h_{\mathbf{w}_b}, \mathcal{D}) = ?$$

$$L_{0/1}(h_{\mathbf{w}_c}, \mathcal{D}) = ?$$

## Back to the perceptron ...

$$L(\mathbf{w}) = \sum_{i=1}^n -y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}$$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} L(\mathbf{w})$$

Note this loss function has a problem, it is not bounded below

## Gradient

With respect to a single  $w_j$

$$\begin{aligned} \frac{\partial}{\partial w_j} L(\mathbf{w}) &= \frac{\partial}{\partial w_j} \left[ - \sum_{i=1}^m y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)} \right] \\ &= - \sum_{i=1}^m y^{(i)} x_j^{(i)} \end{aligned}$$

The training algorithm can focus on batches of misclassified instances, then the batch loss cannot be negative

## Pseudocode