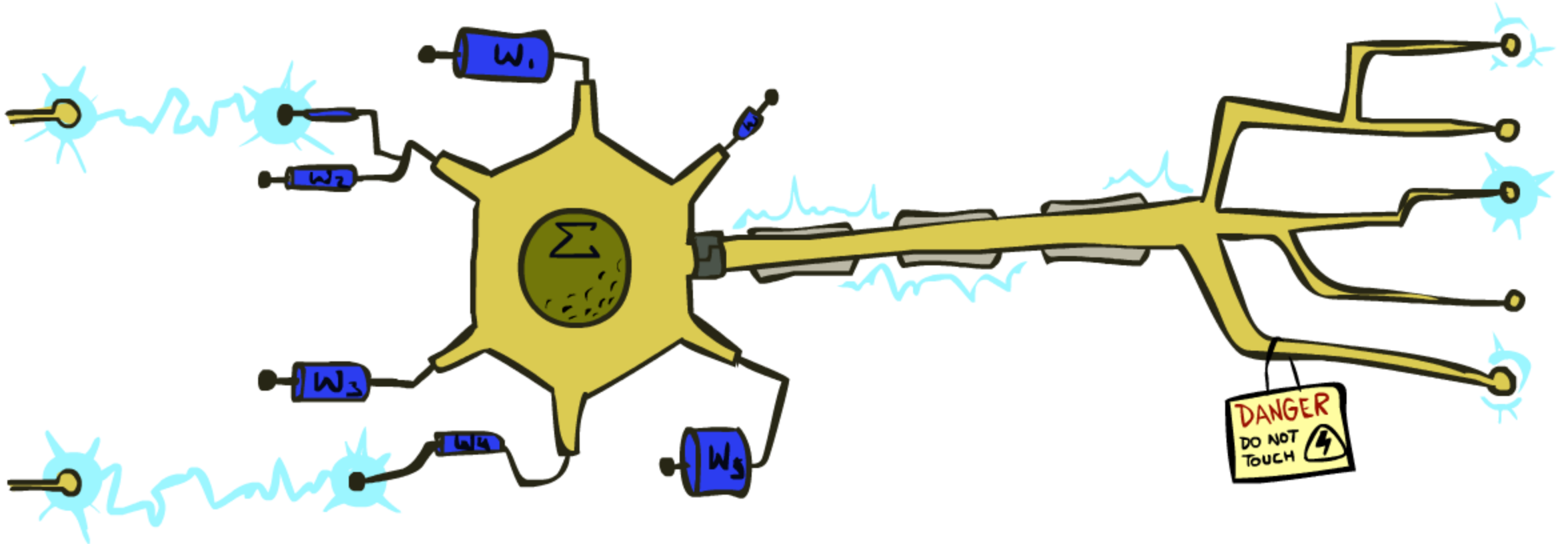


# CS 188: Artificial Intelligence

## Perceptrons and Logistic Regression



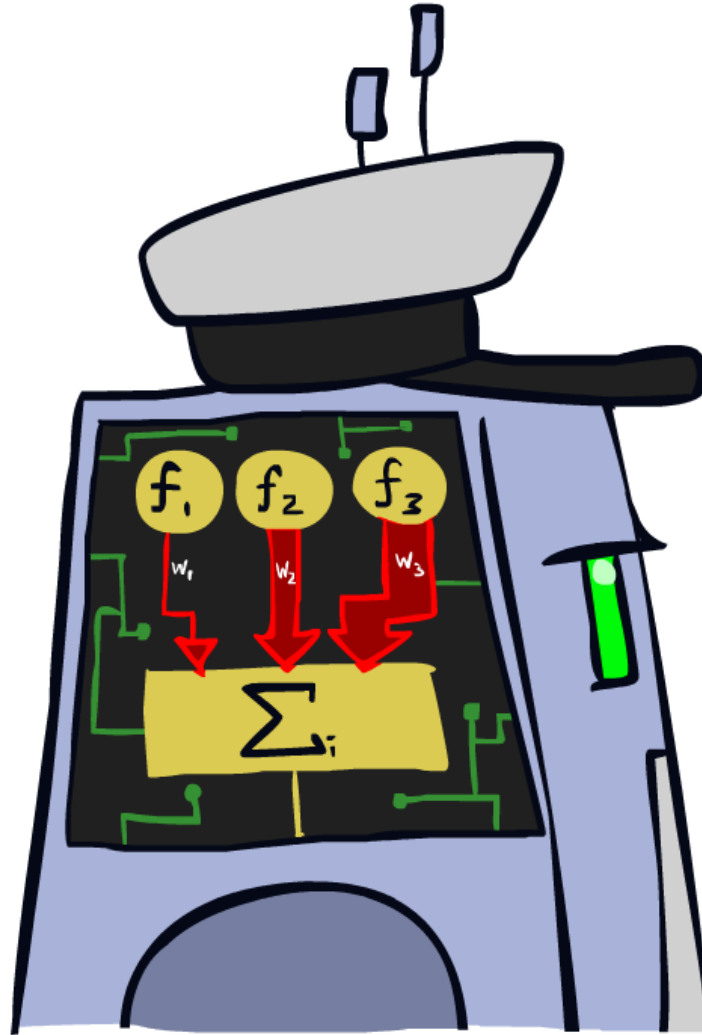
Instructors: Oliver Grillmeyer — University of California, Berkeley

# Announcements

---

- HW8 is due **Thursday, July 31, 11:59 PM PT**
- Project 4 is due **Friday, August 1, 11:59 PM PT**
- HW9 is due **Tuesday, August 5, 11:59 PM PT**
- HW10 is due **Thursday, August 7, 11:59 PM PT**
- Ignore assessment on HWs part B, but please show your work
- Final Exam is **Wednesday, August 13, 7-10 PM PT**
  - Accommodation requests by Wednesday, July 30, 11:59 PM PT

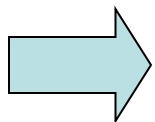
# Linear Classifiers



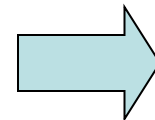
# Feature Vectors

 $x$  $f(x)$  $y$ 

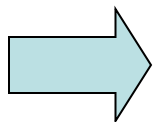
```
Hello,  
  
Do you want free print  
cartridges? Why pay more  
when you can get them  
ABSOLUTELY FREE! Just
```



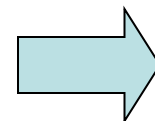
```
# free      : 2  
YOUR_NAME   : 0  
MISPELLED   : 2  
FROM_FRIEND : 0  
...
```



SPAM  
or  
+



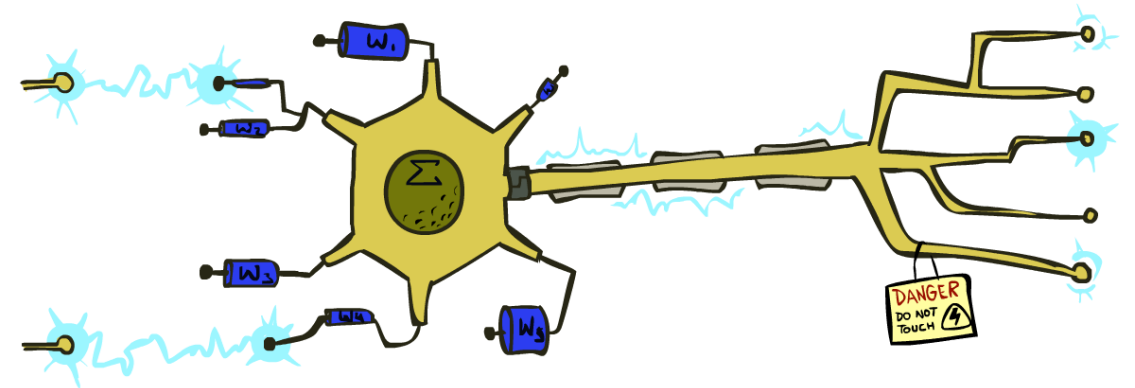
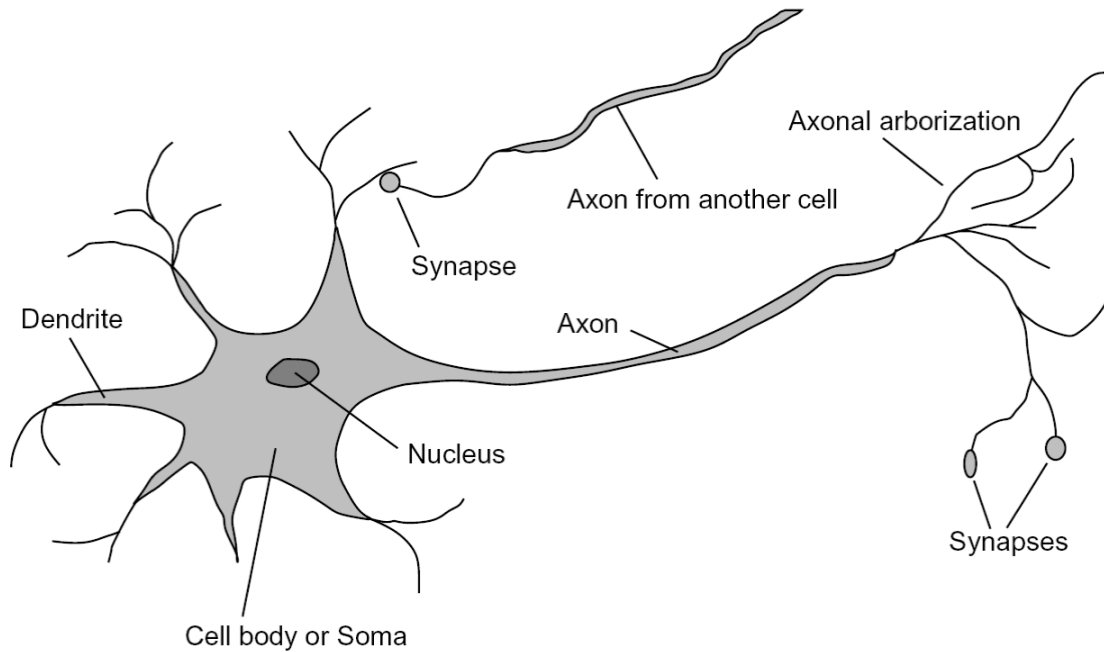
```
PIXEL-7,12 : 1  
PIXEL-7,13 : 0  
...  
NUM_LOOPS  : 1  
...
```



"2"

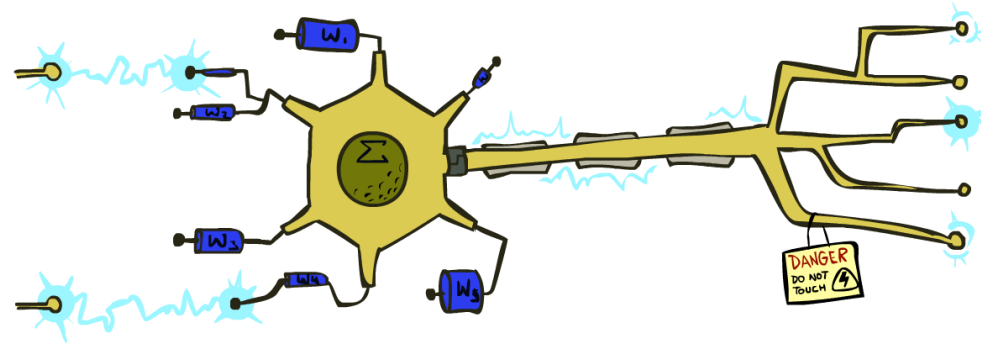
# Some (Simplified) Biology

- Very loose inspiration: human neurons



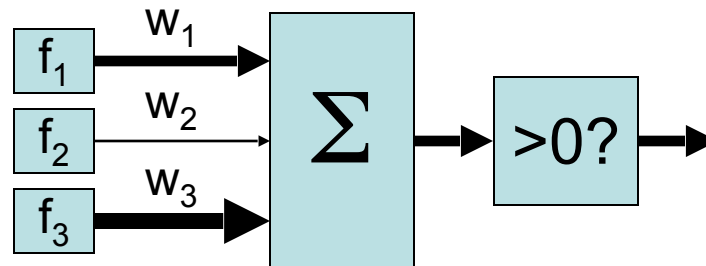
# Linear Classifiers

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



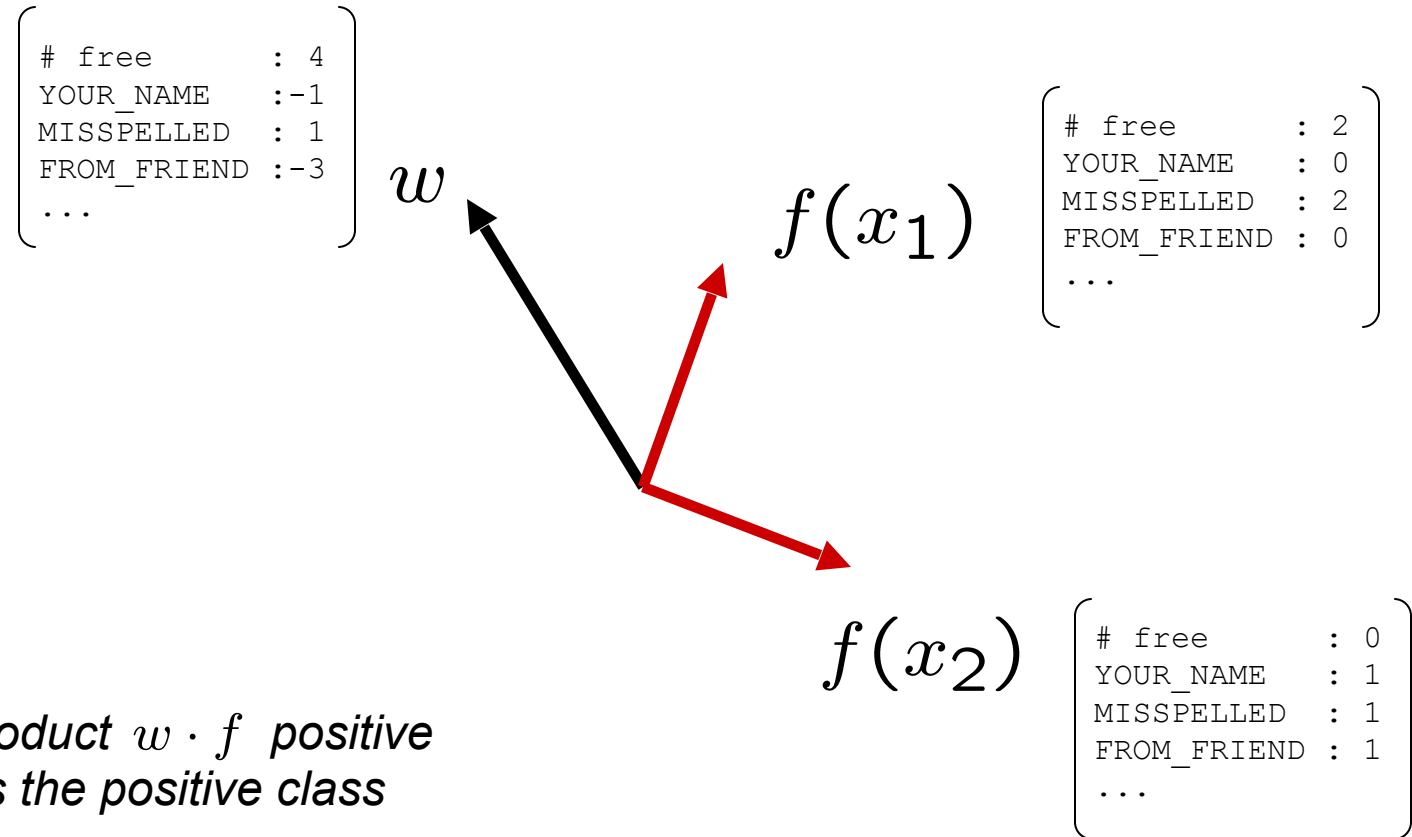
$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

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



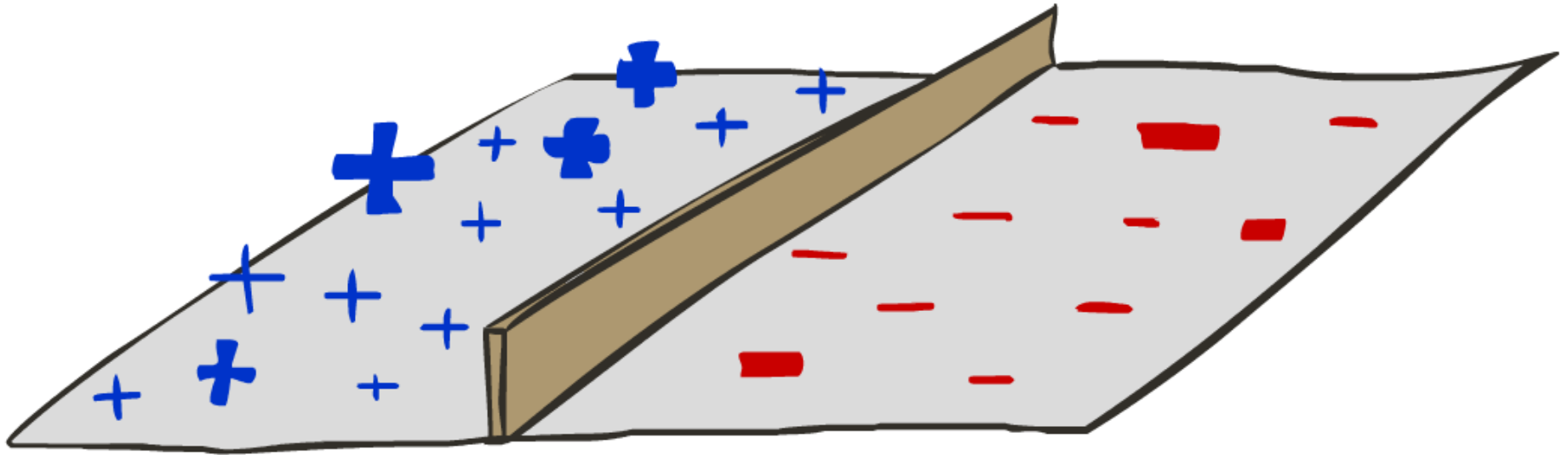
# Weights

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



# Decision Rules

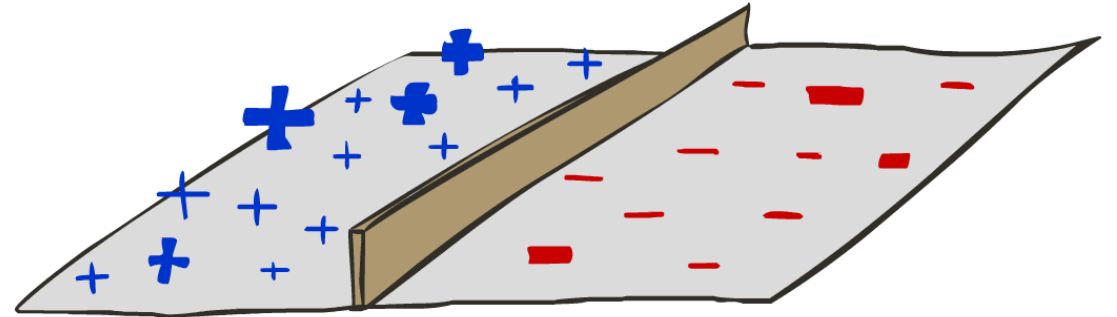
---





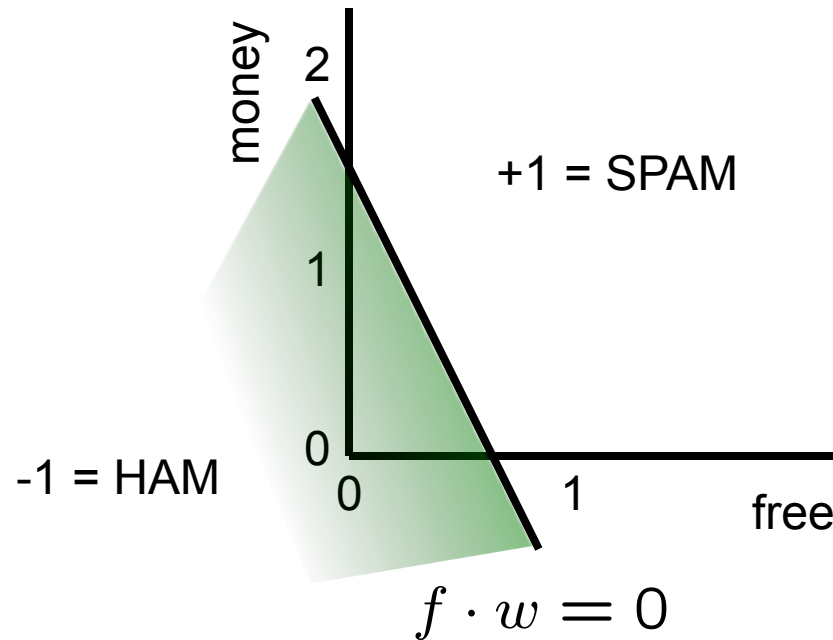
# Binary Decision Rule

- 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
...		



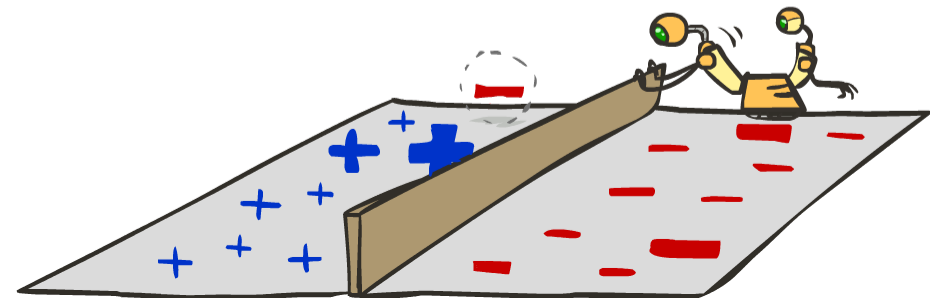
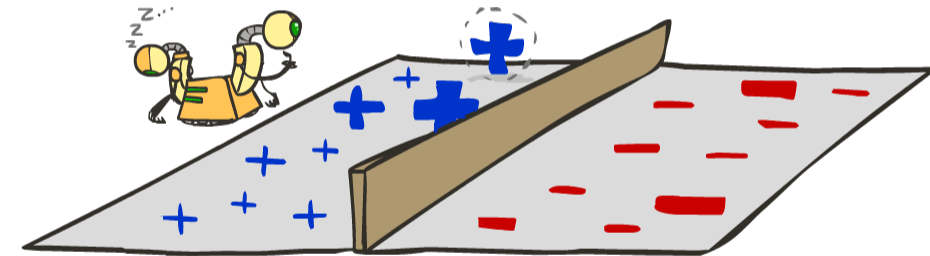
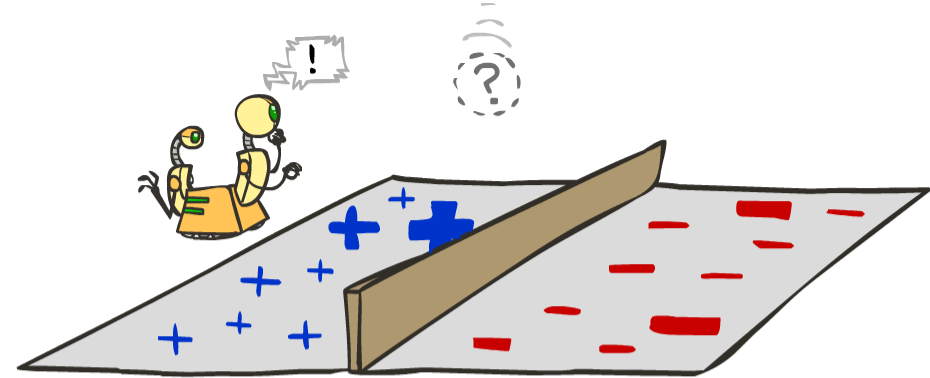
# Weight Updates

---



# 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



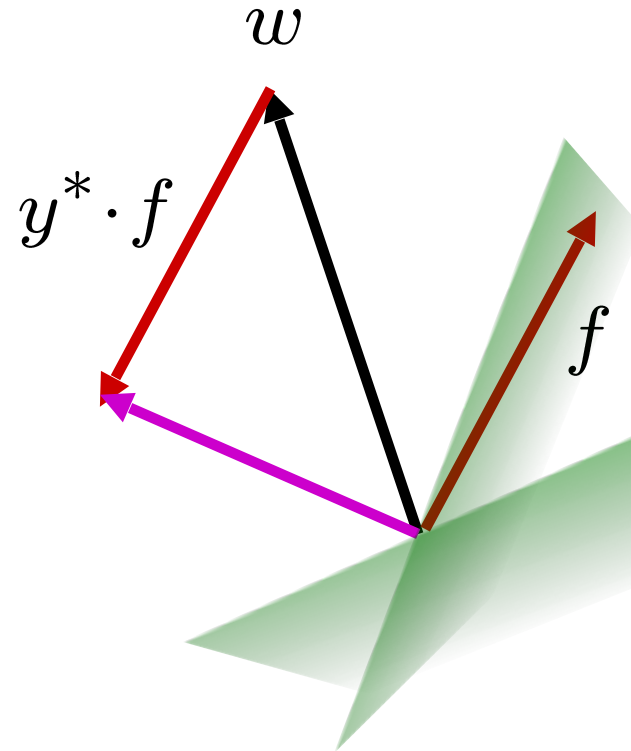
# Learning: Binary Perceptron

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

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

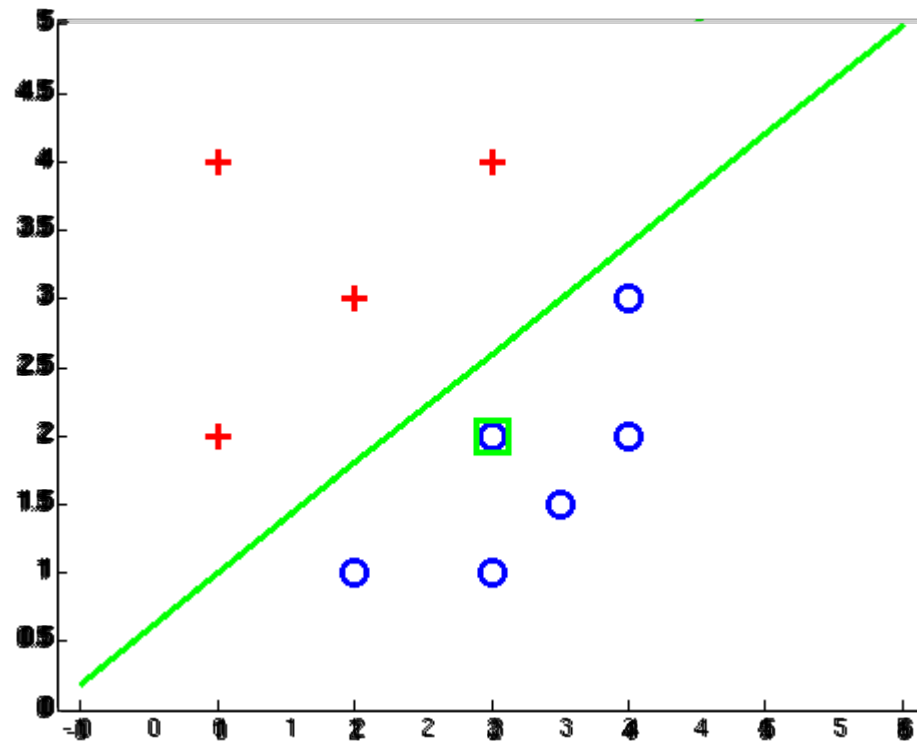
- If correct (i.e.,  $y=y^*$ ), no change!
- If wrong: adjust the weight vector by adding or subtracting the feature vector. Subtract if  $y^*$  is -1.

$$w = w + y^* \cdot f$$



# Examples: Perceptron

## ■ Separable Case



# Multiclass Decision Rule

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

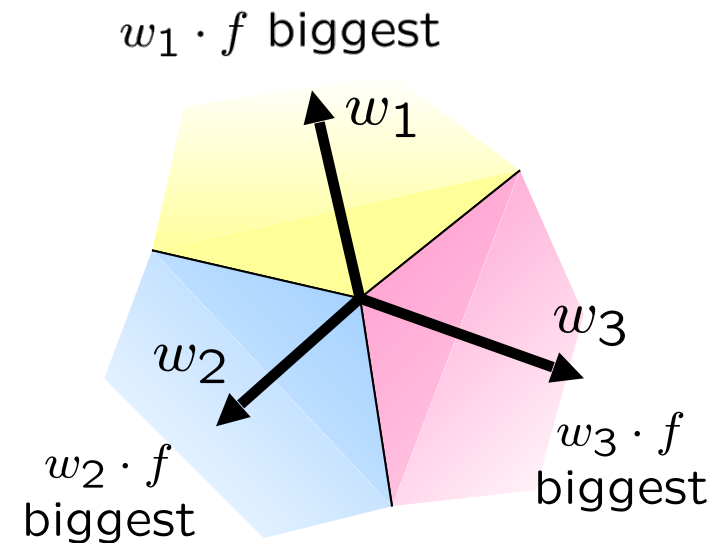
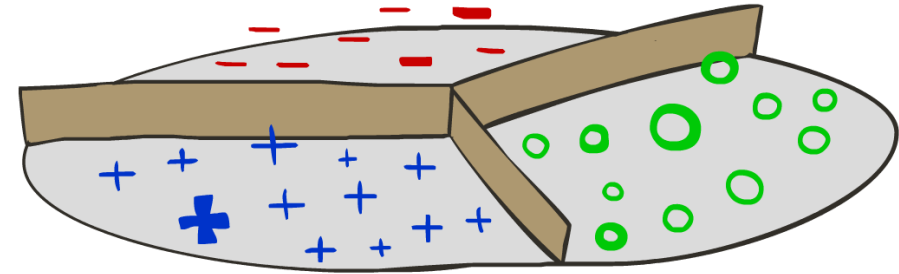
$$w_y$$

- Score (activation) of a class  $y$ :

$$w_y \cdot f(x)$$

- Prediction highest score wins

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



*Binary = multiclass where the negative class has weight zero*

# Learning: Multiclass Perceptron

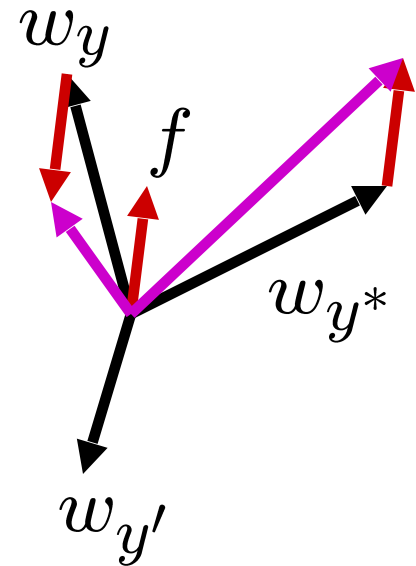
- Start with all weights = 0
- Pick up training examples one by one
- Predict with current weights

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

- If correct, no change!
- If wrong: lower score of wrong answer, raise score of right answer

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

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



# Example: Multiclass Perceptron

“win the vote”

“win the election”

“win the game”

Sample A  $w_S \cdot f_A = 1; w_P \cdot f_A = 0; w_T \cdot f_A = 0$

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

Sample B  $w_S \cdot f_B = -2; w_P \cdot f_B = 3; w_T \cdot f_B = 0$

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

Sample C  $w_S \cdot f_C = -2; w_P \cdot f_C = 3; w_T \cdot f_C = 0$

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

$w_{SPORTS}$

	fA	wS	fC	wS
BIAS : 1	1	0	1	1
win : 0	1	-1	1	0
game : 0	- 0 =	0	+ 1 =	1
vote : 0	1	-1	0	-1
the : 0	1	-1	1	0
...		...		...

$w_{POLITICS}$

	fA	wP	fC	wP
BIAS : 0	1	1	1	0
win : 0	1	1	1	0
game : 0	+ 0 =	0	- 1 =	-1
vote : 0	1	1	0	1
the : 0	1	1	1	0
...		...		...

$w_{TECH}$

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

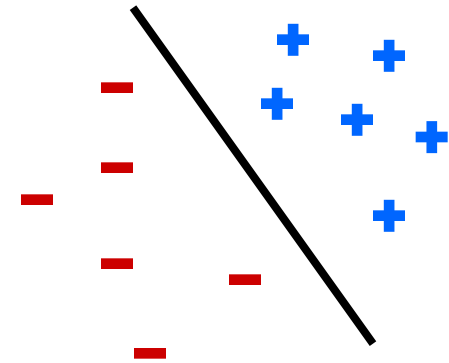


# Properties of Perceptrons

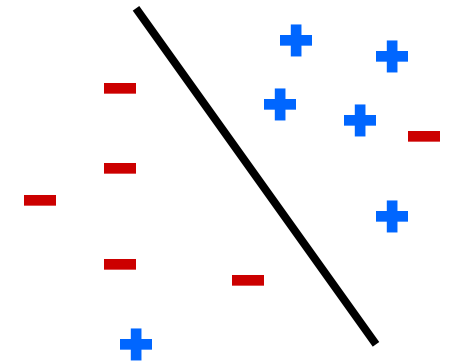
- Separability: true if some parameters get the training set perfectly correct
- Convergence: if the training is separable, perceptron will eventually converge (binary case)
- Mistake Bound: the maximum number of mistakes (binary case) related to the *margin* or degree of separability

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

Separable

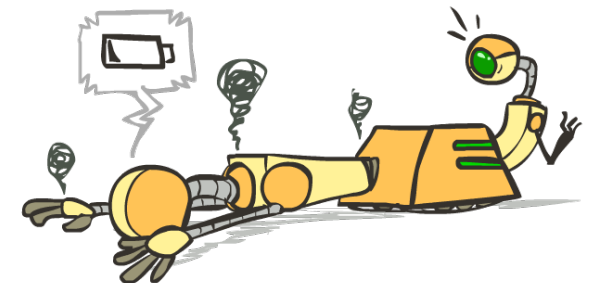
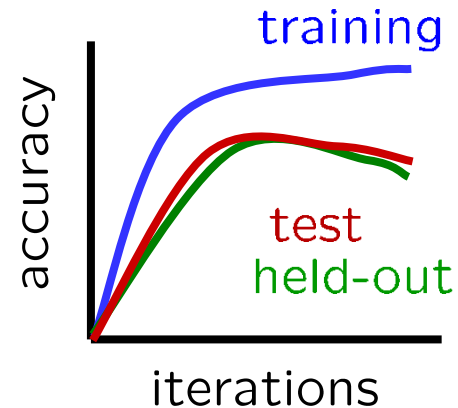
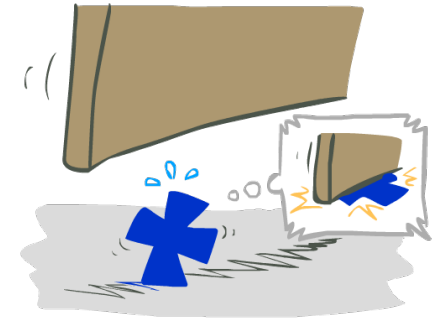
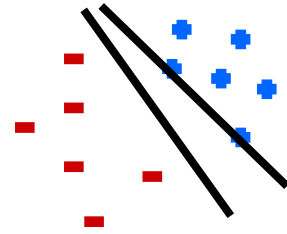
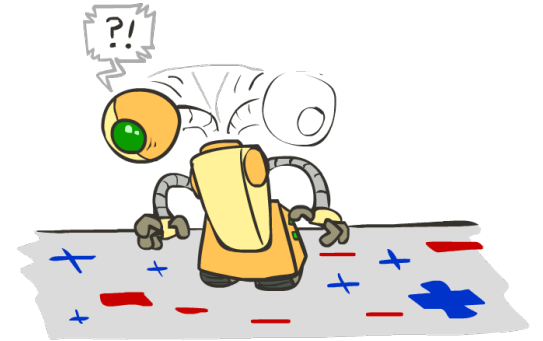
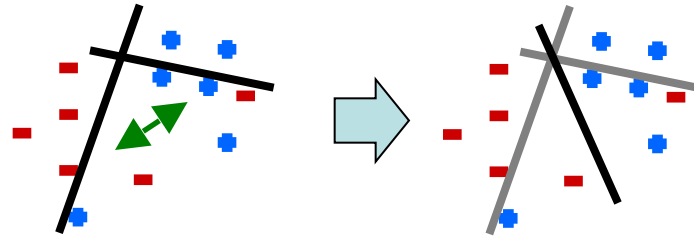


Non-Separable

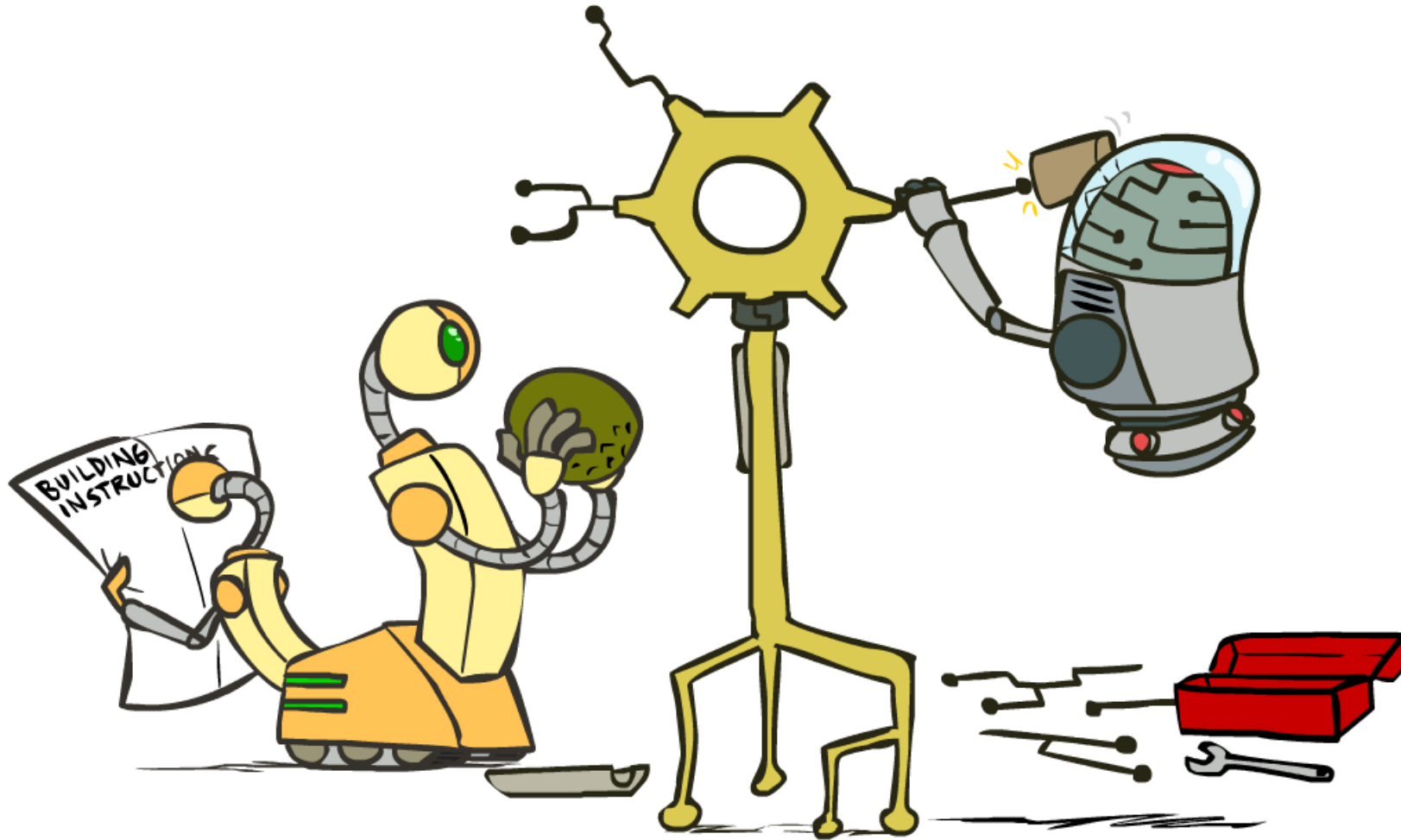


# Problems with the Perceptron

- Noise: if the data isn't separable, weights might thrash
  - 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

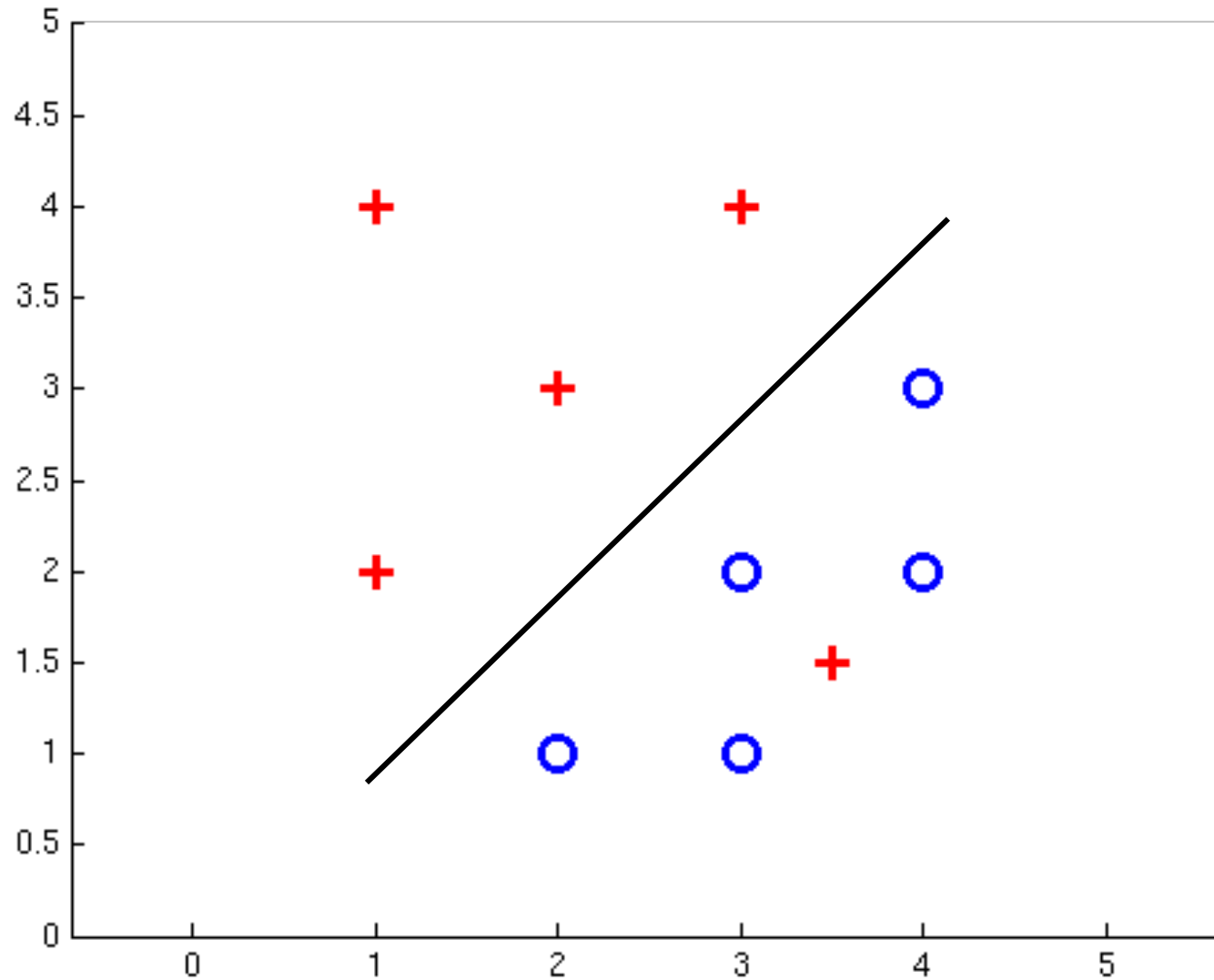


# Improving the Perceptron

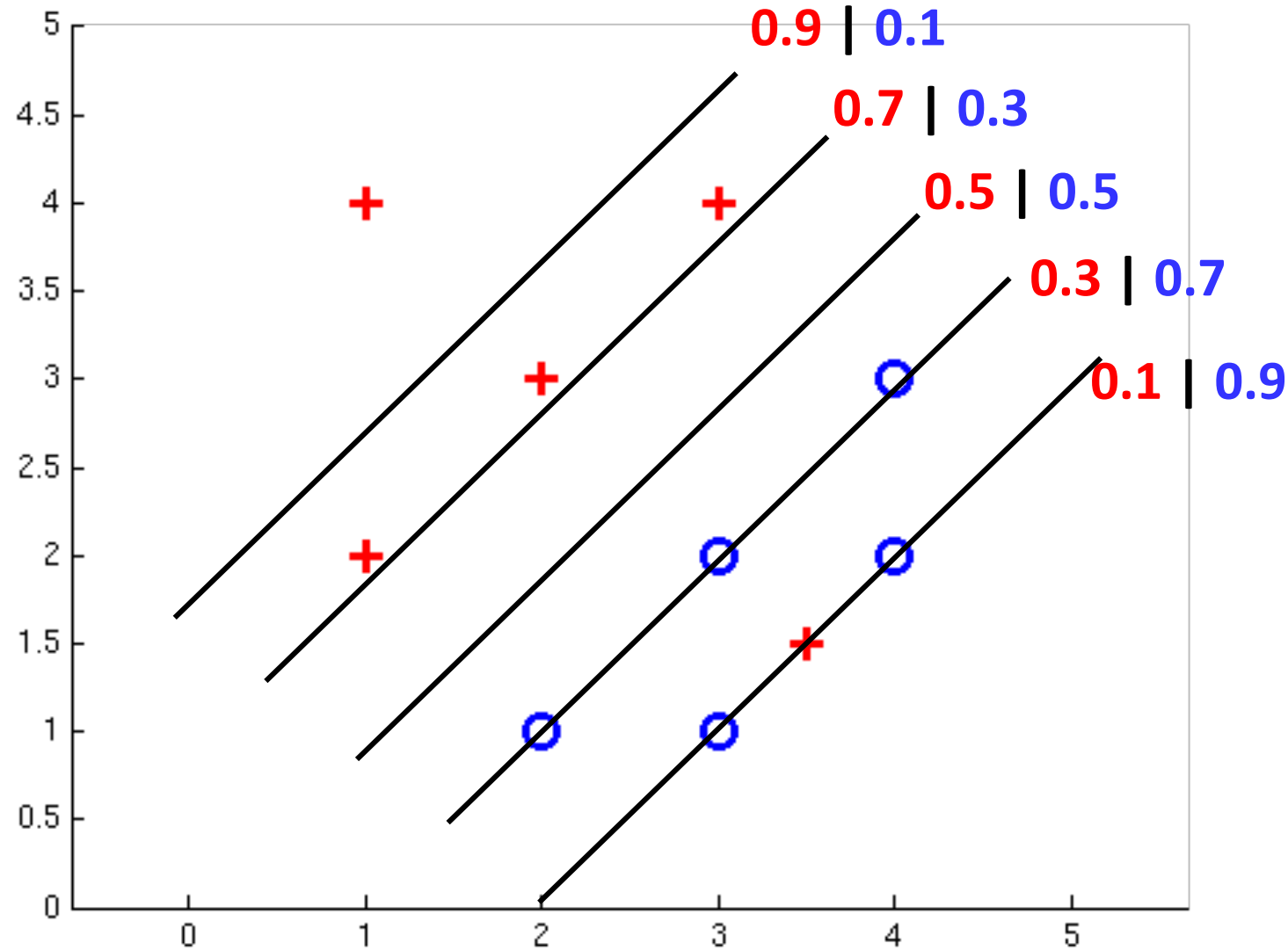


# Non-Separable Case: Deterministic Decision

Even the best linear boundary makes at least one mistake



# Non-Separable Case: Probabilistic Decision

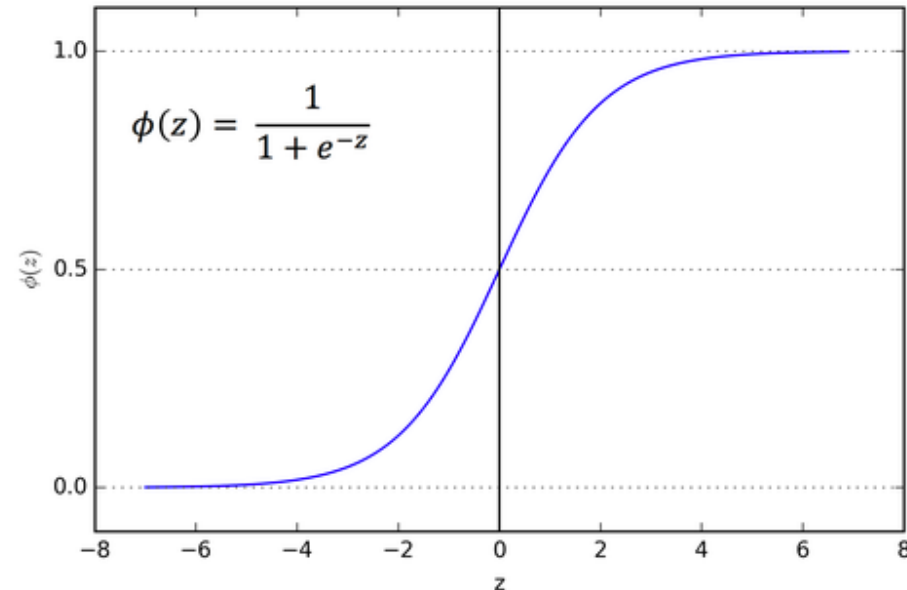


# How to get probabilistic decisions?

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

- Sigmoid function

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



# Best $w$ ?

- Maximum likelihood estimation:

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

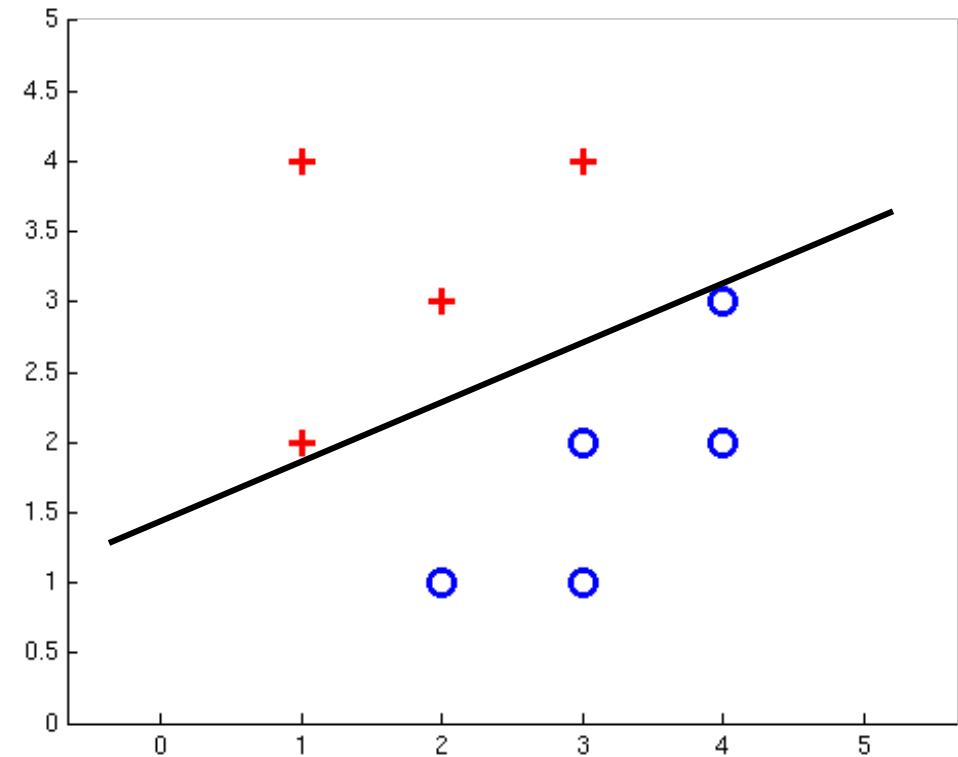
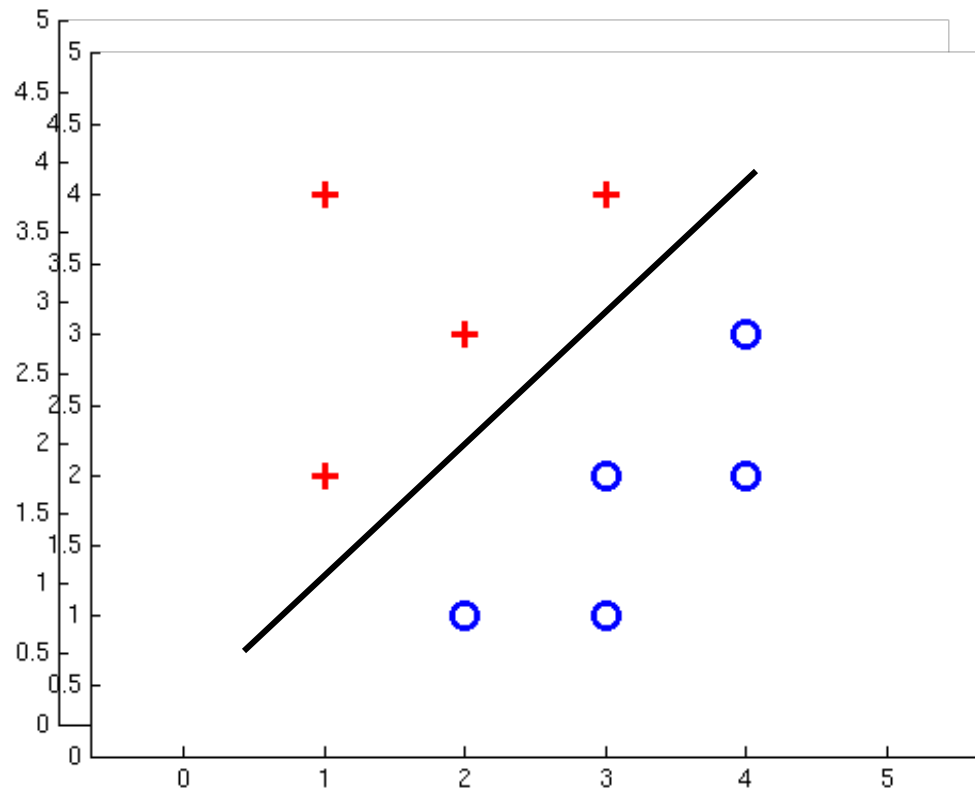
with:

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

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

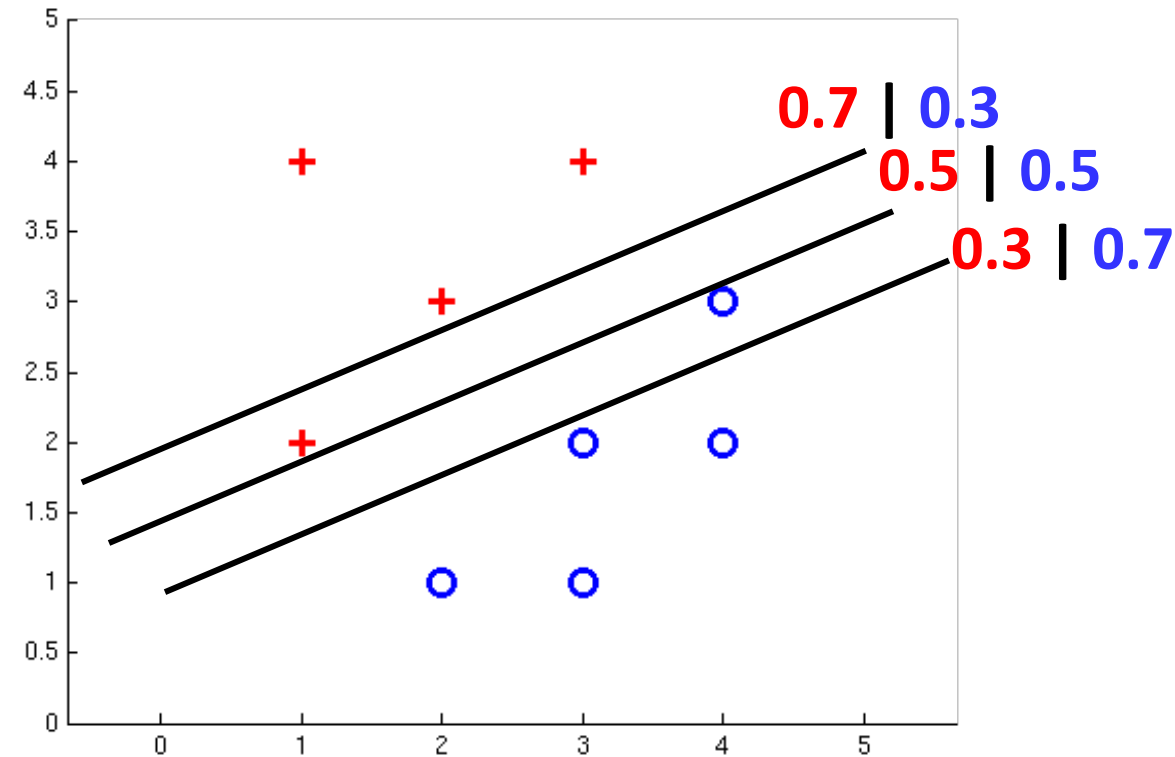
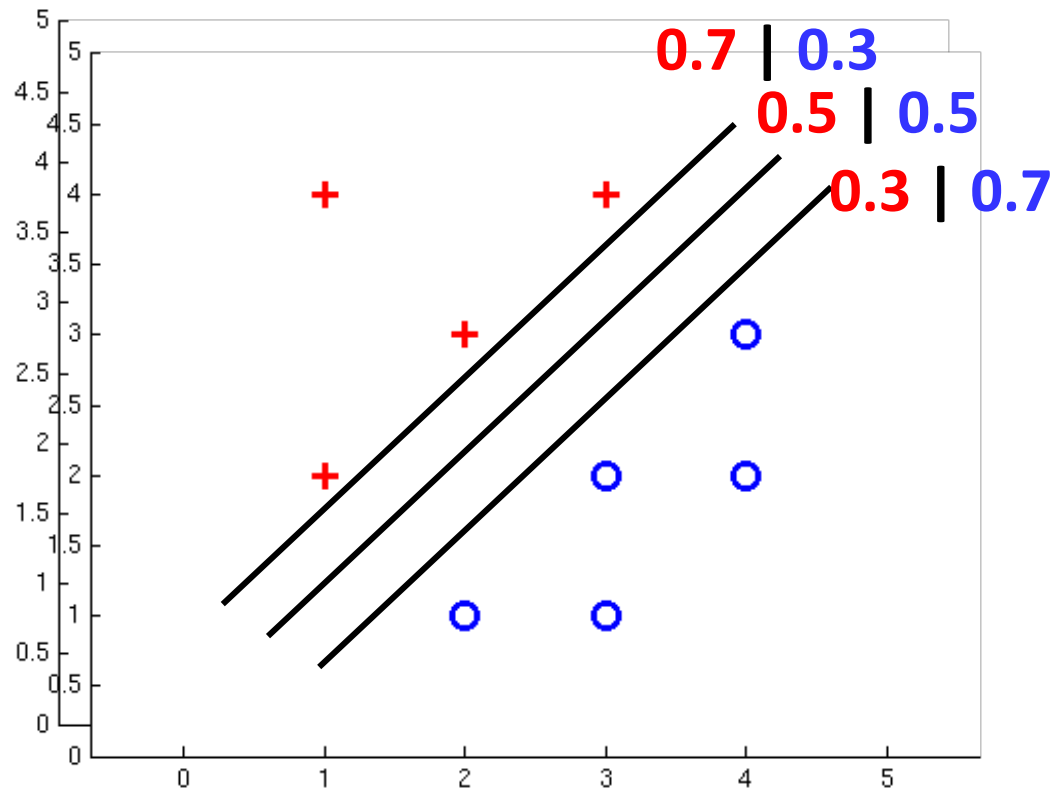
**= Logistic Regression**

# Separable Case: Deterministic Decision – Many Options





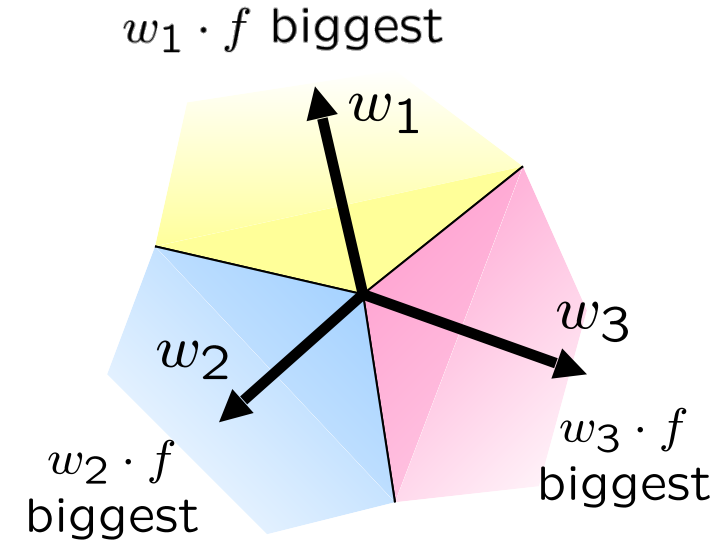
# Separable Case: Probabilistic Decision – Clear Preference



# Multiclass Logistic Regression

- Recall Perceptron:

- A weight vector for each class:  $w_y$
- Score (activation) of a class  $y$ :  $w_y \cdot f(x)$
- Prediction highest score wins  $y = \arg \max_y w_y \cdot f(x)$



- How to make the scores into probabilities?

$$\underbrace{z_1, z_2, z_3}_{\text{original activations}} \rightarrow \underbrace{\frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}}_{\text{softmax activations}}$$

# Best $w$ ?

- Maximum likelihood estimation:

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

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**

# Choosing weights

---

- Optimization

- i.e., how do we solve:

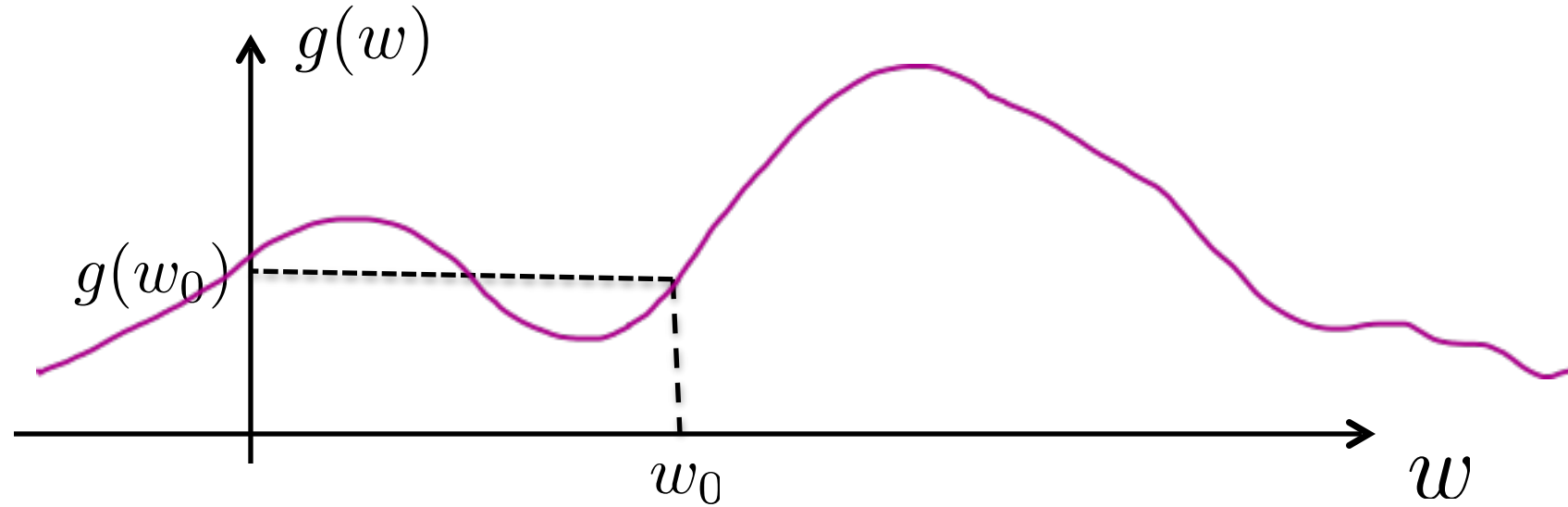
$$\max_w ll(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

# Hill Climbing

- Recall from CSPs lecture: simple, general idea
  - Start wherever
  - Repeat: move to the best neighboring state
  - If no neighbors better than current, quit
- What's particularly tricky when hill-climbing for multiclass logistic regression?
  - Optimization over a continuous space
    - Infinitely many neighbors!
    - How to do this efficiently?



# 1-D Optimization



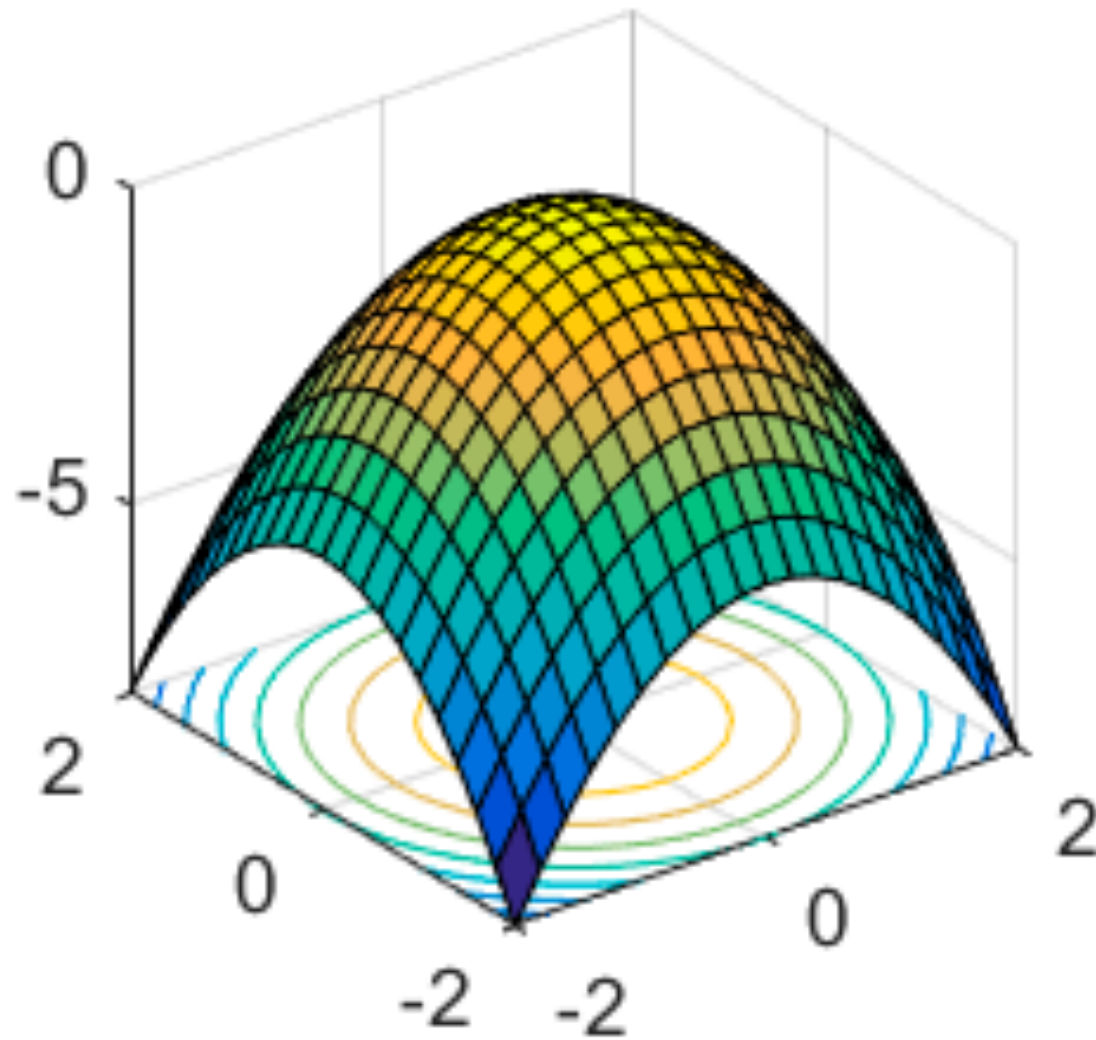
- Could evaluate  $g(w_0 + h)$  and  $g(w_0 - h)$

- Then step in best direction

- Or, evaluate derivative: 
$$\frac{\partial g(w_0)}{\partial w} = \lim_{h \rightarrow 0} \frac{g(w_0 + h) - g(w_0 - h)}{2h}$$

- Tells which direction to step into

# 2-D Optimization



# Gradient Ascent

- Perform update in uphill direction for each coordinate
- The steeper the slope (i.e. the higher the derivative) the bigger the step for that coordinate
- E.g., consider:  $g(w_1, w_2)$

- Updates:

$$w_1 \leftarrow w_1 + \alpha * \frac{\partial g}{\partial w_1}(w_1, w_2)$$

$$w_2 \leftarrow w_2 + \alpha * \frac{\partial g}{\partial w_2}(w_1, w_2)$$

- Updates in vector notation:

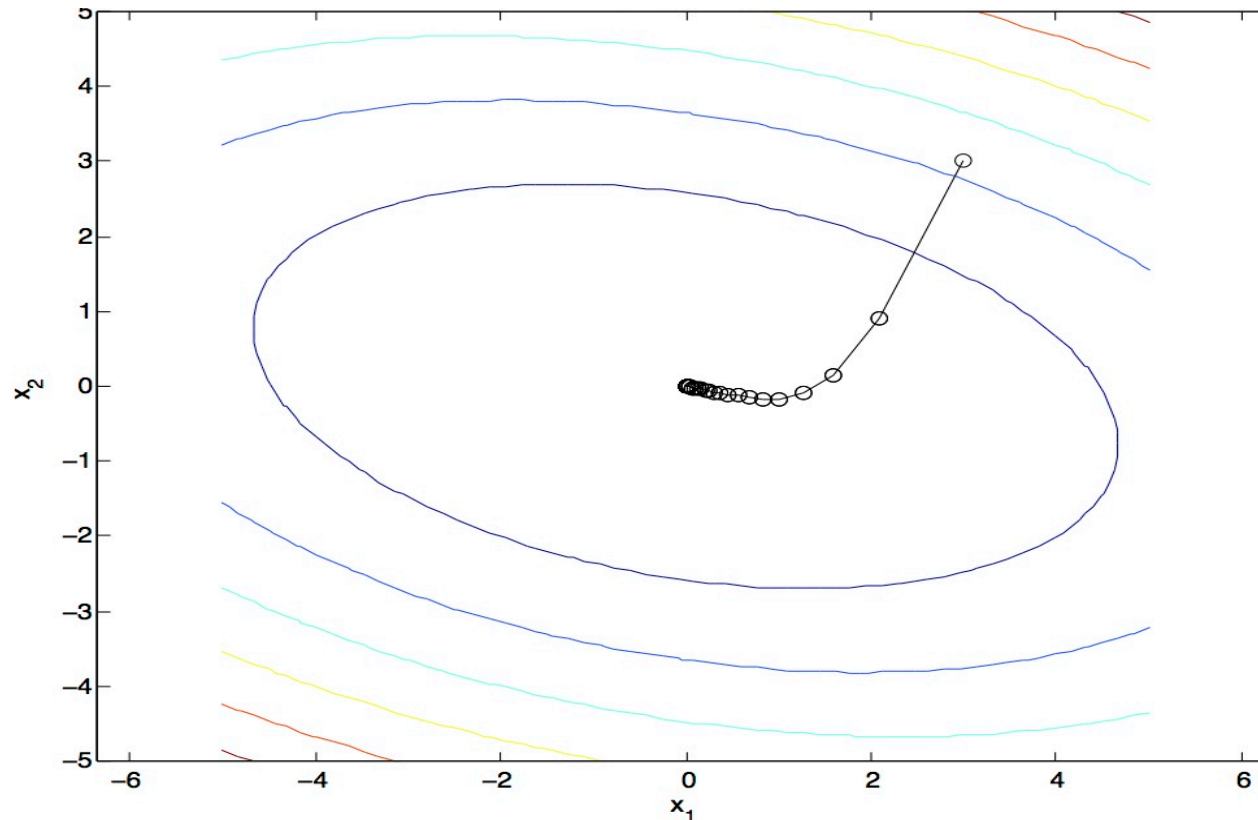
$$w \leftarrow w + \alpha * \nabla_w g(w)$$

$$\text{with: } \nabla_w g(w) = \begin{bmatrix} \frac{\partial g}{\partial w_1}(w) \\ \frac{\partial g}{\partial w_2}(w) \end{bmatrix} = \text{gradient}$$



# Gradient Ascent

- Idea:
  - Start somewhere
  - Repeat: Take a step in the gradient direction



# What is the Steepest Direction?



$$\max_{\Delta: \Delta_1^2 + \Delta_2^2 \leq \varepsilon} g(w + \Delta)$$

- First-Order Taylor Expansion:

$$g(w + \Delta) \approx g(w) + \frac{\partial g}{\partial w_1} \Delta_1 + \frac{\partial g}{\partial w_2} \Delta_2$$

- Steepest Ascent Direction:

$$\max_{\Delta: \Delta_1^2 + \Delta_2^2 \leq \varepsilon} g(w) + \frac{\partial g}{\partial w_1} \Delta_1 + \frac{\partial g}{\partial w_2} \Delta_2$$

- Recall:  $\max_{\Delta: \|\Delta\| \leq \varepsilon} \Delta^\top a \rightarrow$

$$\Delta = \varepsilon \frac{a}{\|a\|}$$

- Hence, solution:  $\Delta = \varepsilon \frac{\nabla g}{\|\nabla g\|}$

**Gradient direction = steepest direction!**

$$\nabla g = \begin{bmatrix} \frac{\partial g}{\partial w_1} \\ \frac{\partial g}{\partial w_2} \end{bmatrix}$$

# Gradient in n dimensions

---

$$\nabla g = \begin{bmatrix} \frac{\partial g}{\partial w_1} \\ \frac{\partial g}{\partial w_2} \\ \dots \\ \frac{\partial g}{\partial w_n} \end{bmatrix}$$

# Optimization Procedure: Gradient Ascent

---

```
■ init  $w$   
■ for iter = 1, 2, ...  
 $w \leftarrow w + \alpha * \nabla g(w)$ 
```

- $\alpha$ : learning rate --- hyperparameter that needs to be chosen carefully
- How? Try multiple choices
  - Crude rule of thumb: update changes  $w$  about 0.1 – 1 %

# Batch Gradient Ascent on the Log Likelihood Objective

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

- `init`  $w$
- `for`  $iter = 1, 2, \dots$

$$w \leftarrow w + \alpha * \sum_i \nabla \log P(y^{(i)} | x^{(i)}; w)$$

# Stochastic Gradient Ascent on the Log Likelihood Objective

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

**Observation:** once gradient on one training example has been computed, might as well incorporate before computing next one

- `init`  $w$
- `for`  $iter = 1, 2, \dots$ 
  - pick random  $j$

$$w \leftarrow w + \alpha * \nabla \log P(y^{(j)} | x^{(j)}; w)$$

# Mini-Batch Gradient Ascent on the Log Likelihood Objective

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

**Observation:** gradient over small set of training examples (=mini-batch) can be computed in parallel, might as well do that instead of a single one

- `init`  $w$
- `for`  $iter = 1, 2, \dots$ 
  - pick random subset of training examples  $J$

$$w \leftarrow w + \alpha * \sum_{j \in J} \nabla \log P(y^{(j)} | x^{(j)}; w)$$

# Learning Rate Finder

- Calculate a good learning rate by trying learning rates over a range of possible values
- Plot the training loss at each of these epochs
- Pick a learning rate where the loss is declining the most before it hits the minimum:  
 $5 \times 10^{-5}$  -  $3 \times 10^{-4}$

