# Perceptrons & Neural Networks

Russell and Norvig:
Chapter 18 (18.6.3,18.7)

# Outline

- ◆ Perceptrons
- ◆ Neural Networks
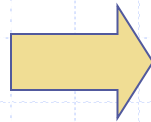
# Classification: Feature Vectors
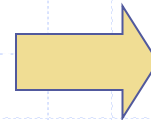
$$x \qquad\qquad f(x) \qquad\qquad y$$

```
Hello,

Do you want free printr
cartriges?  Why pay more
when you can get them
ABSOLUTELY FREE!  Just
```
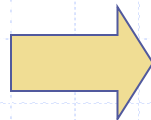
➡️

```
# free        : 2
YOUR_NAME     : 0
MISSPELLED    : 2
FROM_FRIEND   : 0
...
```
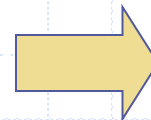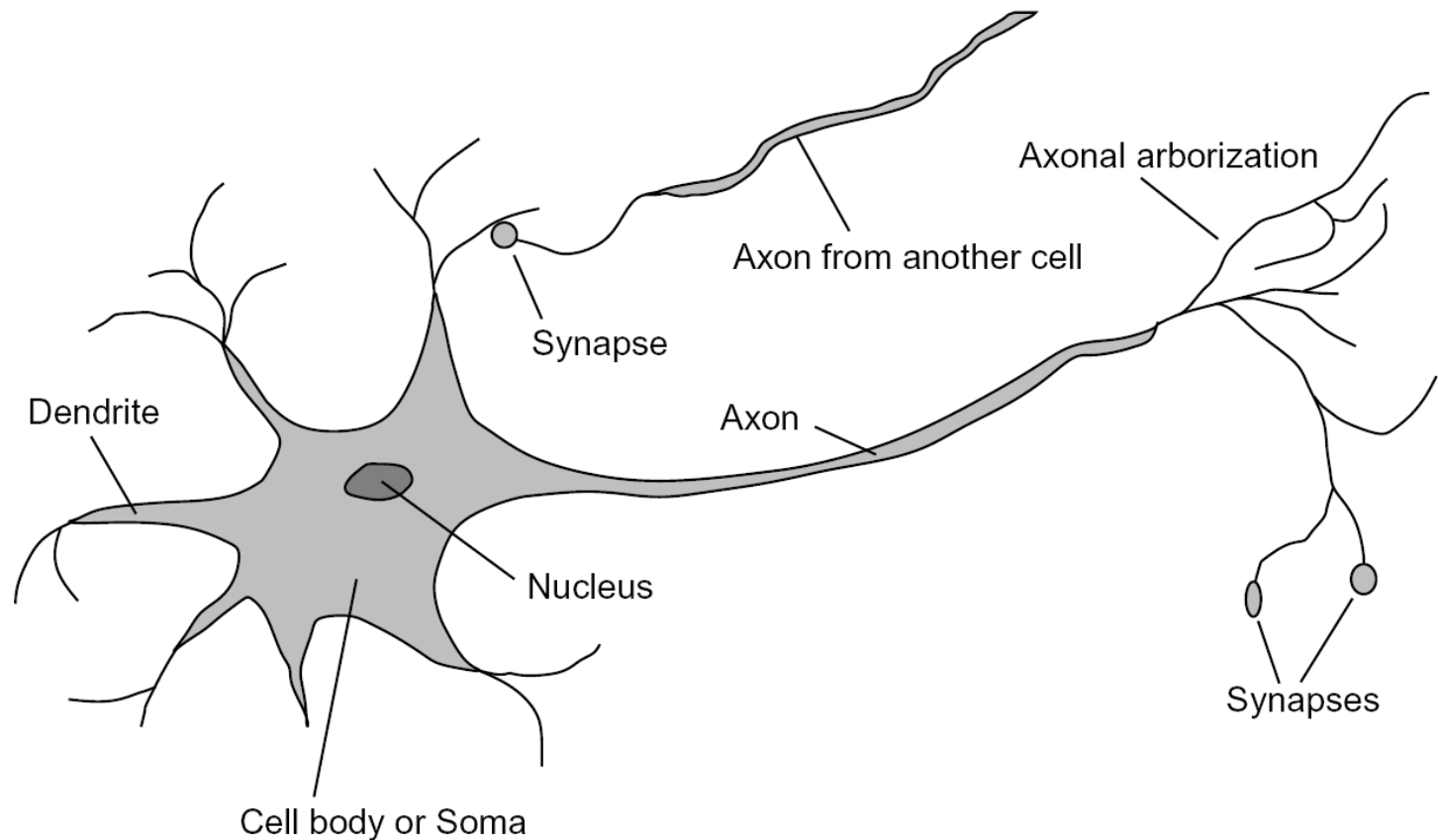
➡️

SPAM
or
+



➡️

```
PIXEL-7,12  : 1
PIXEL-7,13  : 0
...
NUM_LOOPS   : 1
...
```
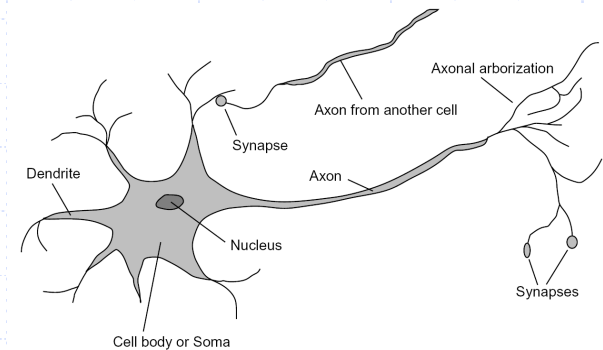
➡️

"2"

# Some (Simplified) Biology

◆ Very loose inspiration: human neurons



Dendrite

Synapse

Axon from another cell
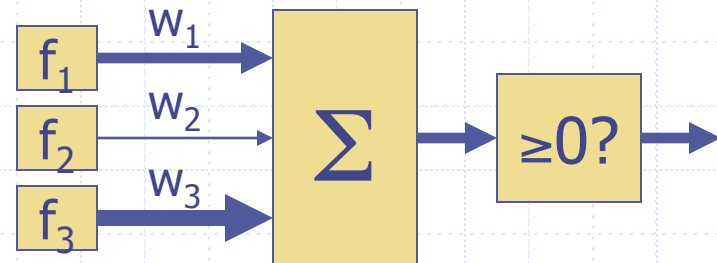
Axonal arborization

Axon

Nucleus

Cell body or Soma

Synapses

# Linear Classifiers

◆ Inputs are feature values
◆ Each feature has a weight
◆ Sum is the activation

$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

◆ If the activation is:
  - Positive, output +1
  - Negative, output -1

# Classification: Weights

- Binary case: compare features to a weight vector
- Learning: figure out the weight vector from examples

$$
\begin{bmatrix}
\texttt{\# free} & \texttt{: 4} \\
\texttt{YOUR\_NAME} & \texttt{:-1} \\
\texttt{MISSPELLED} & \texttt{: 1} \\
\texttt{FROM\_FRIEND} & \texttt{:-3} \\
\texttt{...} &
\end{bmatrix}
$$

$w$

$$
\begin{bmatrix}
\texttt{\# free} & \texttt{: 2} \\
\texttt{YOUR\_NAME} & \texttt{: 0} \\
\texttt{MISSPELLED} & \texttt{: 2} \\
\texttt{FROM\_FRIEND} & \texttt{: 0} \\
\texttt{...} &
\end{bmatrix}
$$

$f(x_1)$

$f(x_2)$

$$
\begin{bmatrix}
\texttt{\# free} & \texttt{: 0} \\
\texttt{YOUR\_NAME} & \texttt{: 1} \\
\texttt{MISSPELLED} & \texttt{: 1} \\
\texttt{FROM\_FRIEND} & \texttt{: 1} \\
\texttt{...} &
\end{bmatrix}
$$

*Dot product $w \cdot f$ positive means the positive class*

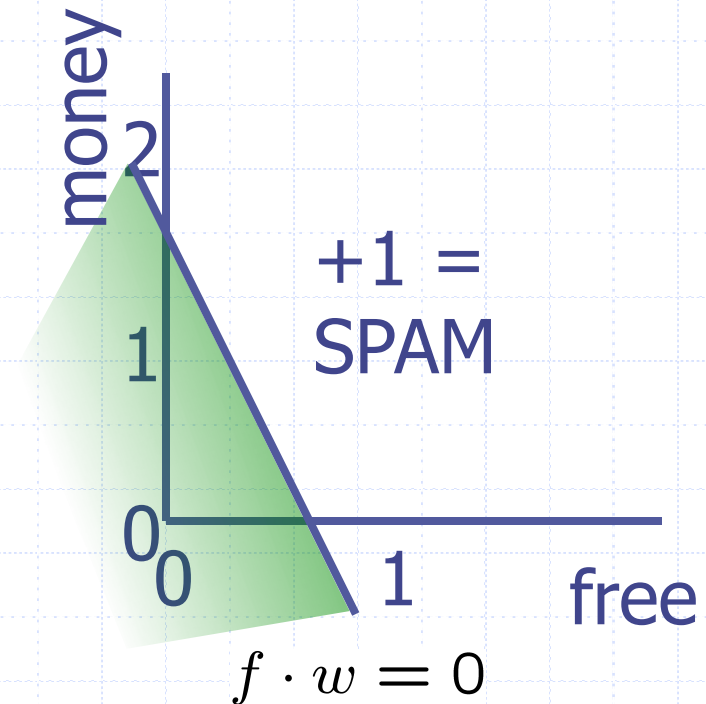# Binary Decision Rule

In the space of feature vectors

- Examples are points
- Any weight vector is a hyperplane
- One side corresponds to Y=+1
- Other corresponds to Y=-1

$w$

```
BIAS  :  -3
free  :   4
money :   2
...
```

money

2

1

0

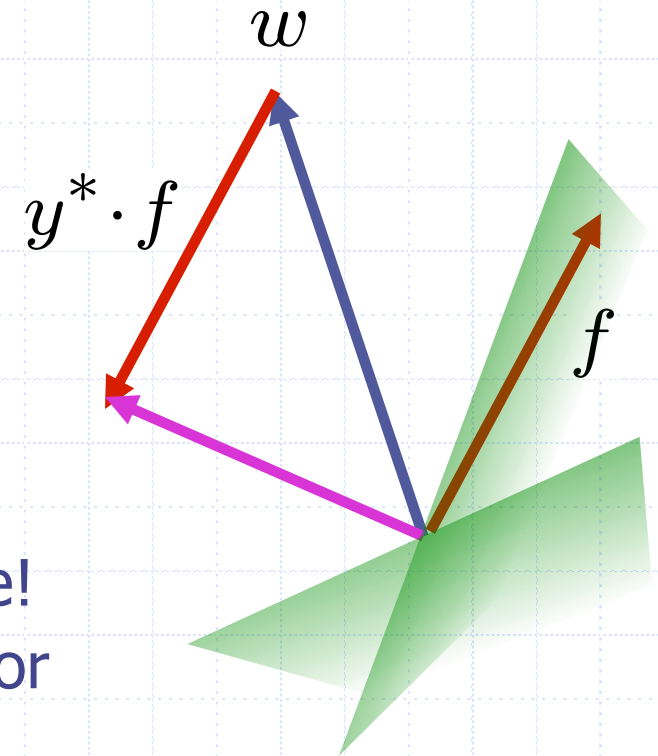0          1          free

+1 =
SPAM

-1 =
HAM

$f \cdot w = 0$

# Learning: Binary Perceptron

◆ Start with weights = 0

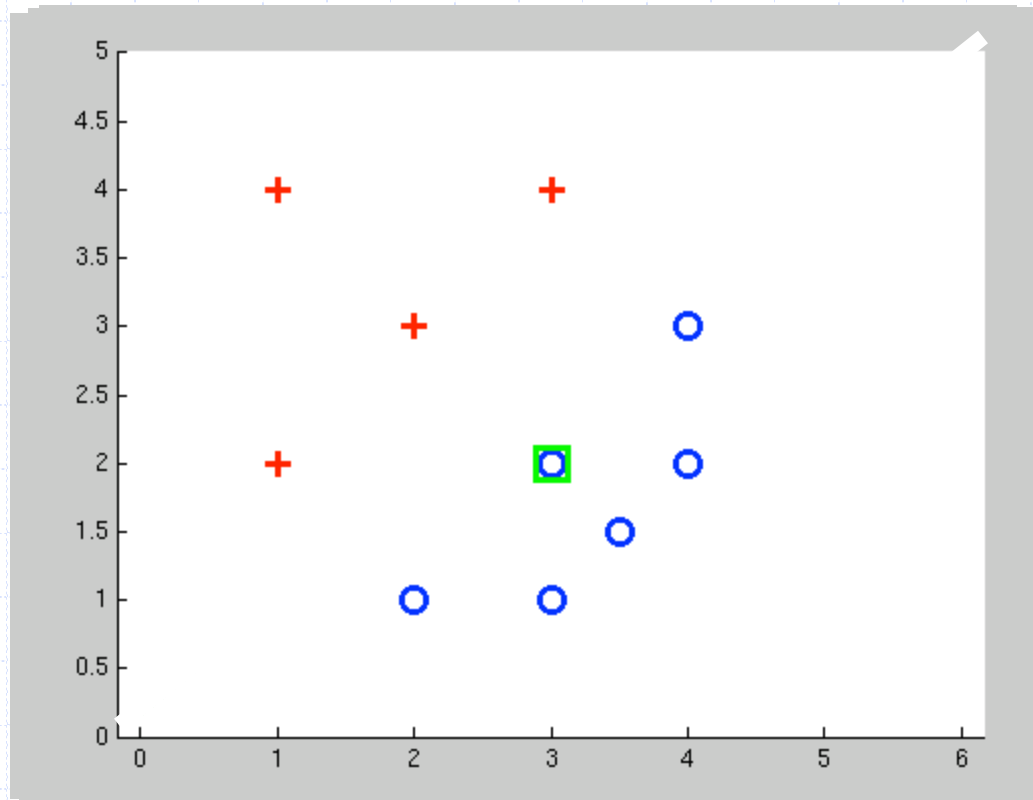◆ For each training instance:

  ▪ Classify with current weights

$$y = \begin{cases} +1 & \text{if } w \cdot f(x) \geq 0 \\ -1 & \text{if } w \cdot f(x) < 0 \end{cases}$$

  ▪ If correct (i.e., y=y*), no change!
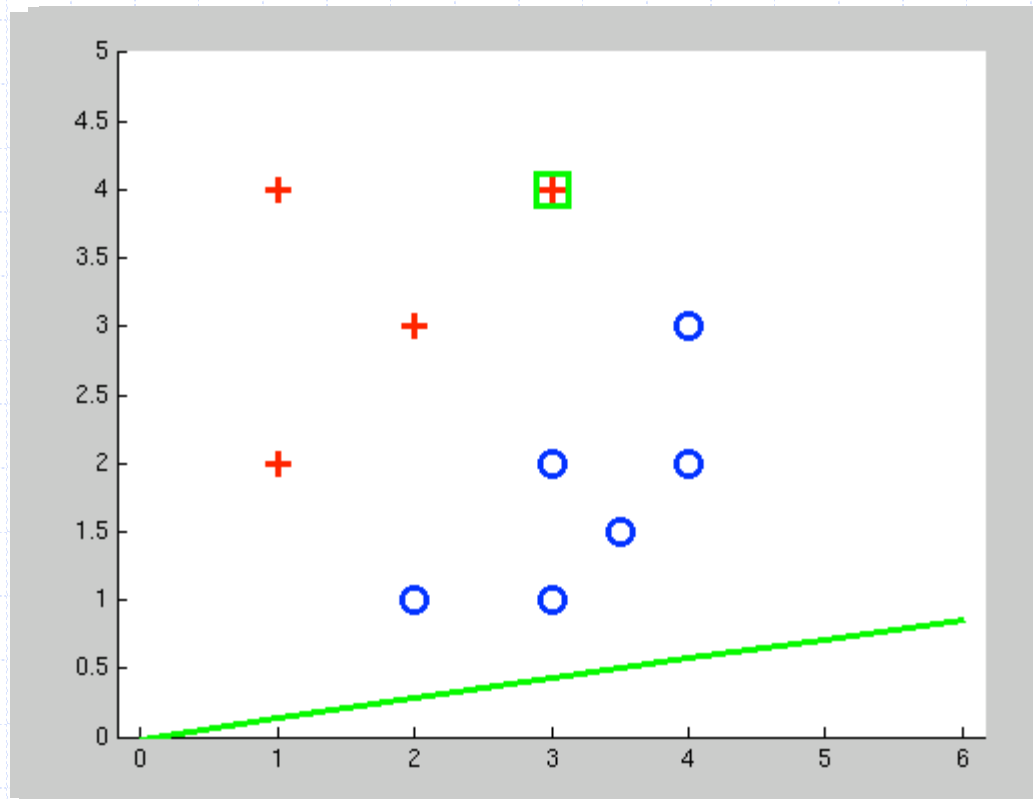  ▪ If wrong: adjust the weight vector by adding or subtracting the feature vector. Subtract if y* is -1.
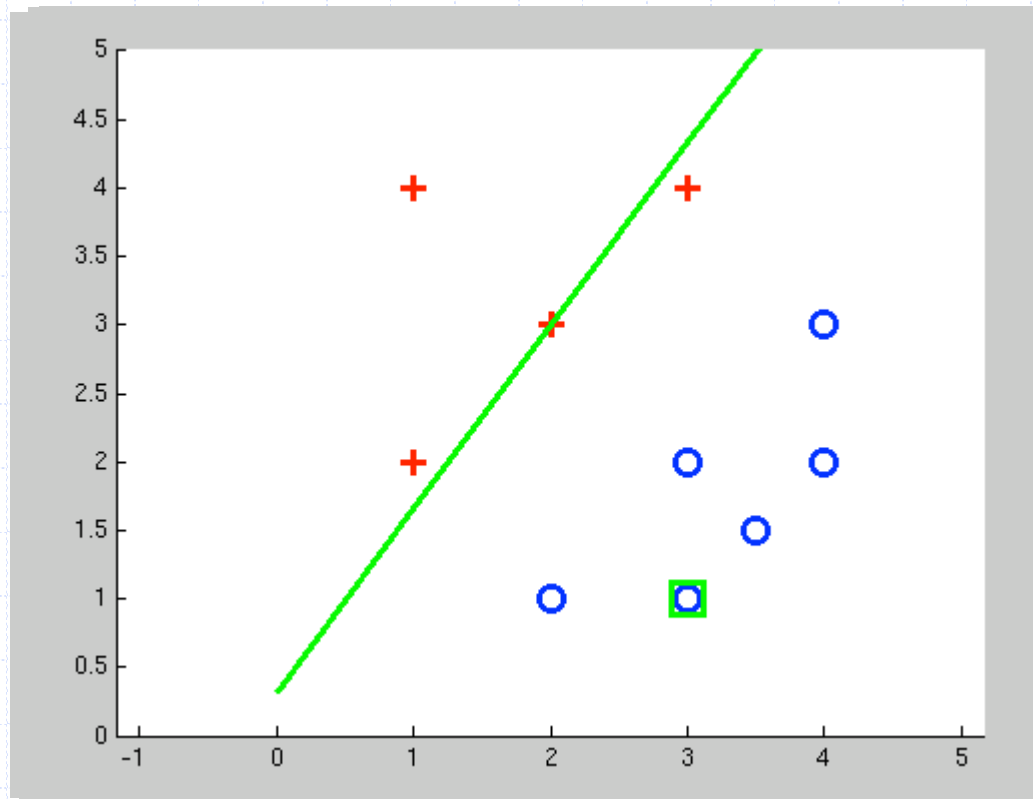
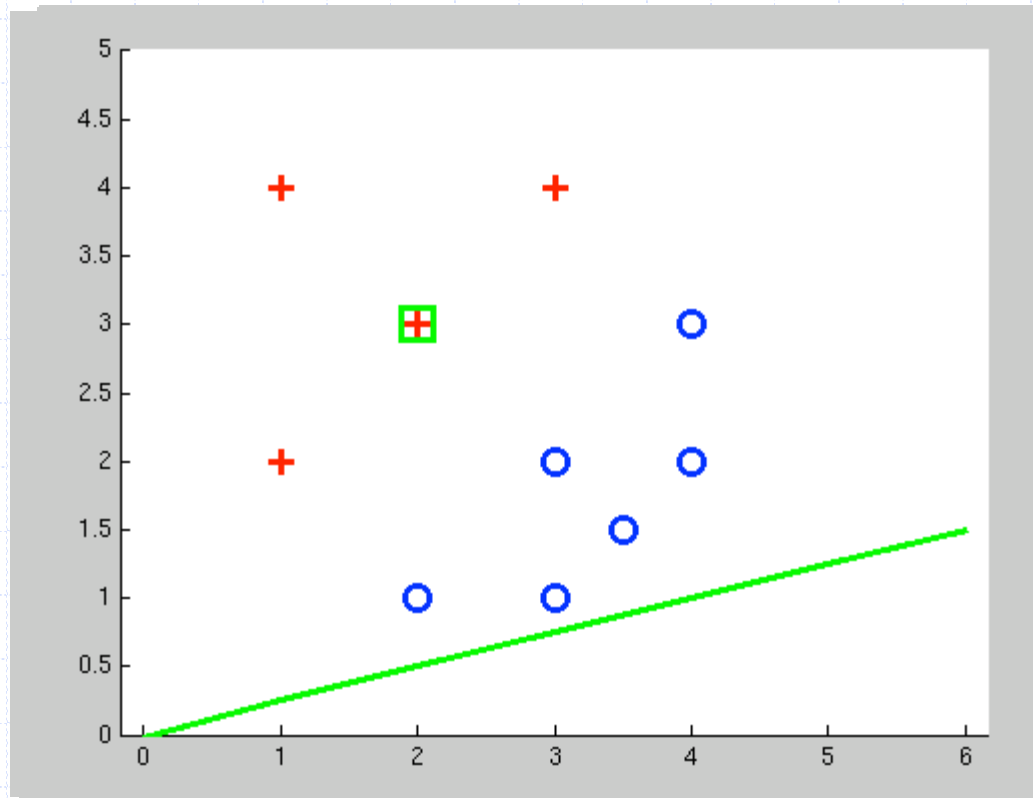$$w = w + y^* \cdot f$$

$w$

$y^* \cdot f$

$f$

# Examples: Perceptron
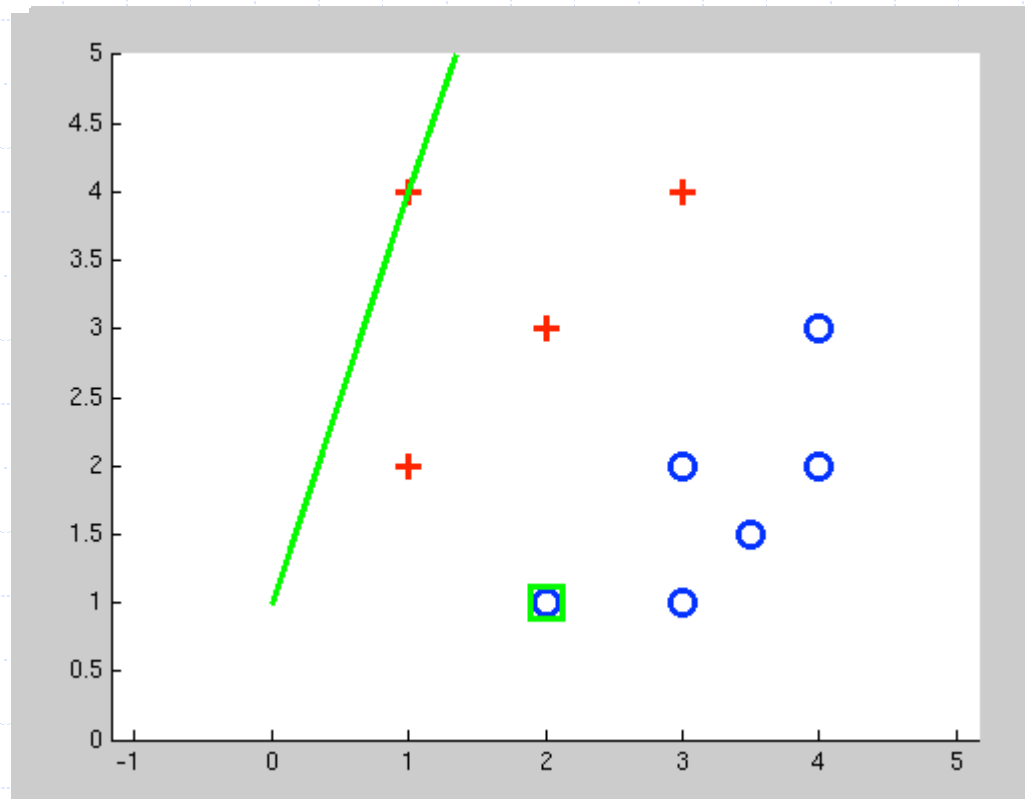
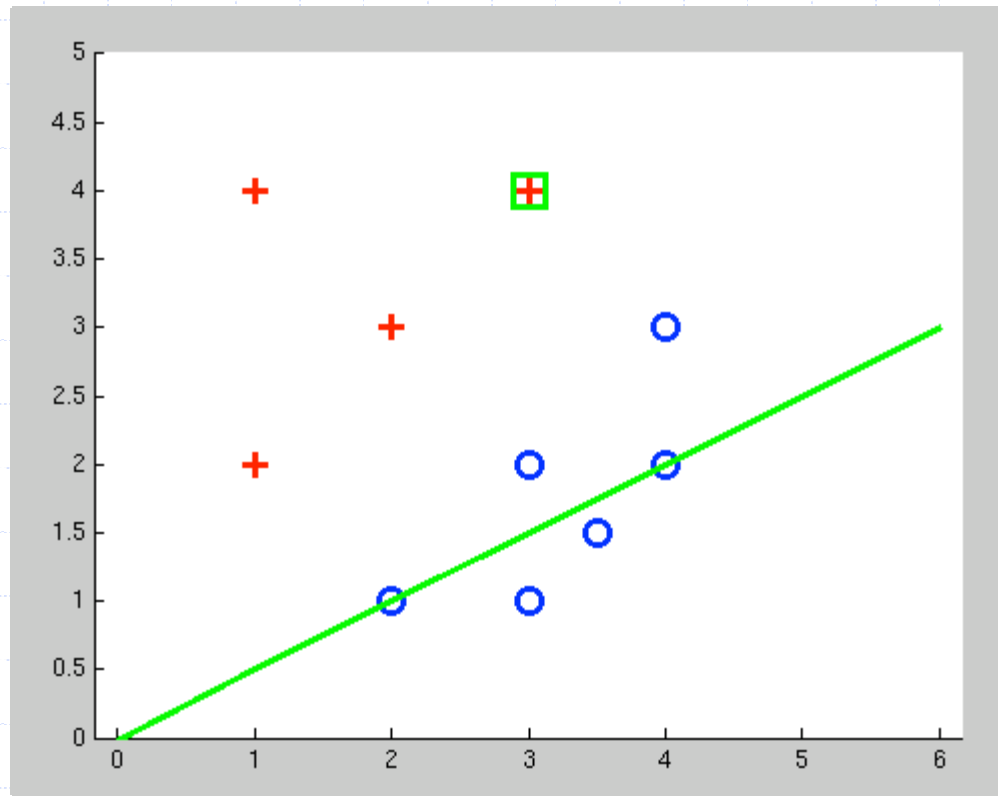# Examples: Perceptron

# Examples: Perceptron
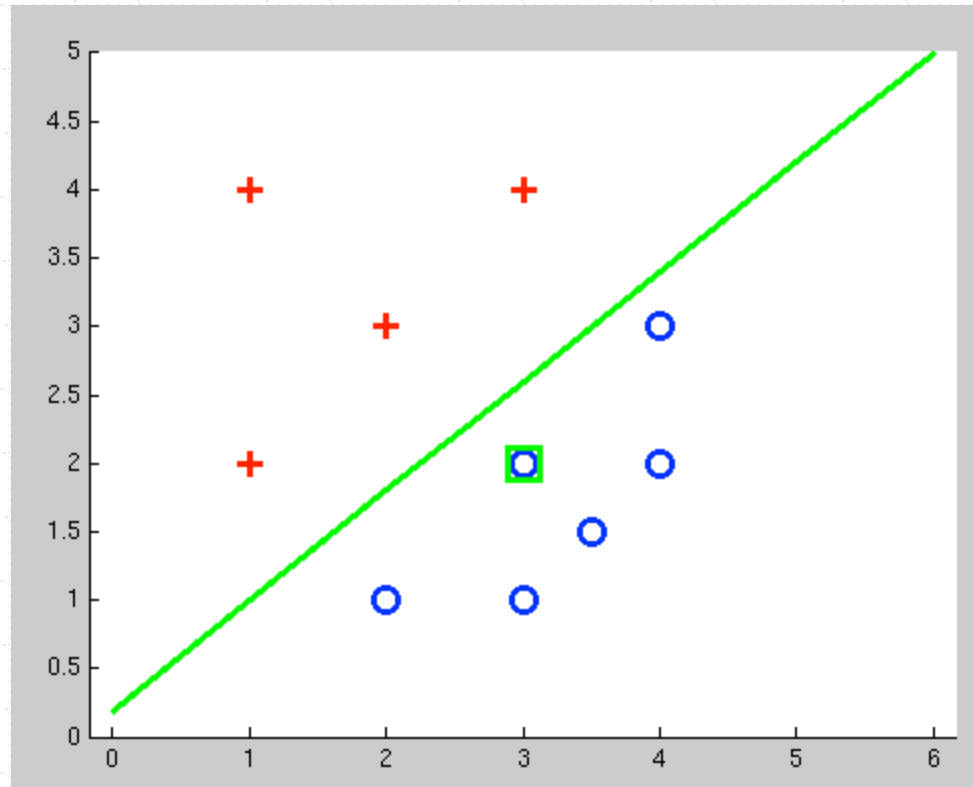
# Examples: Perceptron

# Examples: Perceptron

# Examples: Perceptron

# Examples: Perceptron

# Multiclass Decision Rule

- ◆ If we have multiple classes:
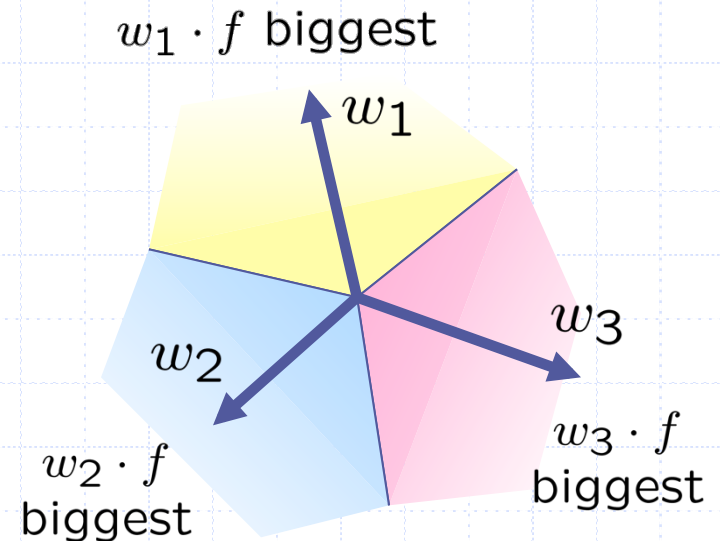  - A weight vector for each class:

    $$w_y$$

  - Score (activation) of a class y:

    $$w_y \cdot f(x)$$

  - Prediction highest score wins

    $$y = \arg\max_y \ w_y \cdot f(x)$$

$w_1 \cdot f$ biggest

$w_1$

$w_3$

$w_2$

$w_3 \cdot f$ biggest

$w_2 \cdot f$ biggest

*Binary = multiclass where the negative class has weight zero*
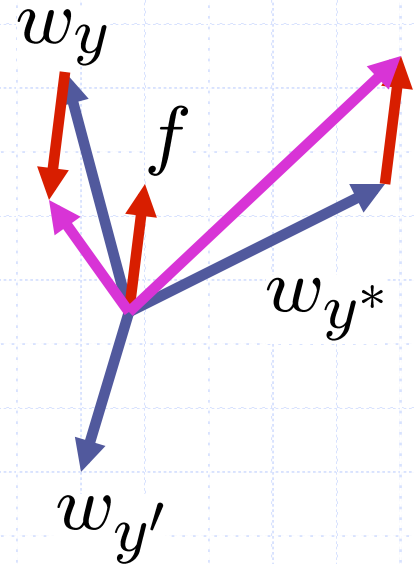
# Learning: Multiclass Perceptron

◆ Start with all weights = 0

◆ Pick up training examples one by one

◆ Predict with current weights

$$y = \arg\max_y \ w_y \cdot f(x)$$

◆ If correct, no change!

◆ If wrong: lower score of wrong answer, raise score of right answer

$$w_y = w_y - f(x)$$

$$w_{y^*} = w_{y^*} + f(x)$$



$w_y$

$f$

$w_{y^*}$

$w_{y'}$

# Example: Multiclass Perceptron

"win the vote"

"win the election"

"win the game"

$$w_{SPORTS}$$

$$w_{POLITICS}$$

$$w_{TECH}$$

```
BIAS  : 1
win   : 0
game  : 0
vote  : 0
the   : 0
...
```

```
BIAS  : 0
win   : 0
game  : 0
vote  : 0
the   : 0
...
```
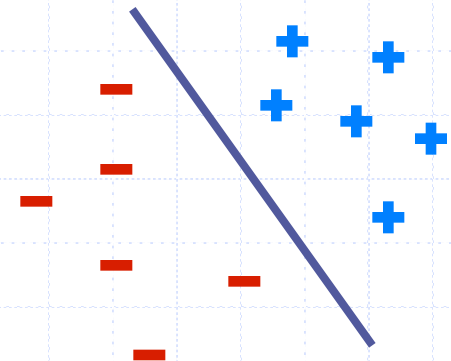
```
BIAS  : 0
win   : 0
game  : 0
vote  : 0
the   : 0
...
```
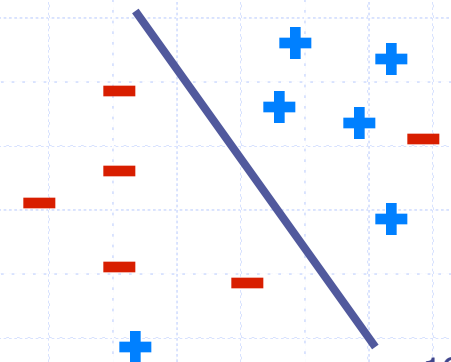
# Properties of Perceptrons

◆ Separability: some parameters get the training set perfectly correct

◆ Convergence: if the training is separable, perceptron will eventually converge (binary case)

◆ Mistake Bound: the maximum number of mistakes (binary case) related to the *margin* or degree of separability

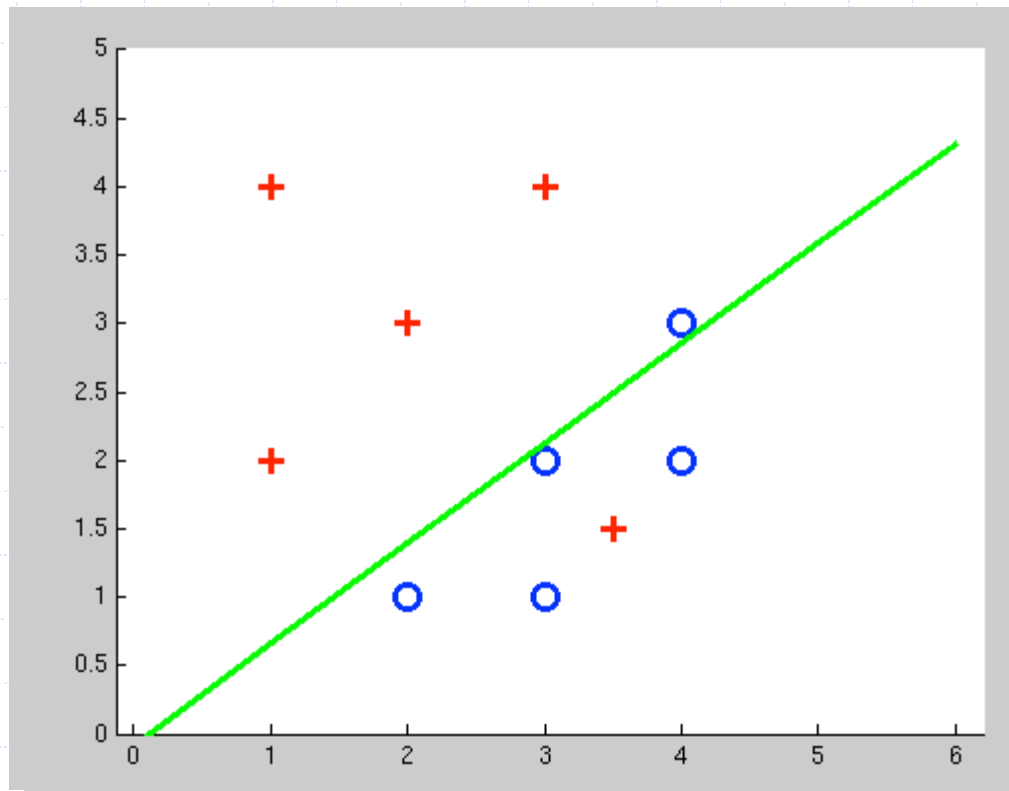$$\text{mistakes} < \frac{k}{\delta^2}$$
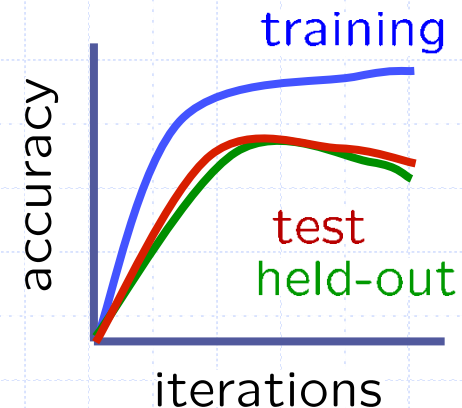
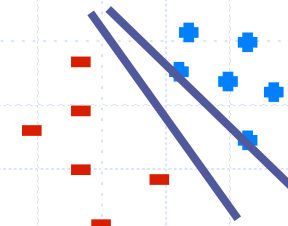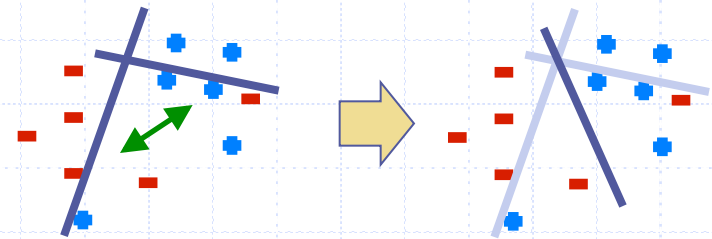Separable

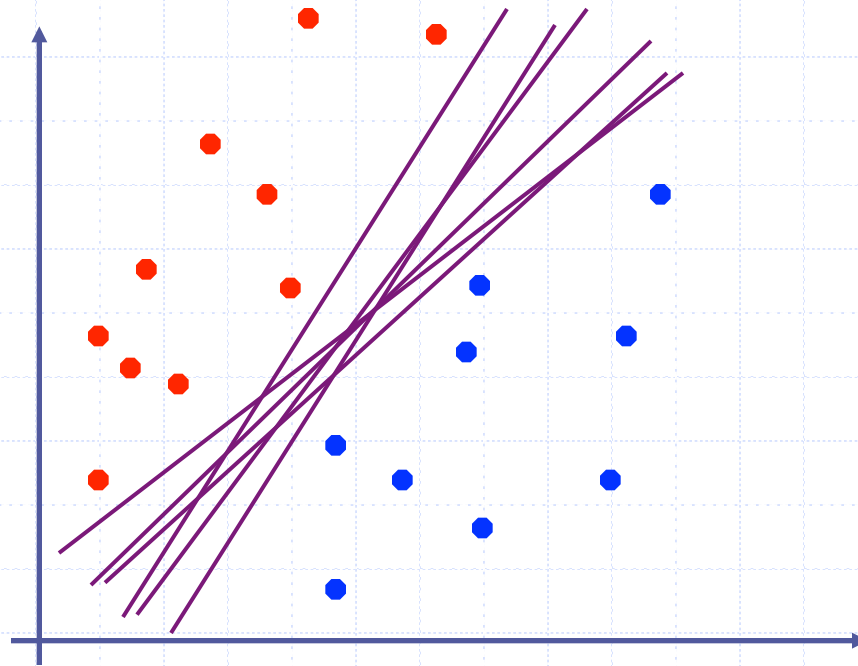Non-Separable

# Examples: Perceptron

◆ Non-Separable Case

# Problems with the Perceptron

- Noise: if the data isn't separable, weights might thrash
  - Averaging weight vectors over time can help (averaged perceptron)

- Mediocre generalization: finds a "barely" separating solution

- Overtraining: test / held-out accuracy usually rises, then falls
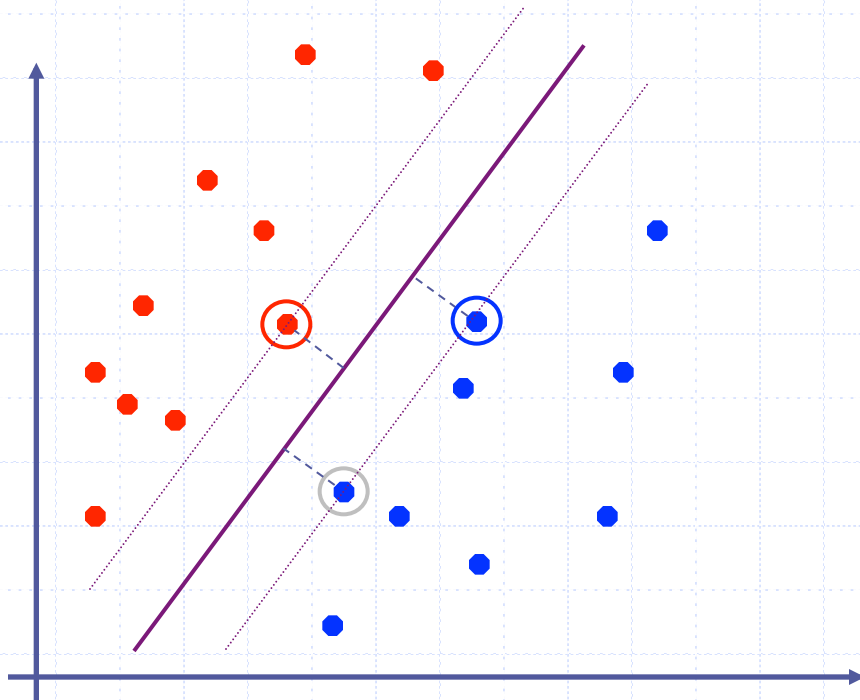  - Overtraining is a kind of overfitting

# Linear Separators

◆ Which of these linear separators is optimal?

# Support Vector Machines

- **Maximizing the margin:** good according to intuition, theory, practice
- Only support vectors matter; other training examples are ignorable
- Support vector machines (SVMs) find the separator with max margin
- Basically, SVMs are MIRA where you optimize over all examples at once

## MIRA

$$\min_{w} \ \frac{1}{2}||w - w'||^2$$

$$w_{y^*} \cdot f(x_i) \geq w_y \cdot f(x_i) + 1$$

## SVM

$$\min_{w} \ \frac{1}{2}||w||^2$$

$$\forall i, y \ w_{y^*} \cdot f(x_i) \geq w_y \cdot f(x_i) + 1$$

# Classification: Comparison

◆ Naïve Bayes
  - Builds a model training data
  - Gives prediction probabilities
  - Strong assumptions about feature independence
  - One pass through data (counting)

◆ Perceptrons :
  - Makes less assumptions about data
  - Mistake-driven learning
  - Multiple passes through data (prediction)
  - Often more accurate