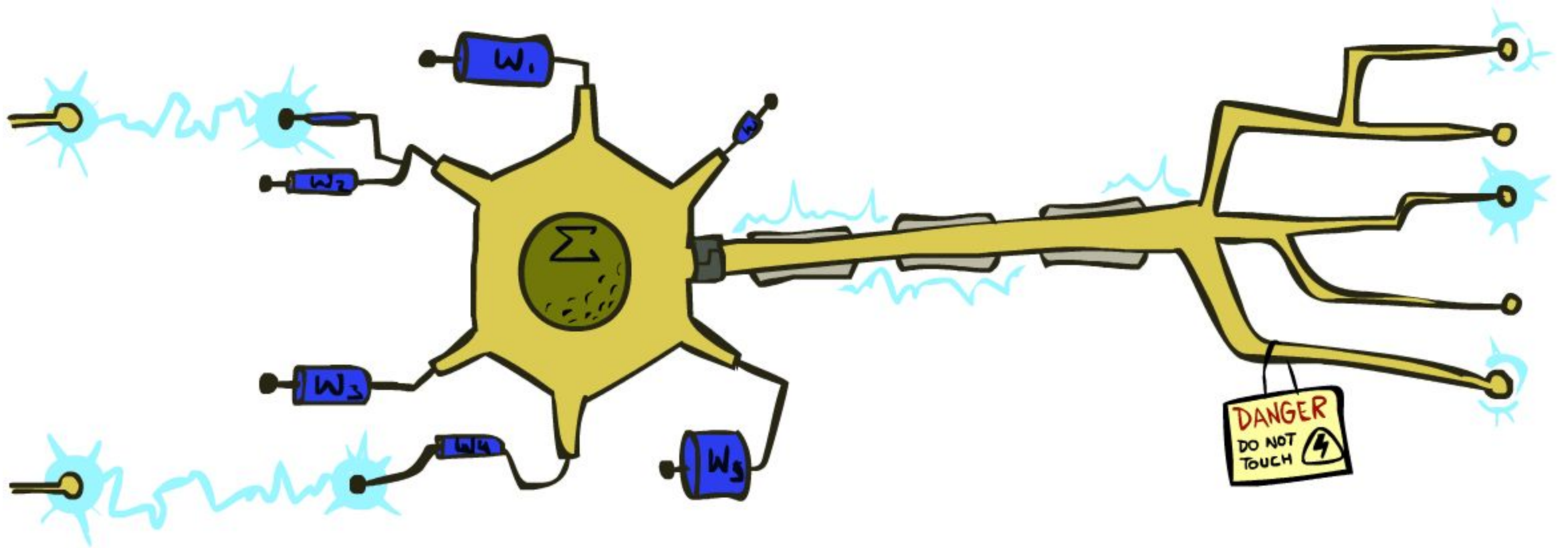# CS 188: Artificial Intelligence
## Perceptrons and Logistic Regression



Emma Pierson and Peyrin Kao

(Original slides from Pieter Abbeel & Dan Klein)

University of California, Berkeley

# Some final concepts from last class: Baselines

- First step: get a baseline
  - Baselines are very simple "straw man" procedures
  - Help determine how hard the task is
  - Help know what a "good" accuracy is

- Weak baseline: most frequent label classifier
  - Gives all test instances whatever label was most common in the training set
  - E.g. for spam filtering, might label everything as ham
  - Accuracy might be very high if the problem is skewed
  - E.g. calling everything "ham" gets 66%, so a classifier that gets 70% isn't very good…

- For real research, usually use previous work as a (strong) baseline

# Some final concepts from last class: Confidences

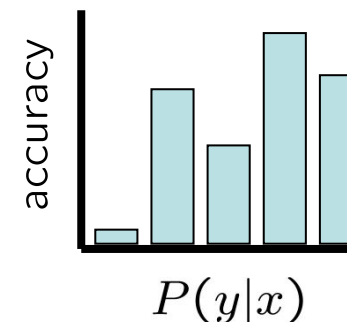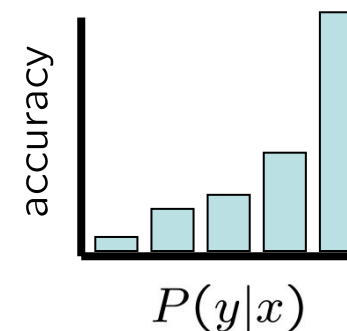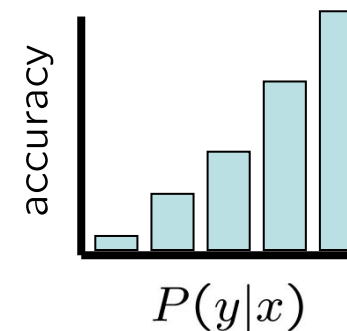- The confidence of a probabilistic classifier:
  - Posterior probability of the top label

  $$\text{confidence}(x) = \max_y P(y|x)$$

  - Represents how sure the classifier is of the classification
  - Any probabilistic model will have confidences
  - No guarantee confidence is correct

- Calibration
  - Weak calibration: higher confidences mean higher accuracy
  - Strong calibration: confidence predicts accuracy rate
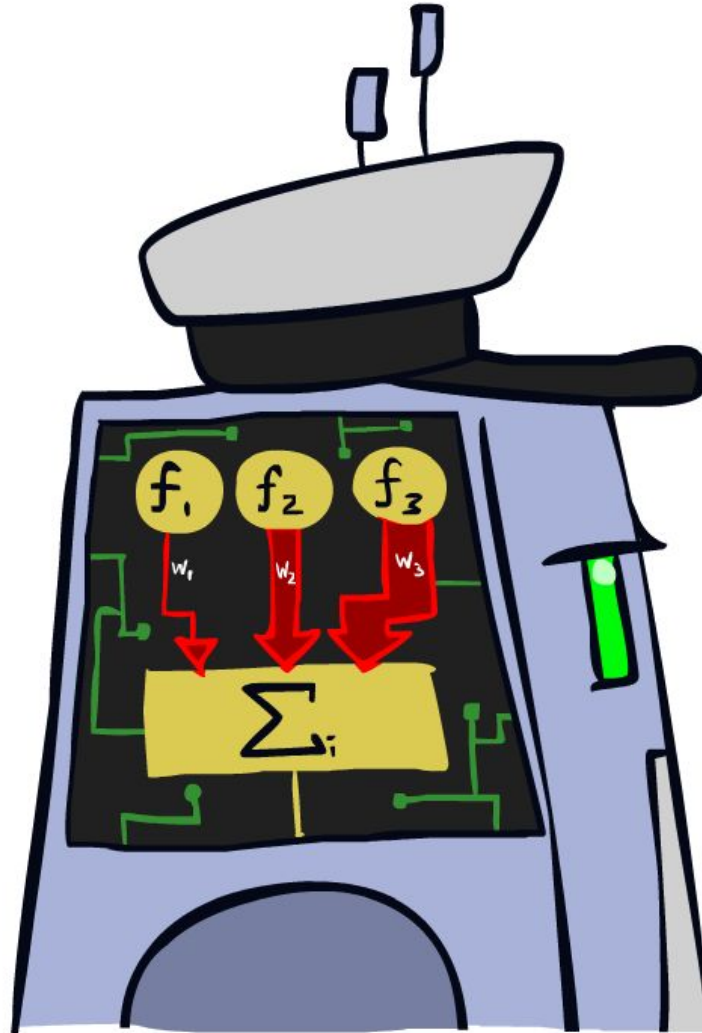  - What's the value of calibration?

# Review from last class

- Machine learning: estimate a model from data
    - e.g., estimate parameters of a Naïve Bayes model
- Specific machine learning task we often want to perform: **classification**
    - Last class we discussed a particular approach to classification: Naïve Bayes
    - This class we'll discuss others: perceptrons, logistic regression
- We also discussed more general machine learning concepts (foundational, not particular to a specific type of model)
    - Features
    - Overfitting
    - Parameters and hyperparameters
    - Train / held-out / test split
    - Regularization / smoothing
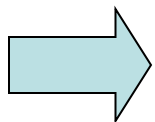    - Maximum likelihood estimate
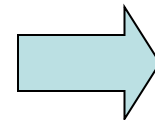
# Linear Classifiers

# 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
```

$\Rightarrow$

```
# free       : 2
YOUR_NAME    : 0
MISSPELLED   : 2
FROM_FRIEND  : 0
...
```
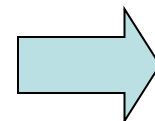
$\Rightarrow$

SPAM
or
+
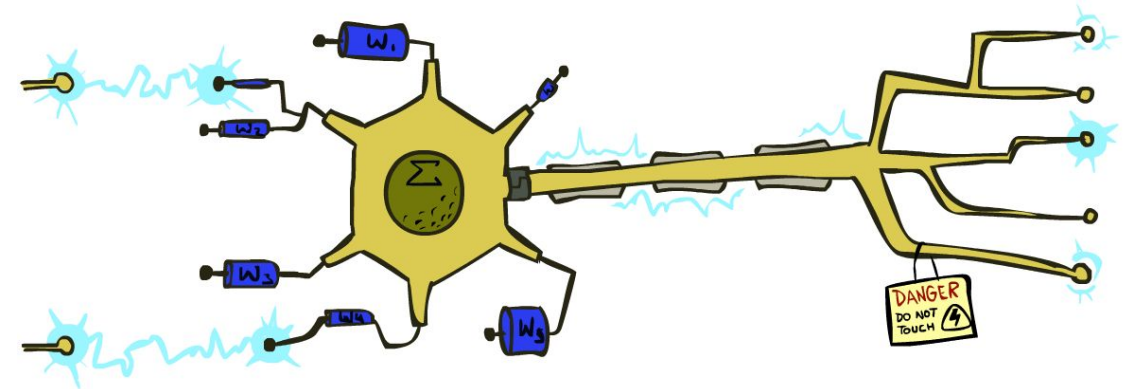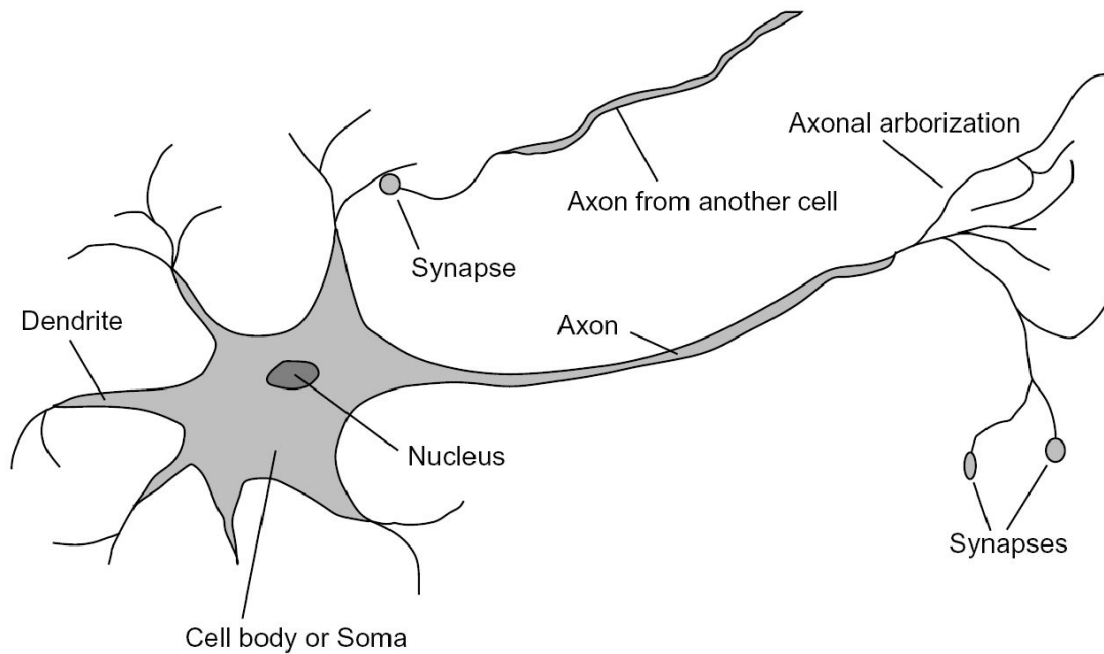


$\Rightarrow$

```
PIXEL-7,12  : 1
PIXEL-7,13  : 0
...
NUM_LOOPS   : 1
...
```
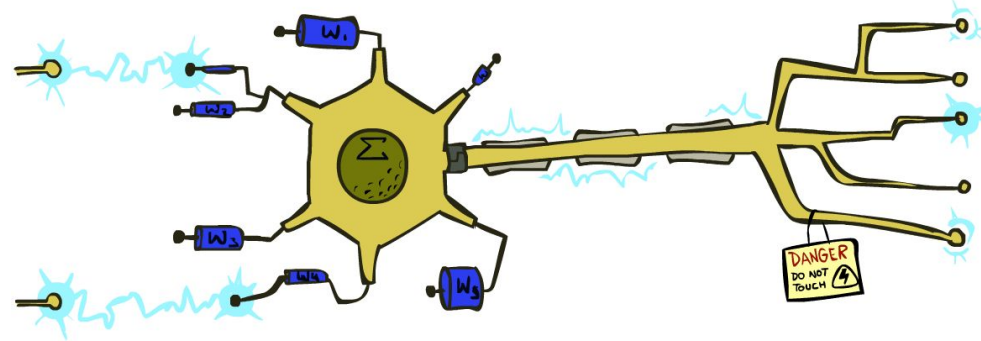
$\Rightarrow$

"2"

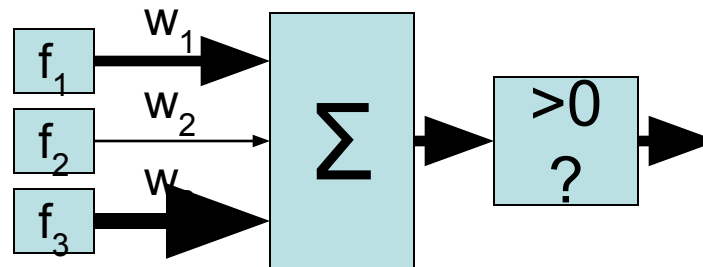# Some (Simplified) Biology

- Very loose inspiration: human neurons

# Linear Classifiers

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

$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$
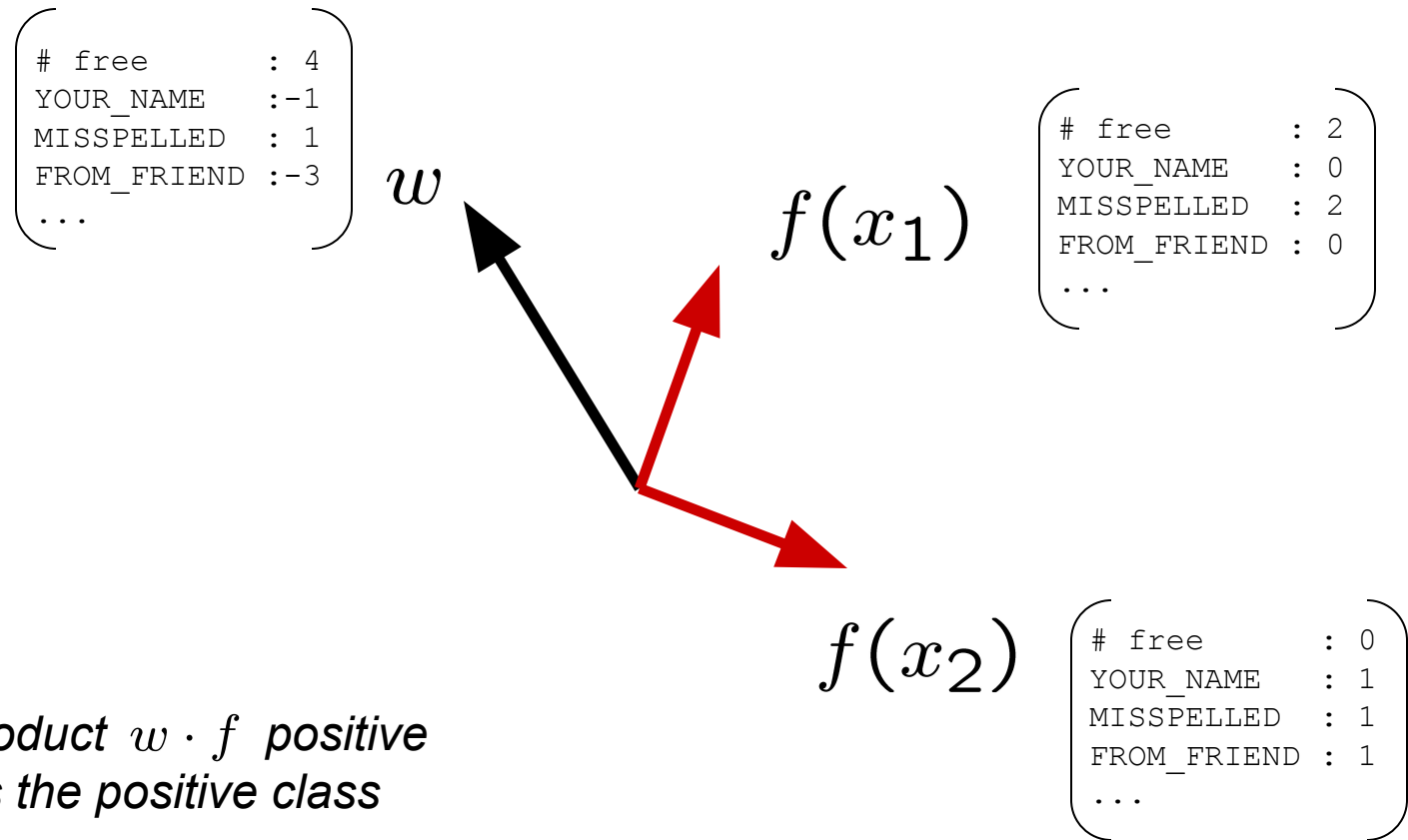
- If the activation is:
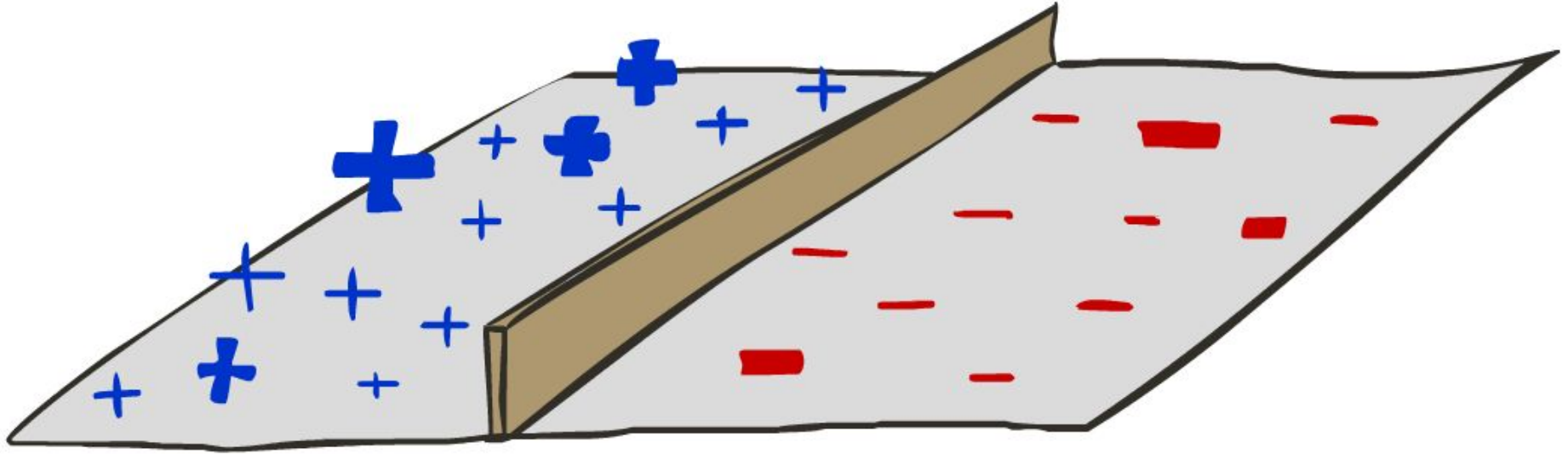  - Positive, output +1
  - Negative, output -1

# Weights

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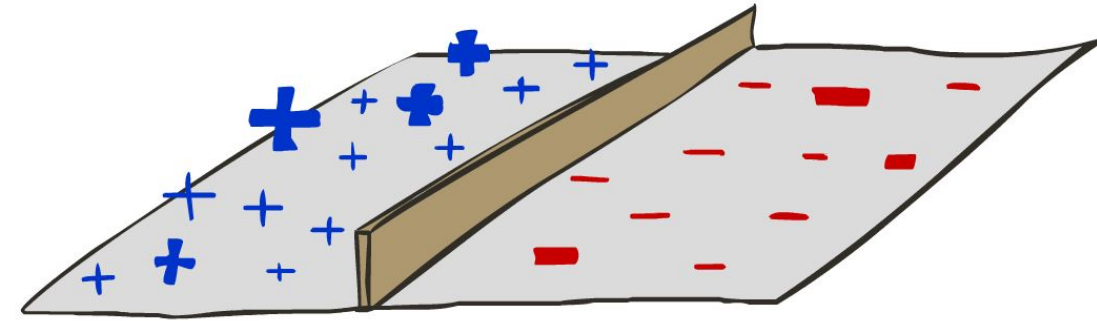- Learning: figure out the weight vector from examples

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

$w$

$f(x_1)$

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

$f(x_2)$

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

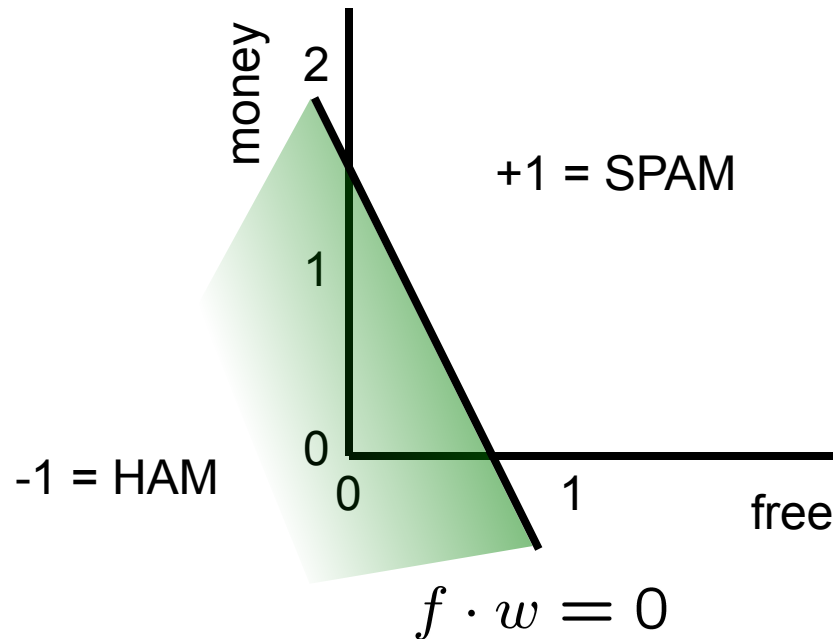*Dot product $w \cdot f$ positive means the positive class*

# Decision Rules

# Binary Decision Rule

- In the space of feature vectors
  - Examples are points
  - Any weight vector encodes a hyperplane
  - One side corresponds to Y=+1
  - Other corresponds to Y=-1

$w$

```
BIAS  : -3
free  :  4
money :  2
...
```
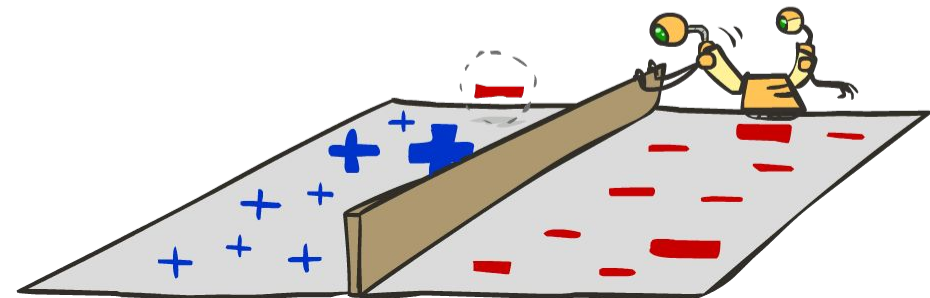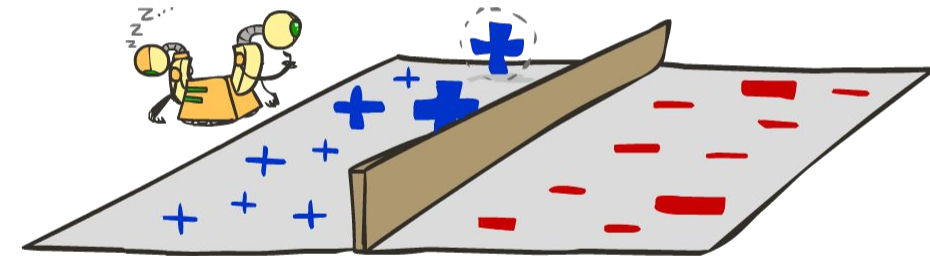
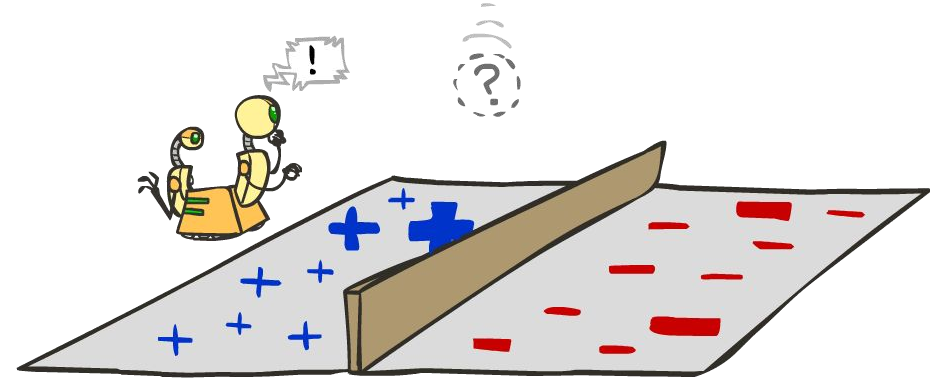+1 = SPAM

-1 = HAM

money

free

$f \cdot w = 0$

# Weight Updates

# Learning: Binary Perceptron

- Start with weights = 0
- For each training instance:
  - Classify with current weights

  - If correct (i.e., y=y*), no change!

  - If wrong: adjust the weight vector
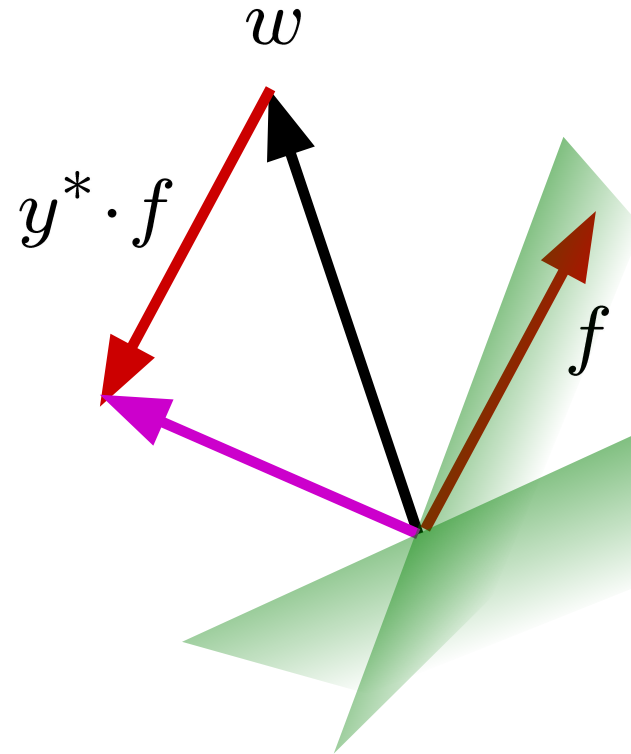
# Learning: Binary Perceptron

- 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.
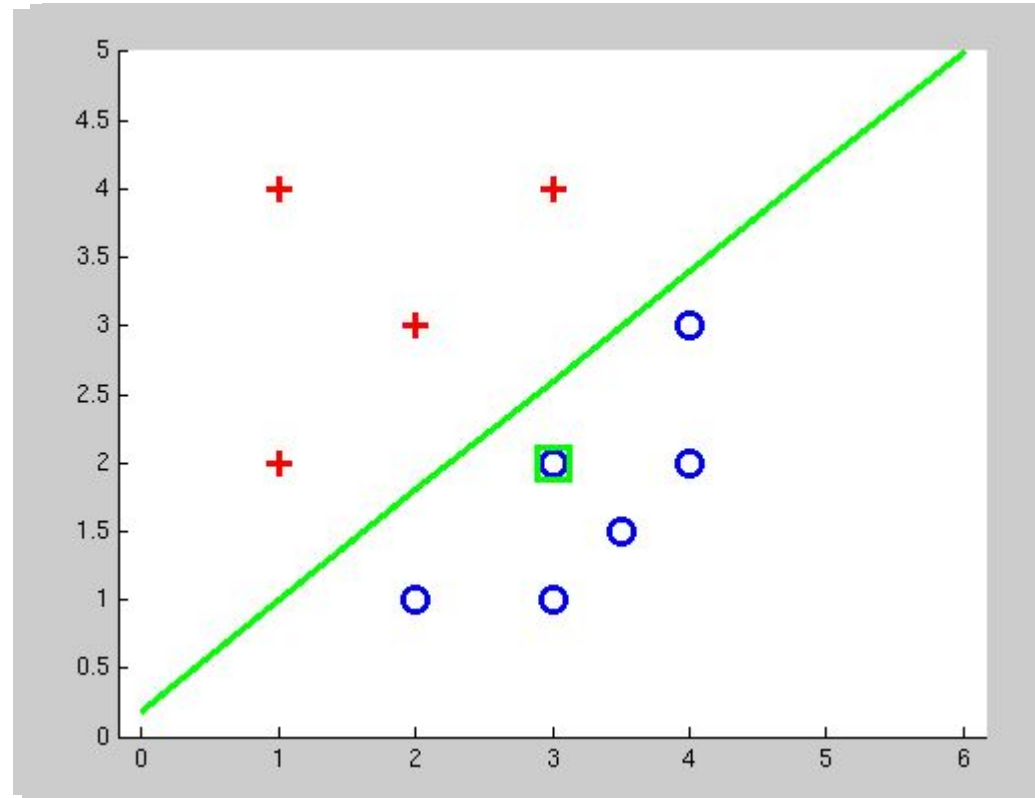
$$w_{\text{new}} = w_{\text{old}} + y^* \cdot f$$

$$\text{score}_{\text{new}} = (w_{\text{old}} + y^* \cdot f) \cdot f$$

# Examples: Perceptron

- Separable Case

# Multiclass Decision Rule

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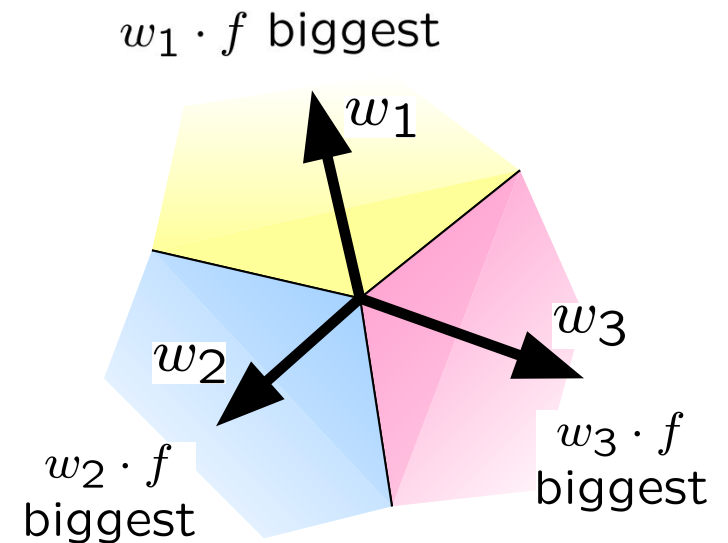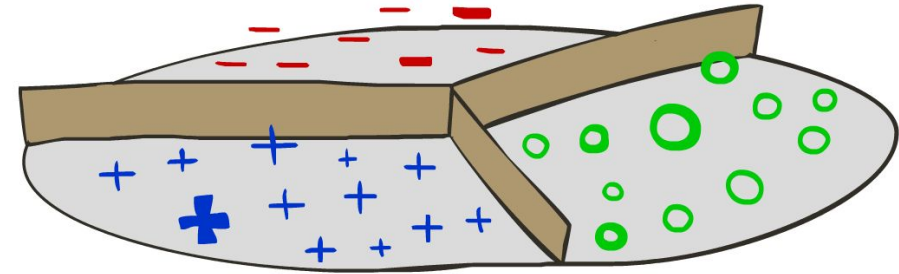  - A weight vector for each class:

  $$w_y$$

  - 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

$w_1$

$w_3$

$w_2$

$w_2 \cdot f$
biggest

$w_3 \cdot f$
biggest

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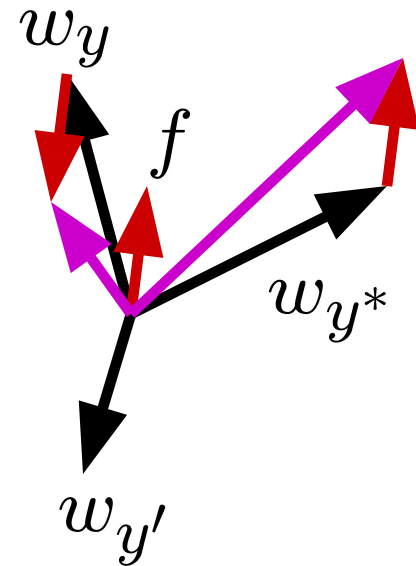
# Learning: Multiclass Perceptron

- 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 (y), raise score of right answer (y*)

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

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

# Example: Multiclass Perceptron

"win the vote"

"win the election"

"win the game"

$$w_{SPORTS}$$

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

$$w_{POLITICS}$$

```
BIAS  : 0
win   : 0
game  : 0
vote  : 0
the   : 0
...
```
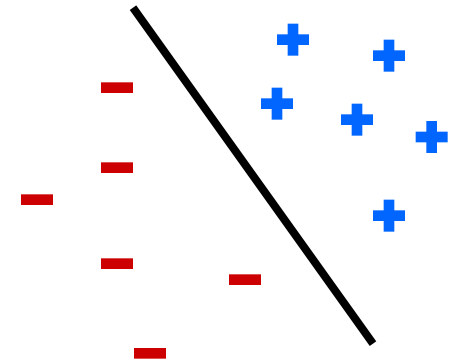
$$w_{TECH}$$

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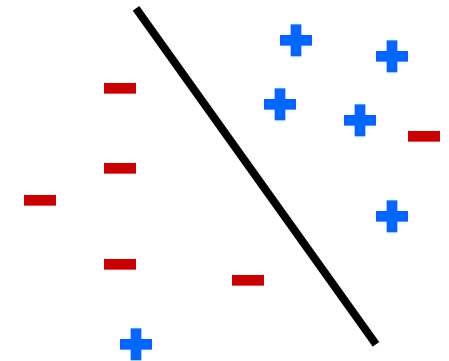
# Properties of Perceptrons

- Separability: true if 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

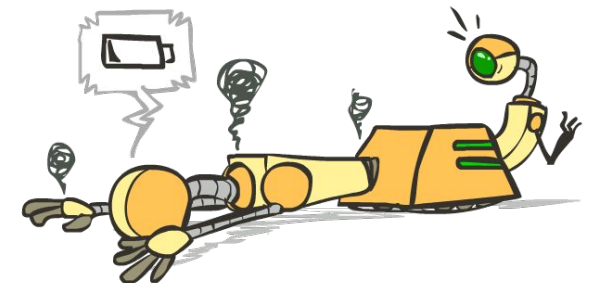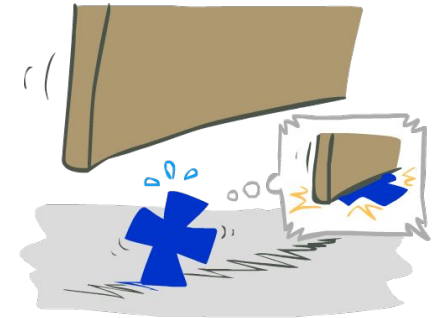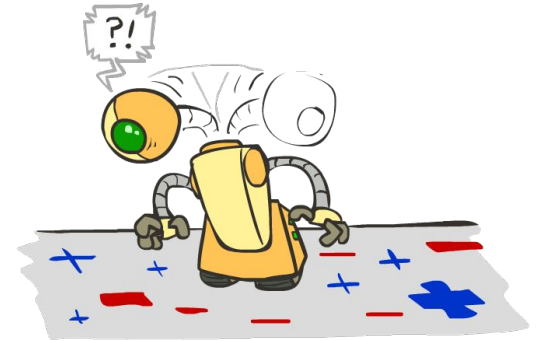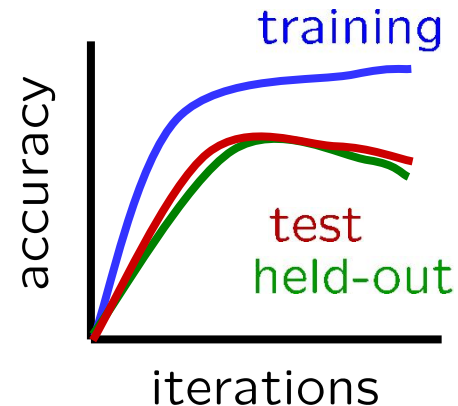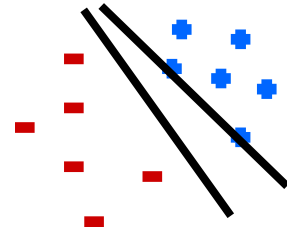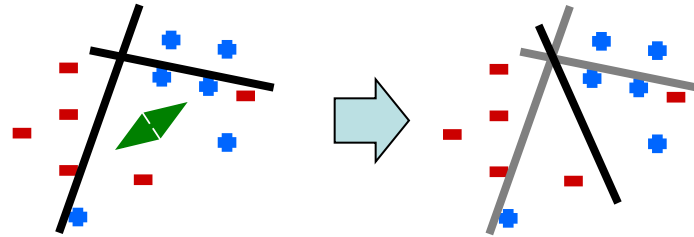$$\text{mistakes} < \frac{k}{\delta^2}$$
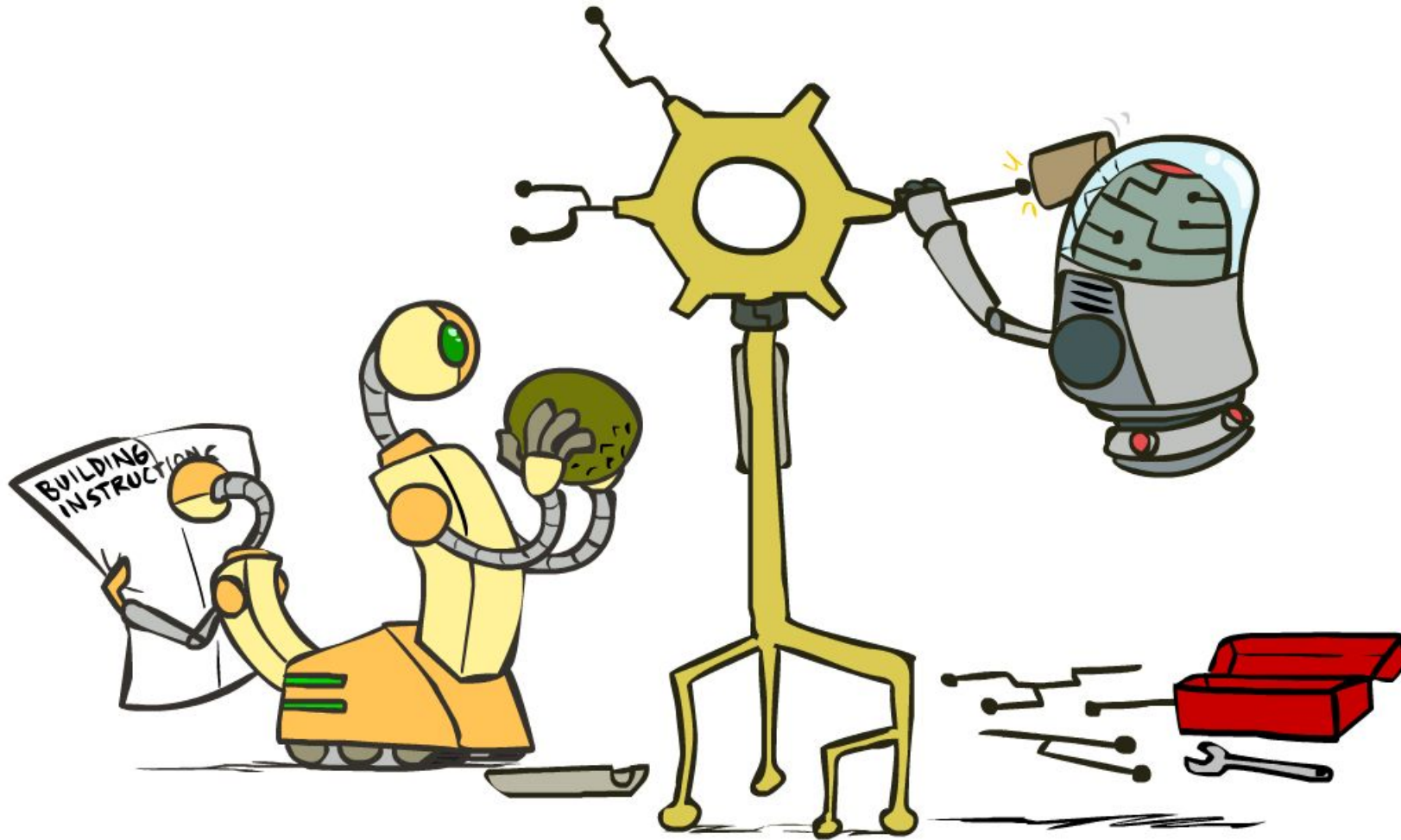
Separable



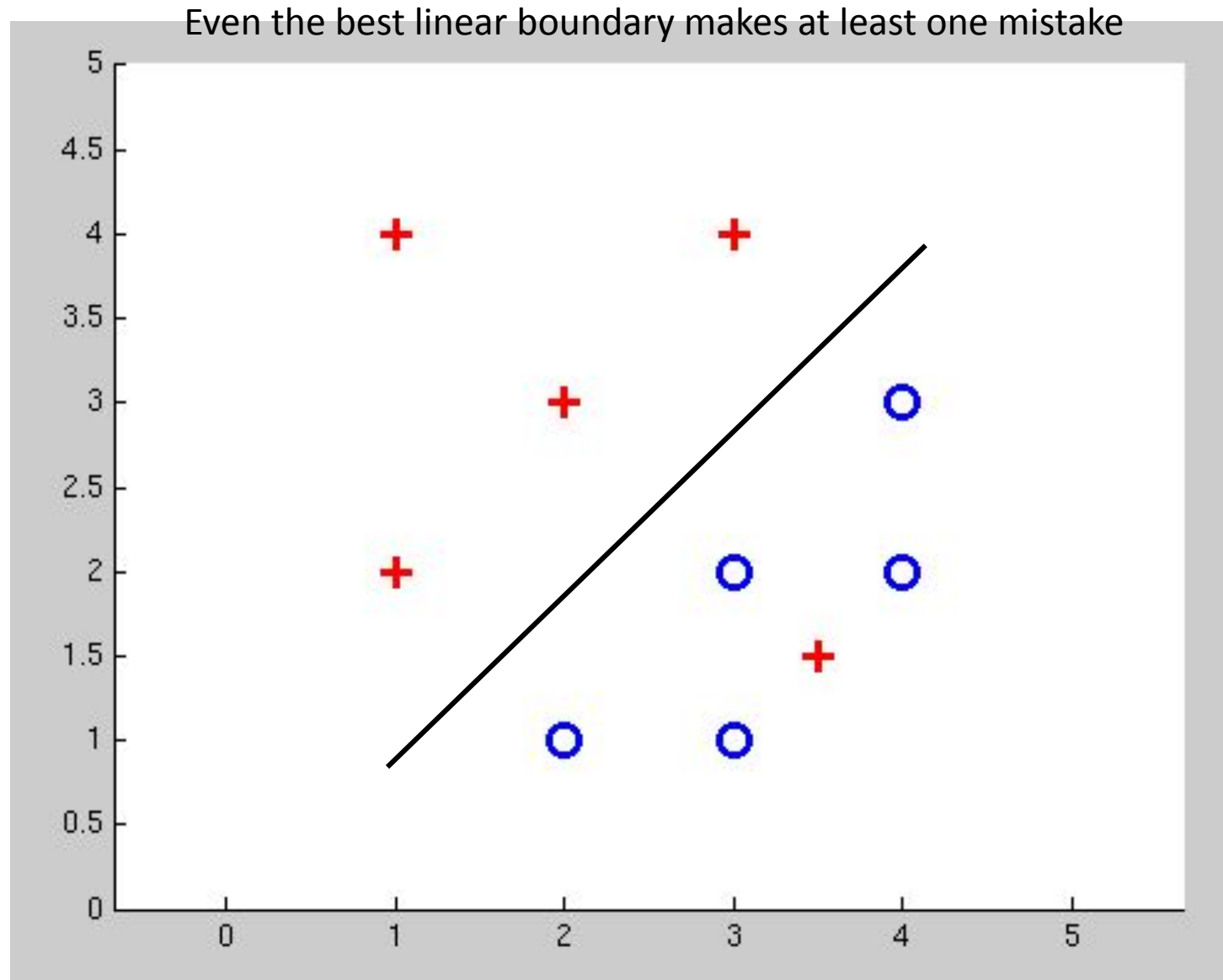Non-Separable

# Problems with the Perceptron

- Noise: if the data isn't separable, weights might thrash (bounce around between solutions)
  - 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
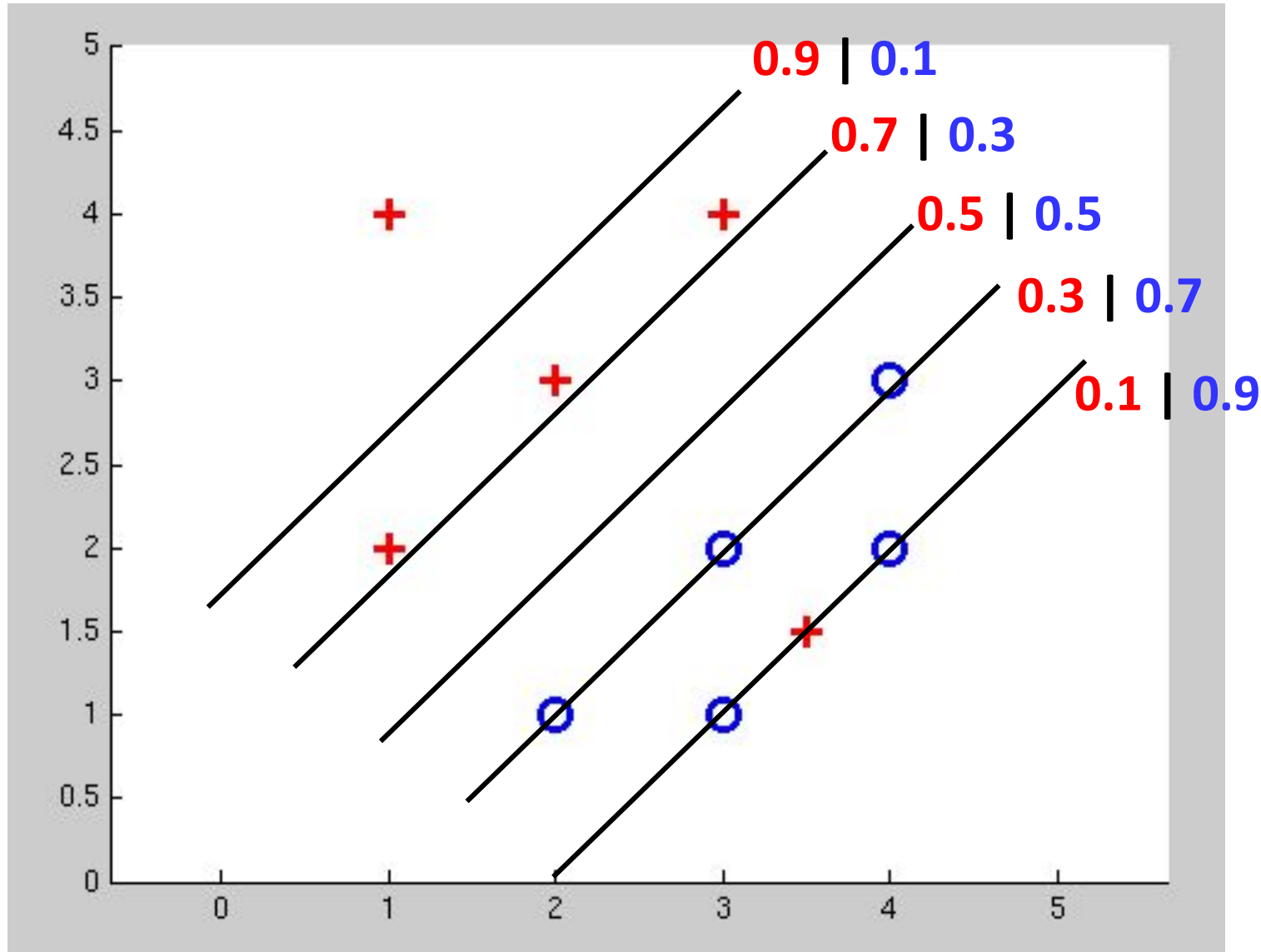  - Overtraining is a kind of overfitting

# Non-Separable Case: Deterministic Decision



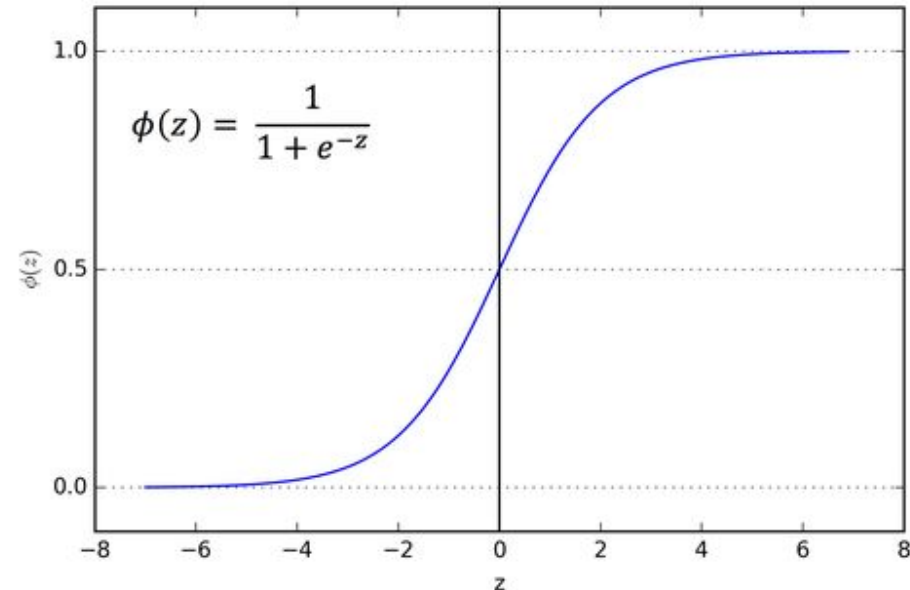Even the best linear boundary makes at least one mistake

# How to get probabilistic decisions?

- Perceptron scoring: $z = w \cdot f(x)$
- If $z = w \cdot f(x)$ very positive □ want probability going to 1
- If $z = w \cdot f(x)$ very negative □ want probability going to 0

- Sigmoid function

$$\phi(z) = \frac{1}{1 + e^{-z}}$$

# Best w?

- Maximum likelihood estimation: maximize probability of labels *y* given features *x*, parameters *w*.

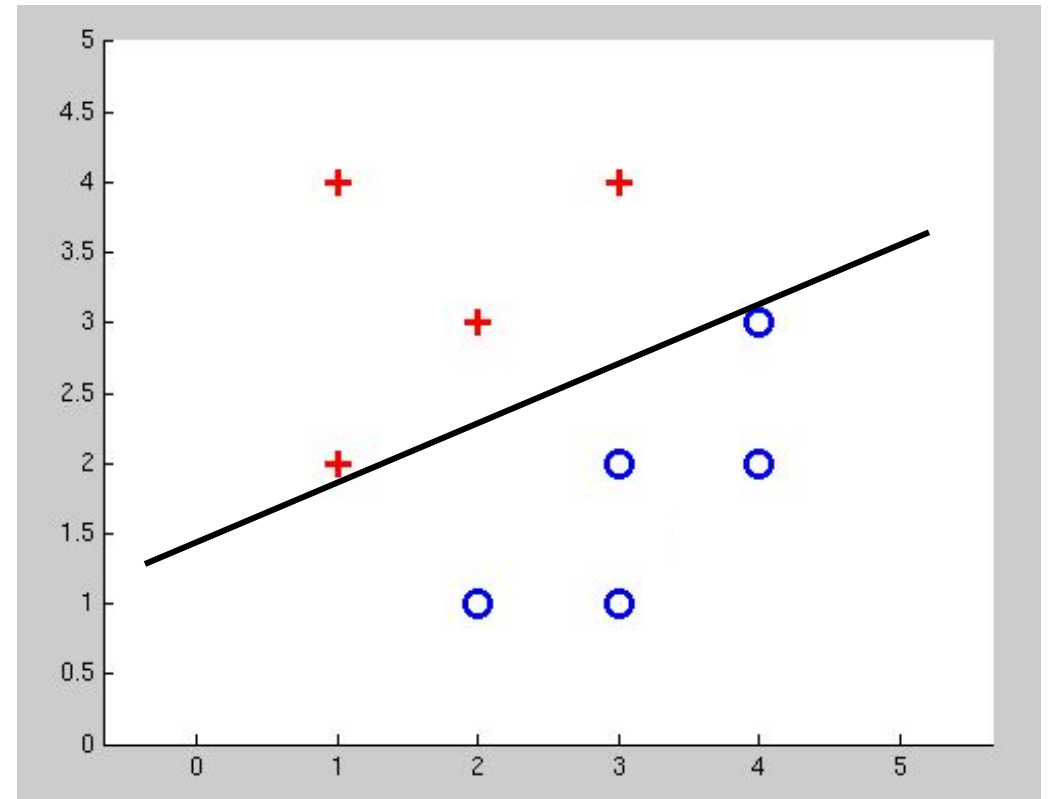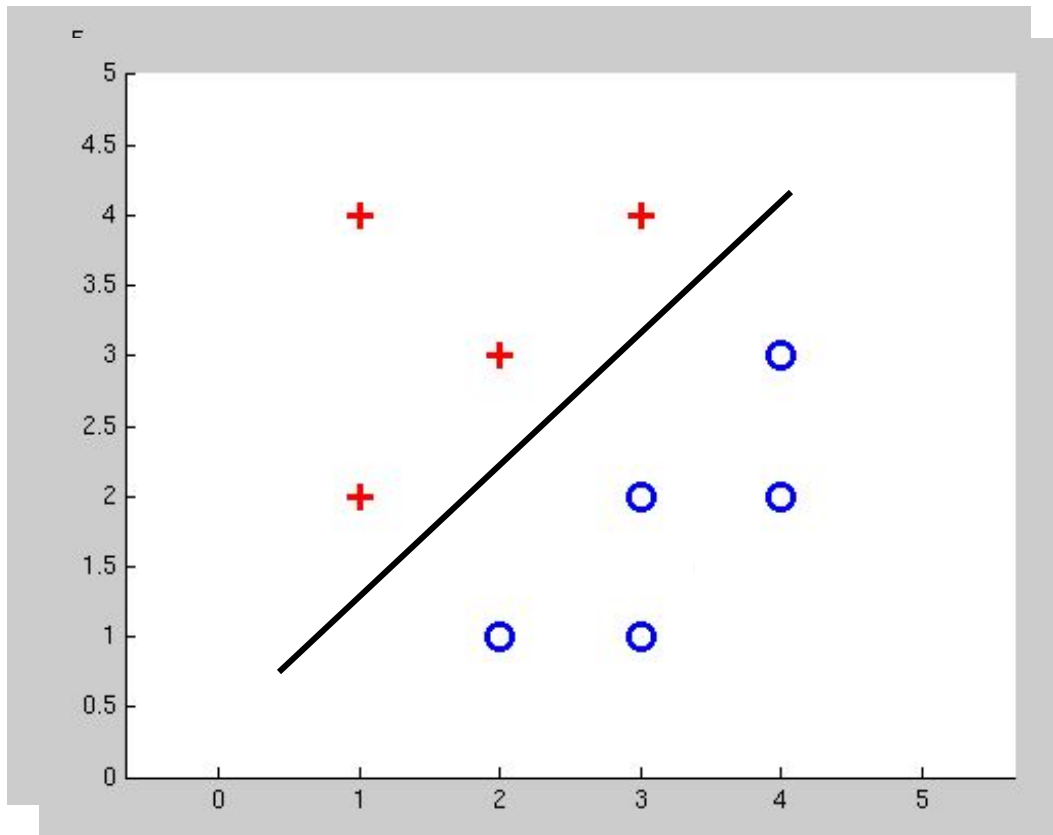$$\max_w \; ll(w) = \max_w \; \sum_i \log P(y^{(i)}|x^{(i)}; w)$$

with:

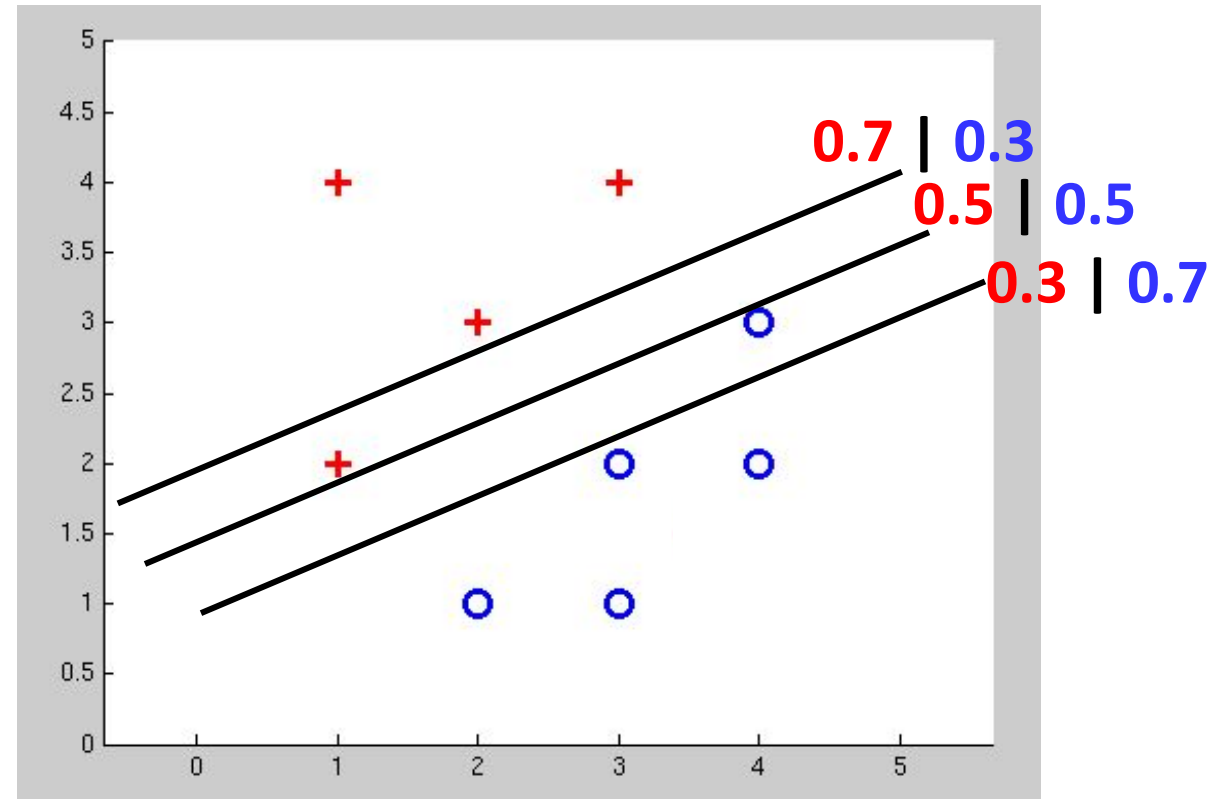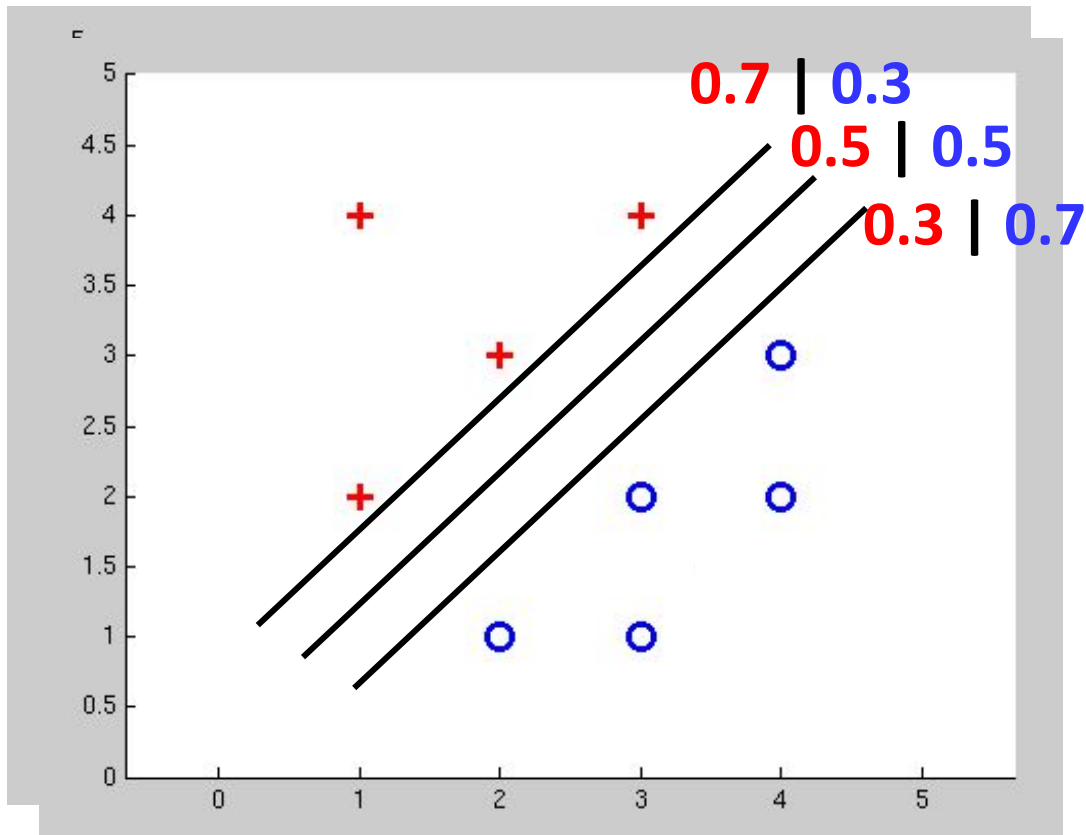$$P(y^{(i)} = +1|x^{(i)}; w) = \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

$$P(y^{(i)} = -1|x^{(i)}; w) = 1 - \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

**= Logistic Regression**
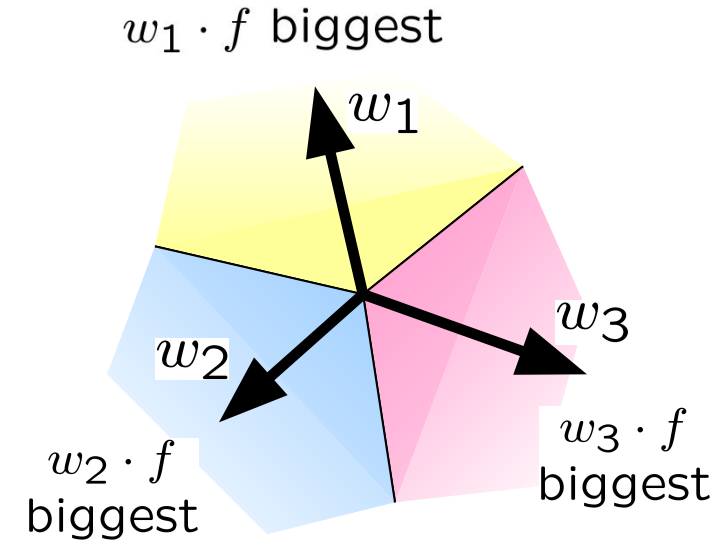
# Separable Case: Probabilistic Decision – Clear Preference

# Multiclass Logistic Regression

$w_1 \cdot f$ biggest

- Recall Perceptron:

  - A weight vector for each class: $\quad w_y$

  - Score (activation) of a class y: $\quad w_y \cdot f(x)$

  - Prediction highest score wins

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

$w_2 \cdot f$ biggest

$w_3 \cdot f$ biggest

- How to make the scores into probabilities?

$$z_1, z_2, z_3 \rightarrow \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

original activations

softmax activations

# Best w?

- Maximum likelihood estimation:

$$\max_{w} \quad ll(w) = \max_{w} \quad \sum_{i} \log P(y^{(i)}|x^{(i)};w)$$

with: $$P(y^{(i)}|x^{(i)};w) = \frac{e^{w_{y^{(i)}} \cdot f(x^{(i)})}}{\sum_{y} e^{w_y \cdot f(x^{(i)})}}$$

**= Multi-Class Logistic Regression**

# Next Lecture

- ## Optimization

  - i.e., how do we solve:

$$\max_w \; ll(w) = \max_w \; \sum_i \log P(y^{(i)}|x^{(i)}; w)$$