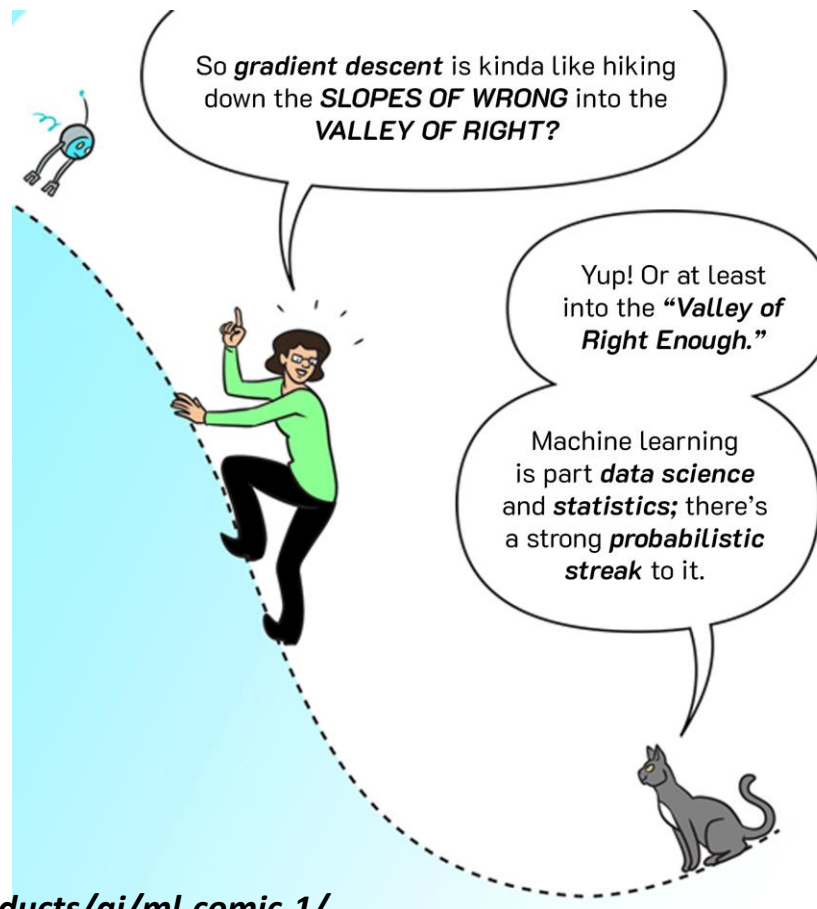


# Machine Learning

## L4\_Perceptron Learning



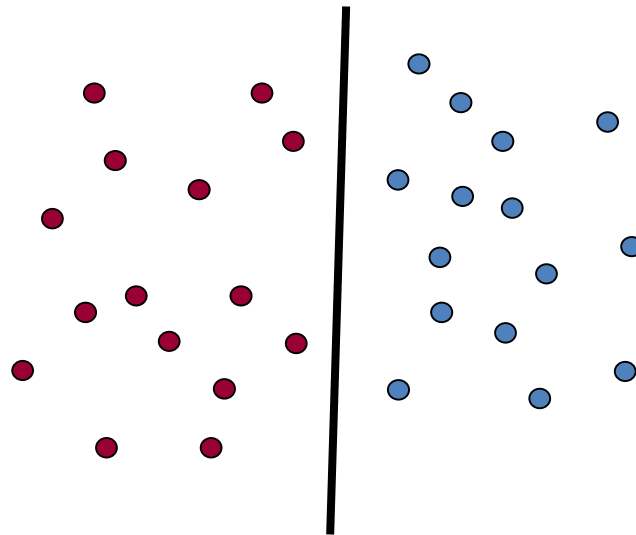
# Today's outline

- Linear models
- Perceptron

# Linear models

A strong high-bias assumption is *linear separability*:  
in 2 dimensions, can separate classes by a line  
in higher dimensions, need hyperplanes

A *linear model* is a model that assumes the data is linearly separable



# Hyperplanes

A hyperplane is line/plane in a high-dimensional space



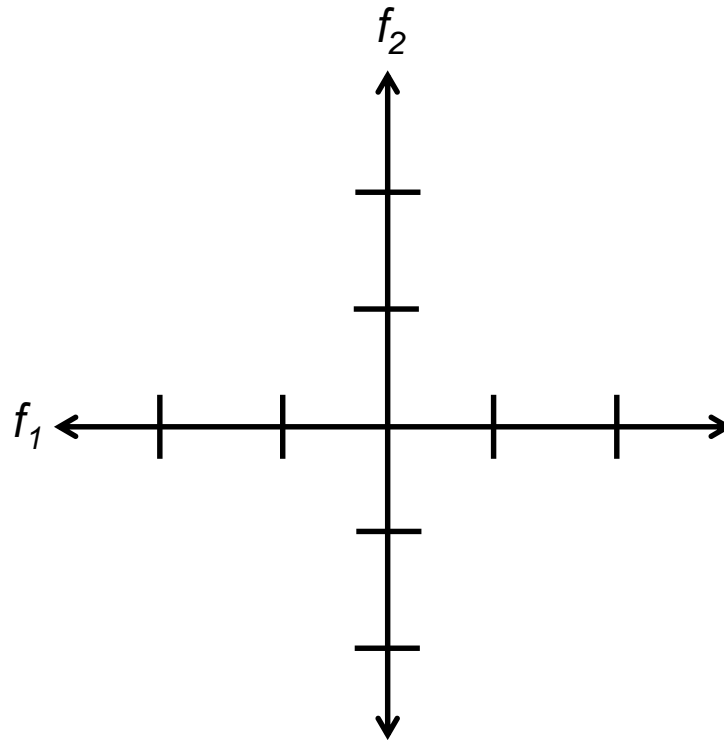
What defines a line?

What defines a hyperplane?

# Defining a line

Any pair of values  $(w_1, w_2)$  defines a line through the origin:

$$0 = w_1 f_1 + w_2 f_2$$



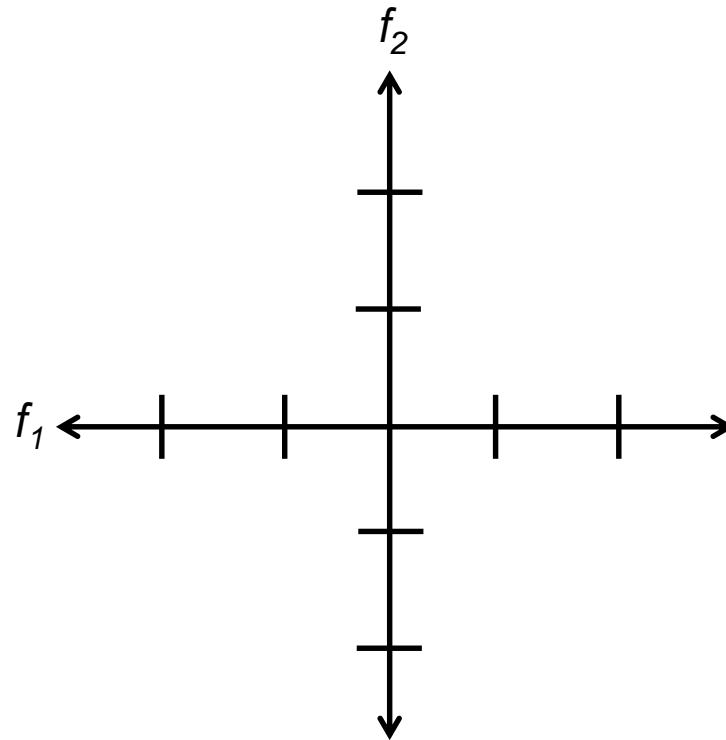
# Defining a line

Any pair of values  $(w_1, w_2)$  defines a line through the origin:

$$0 = w_1 f_1 + w_2 f_2$$

$$0 = 1f_1 + 2f_2$$

<b>-2</b>	<b>1</b>
<b>-1</b>	<b>0.5</b>
<b>0</b>	<b>0</b>
<b>1</b>	<b>-0.5</b>
<b>2</b>	<b>-1</b>



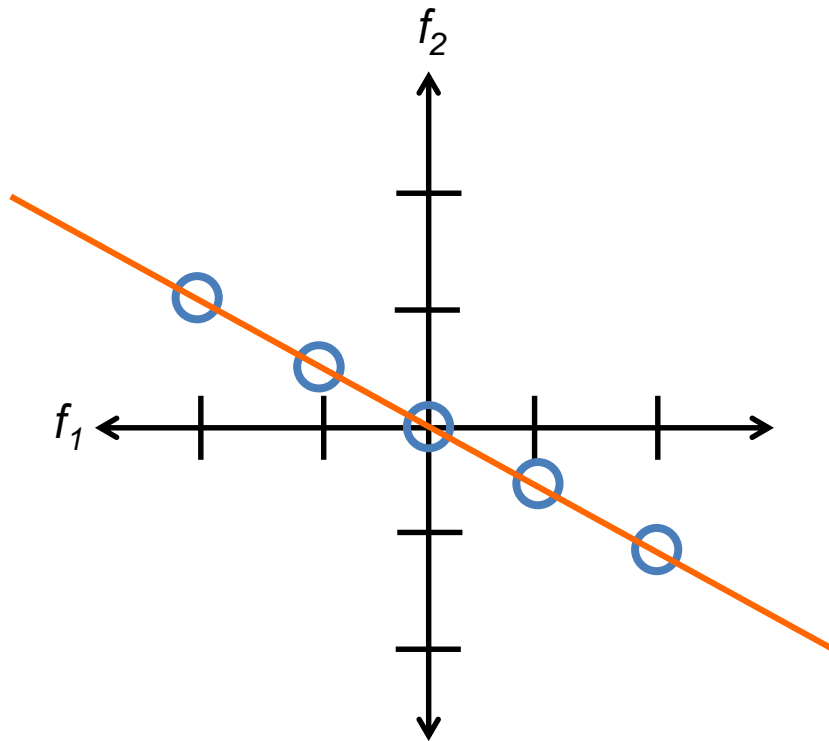
# Defining a line

Any pair of values  $(w_1, w_2)$  defines a line through the origin:

$$0 = w_1 f_1 + w_2 f_2$$

$$0 = 1f_1 + 2f_2$$

<b>-2</b>	<b>1</b>
<b>-1</b>	<b>0.5</b>
<b>0</b>	<b>0</b>
<b>1</b>	<b>-0.5</b>
<b>2</b>	<b>-1</b>



# Defining a line

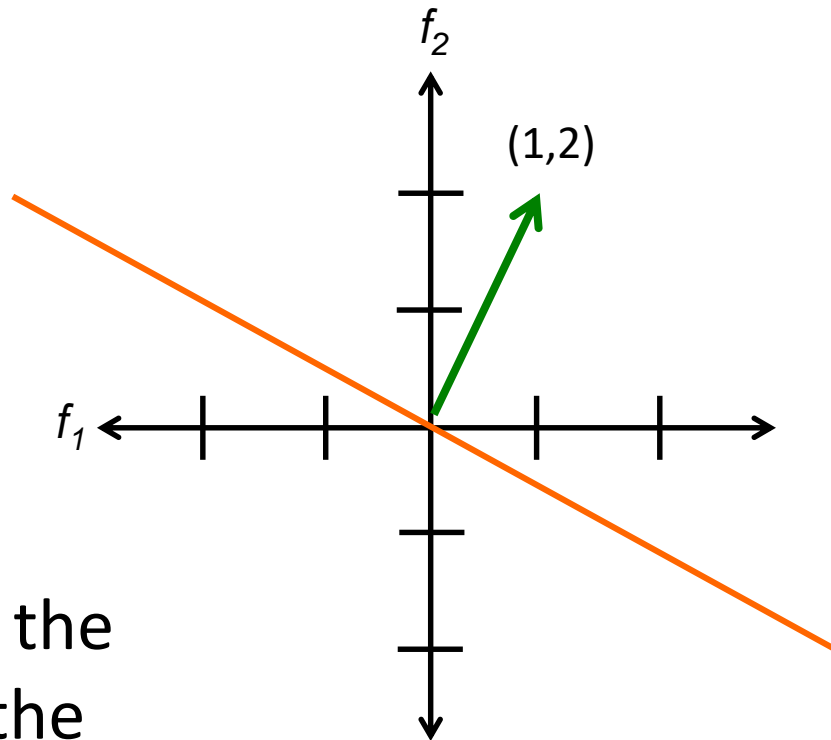
Any pair of values  $(w_1, w_2)$  defines a line through the origin:

$$0 = w_1 f_1 + w_2 f_2$$

$$0 = 1f_1 + 2f_2$$

$$w=(1,2)$$

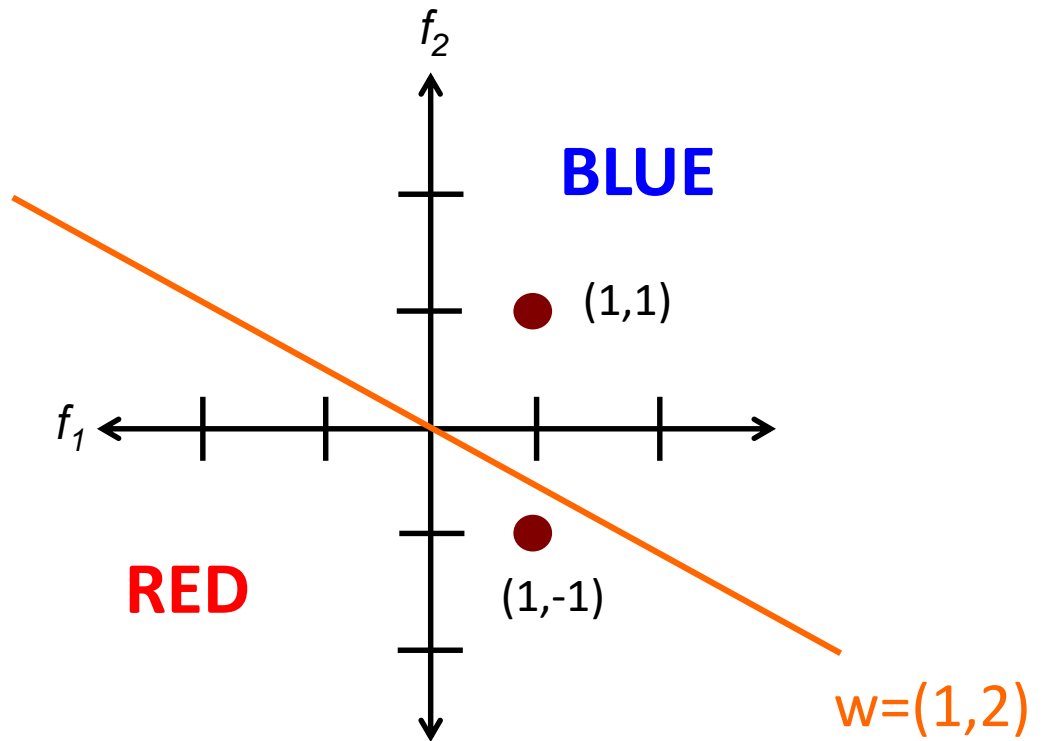
We can also view it as the line perpendicular to the *weight vector*



# Classifying with a line

Mathematically, how can we classify points based on a line?

$$0 = 1f_1 + 2f_2$$



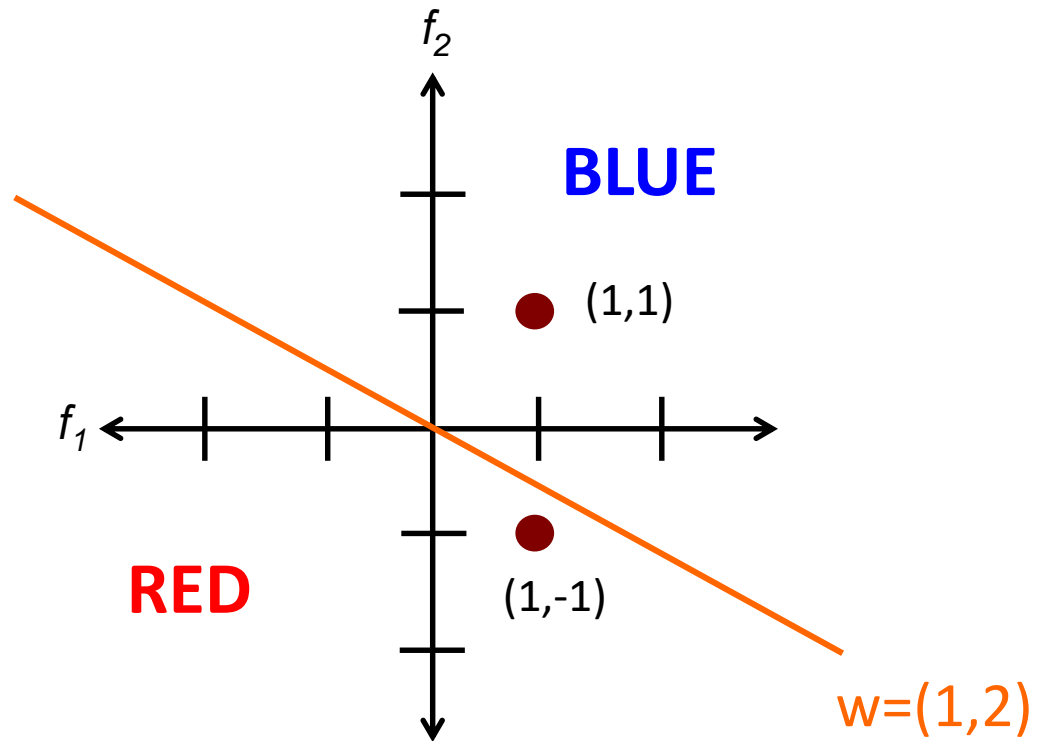
# Classifying with a line

Mathematically, how can we classify points based on a line?

$$0 = 1f_1 + 2f_2$$

$$(1,1): 1 * 1 + 2 * 1 = 3$$

$$(1,-1): 1 * 1 + 2 * -1 = -1$$



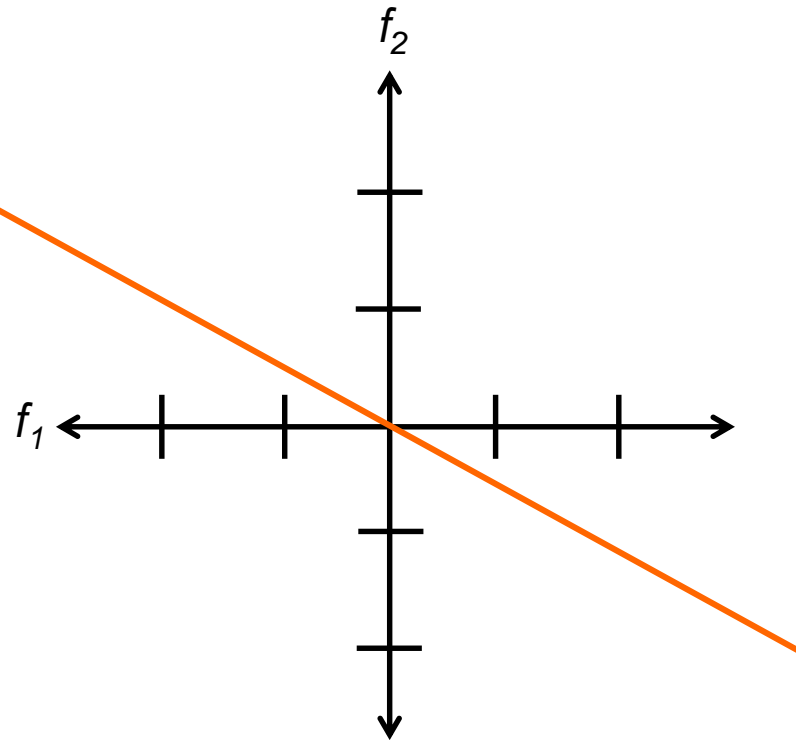
The sign indicates which side of the line

# Defining a line

Any pair of values  $(w_1, w_2)$  defines a line through the origin:

$$0 = w_1 f_1 + w_2 f_2$$

$$0 = 1f_1 + 2f_2$$



How do we move the line off of the origin?

# Defining a line

Any pair of values  $(w_1, w_2)$  defines a line through the origin:

$$a = w_1 f_1 + w_2 f_2$$

$$-1 = 1f_1 + 2f_2$$

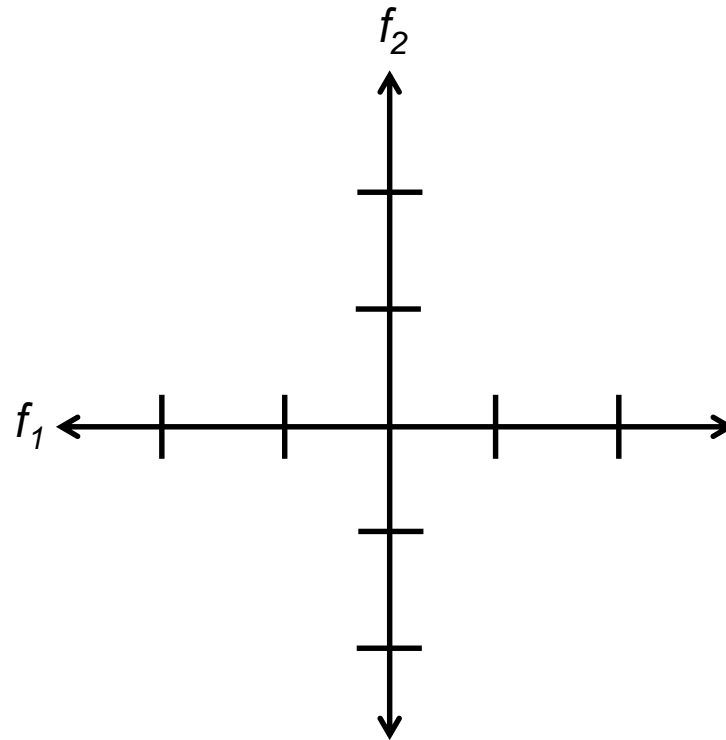
**-2**

**-1**

**0**

**1**

**2**



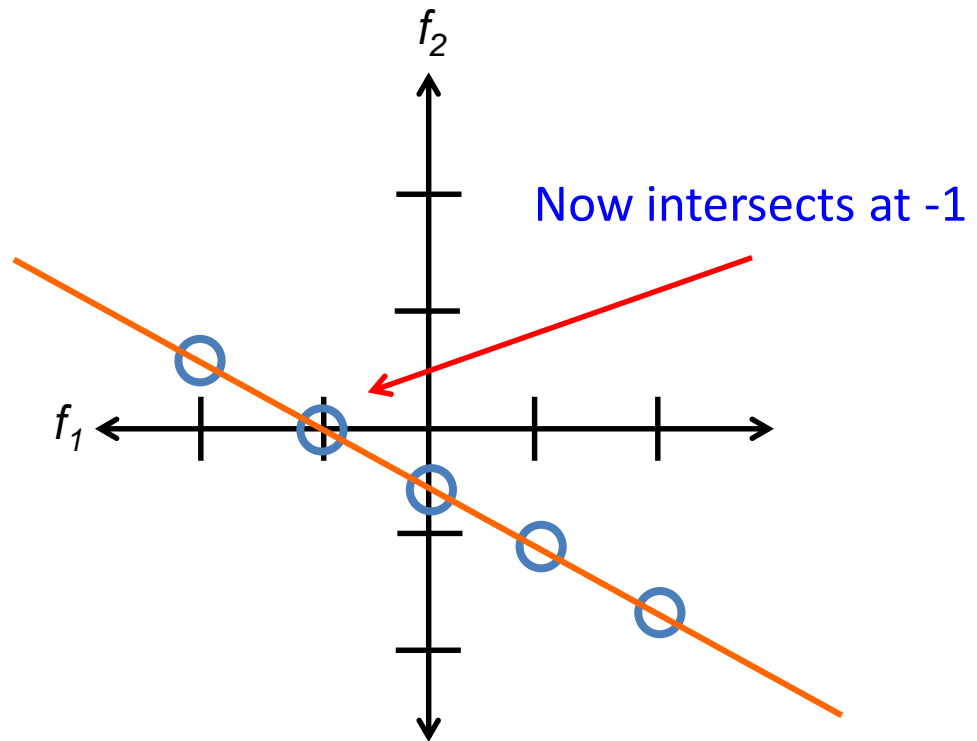
# Defining a line

Any pair of values  $(w_1, w_2)$  defines a line through the origin:

$$a = w_1 f_1 + w_2 f_2$$

$$-1 = 1f_1 + 2f_2$$

<b>-2</b>	<b>0.5</b>
<b>-1</b>	<b>0</b>
<b>0</b>	<b>-0.5</b>
<b>1</b>	<b>-1</b>
<b>2</b>	<b>-1.5</b>



# Linear models

A linear model in  $n$ -dimensional space (i.e.  $n$  features) is defined by  $n+1$  weights:

- In **two dimensions**, a **line**:

$$0 = w_1 f_1 + w_2 f_2 + b \quad (\text{where } b = -a)$$

- In **three dimensions**, a **plane**:

$$0 = w_1 f_1 + w_2 f_2 + w_3 f_3 + b$$

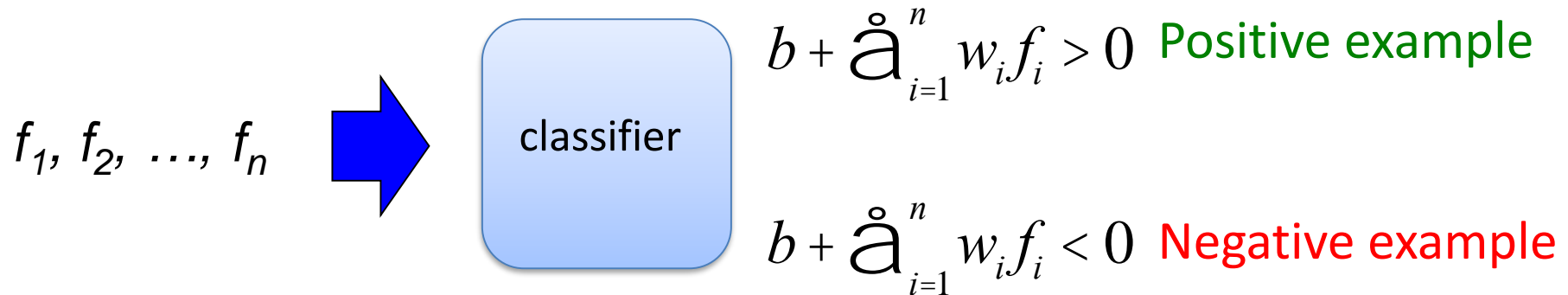
- In  **$n$ -dimensions**, a **hyperplane**

$$0 = b + \mathring{a}_{i=1}^n w_i f_i$$



# Classifying with a linear model

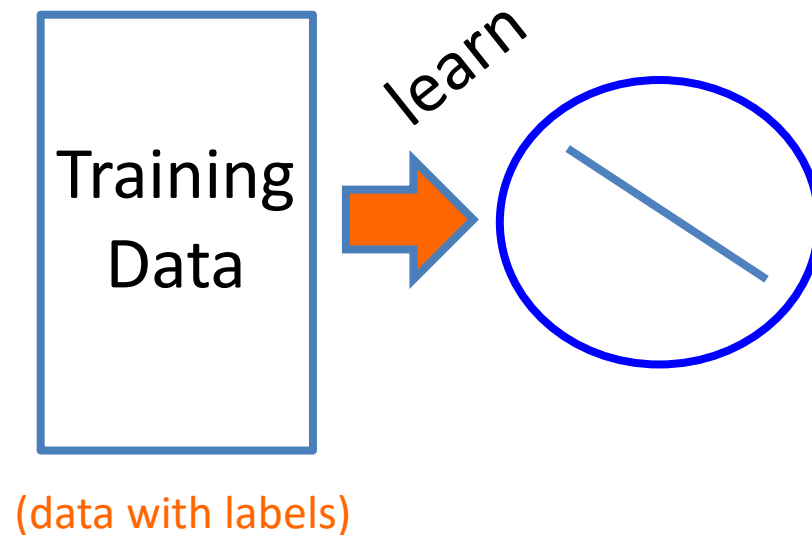
We can classify with a linear model by checking the sign:



# Learning a linear model

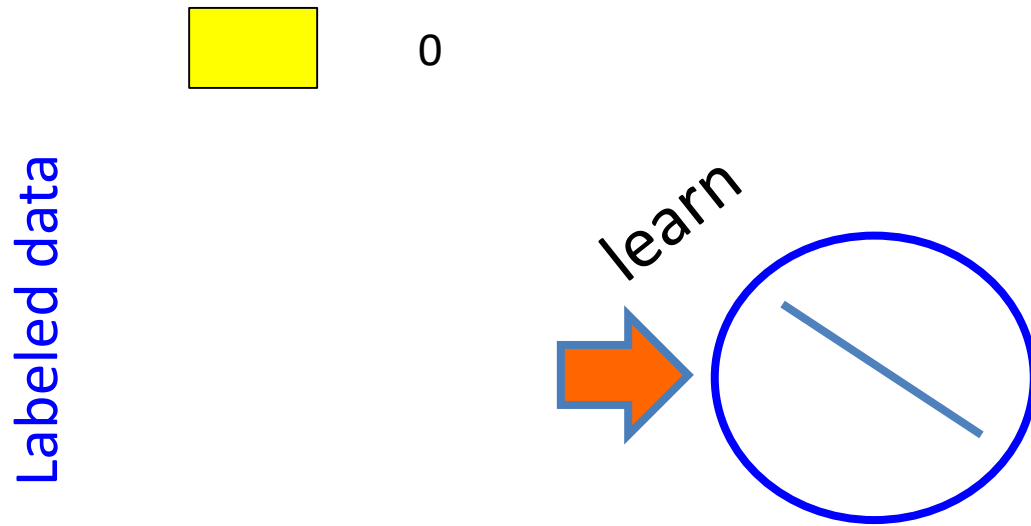
Geometrically, we know what a linear model represents

Given a linear model (i.e. a set of weights and  $b$ ) we can classify examples



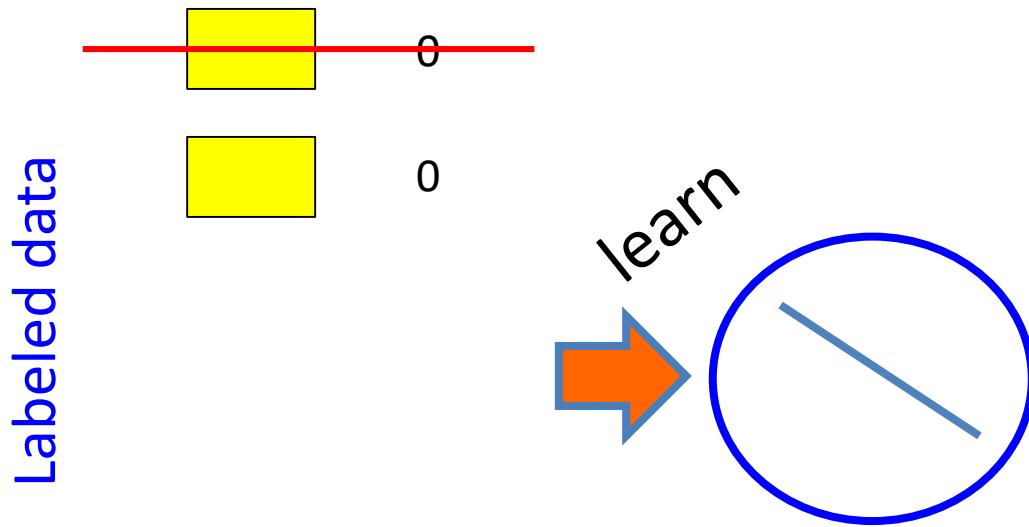
How do we learn a linear model?

# Online learning algorithm



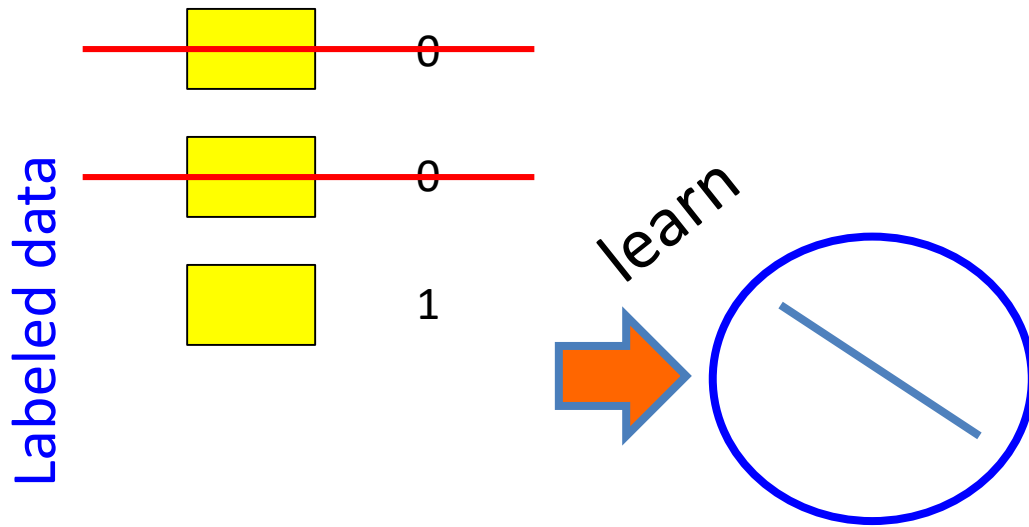
Only get to see one example at a time!

# Online learning algorithm



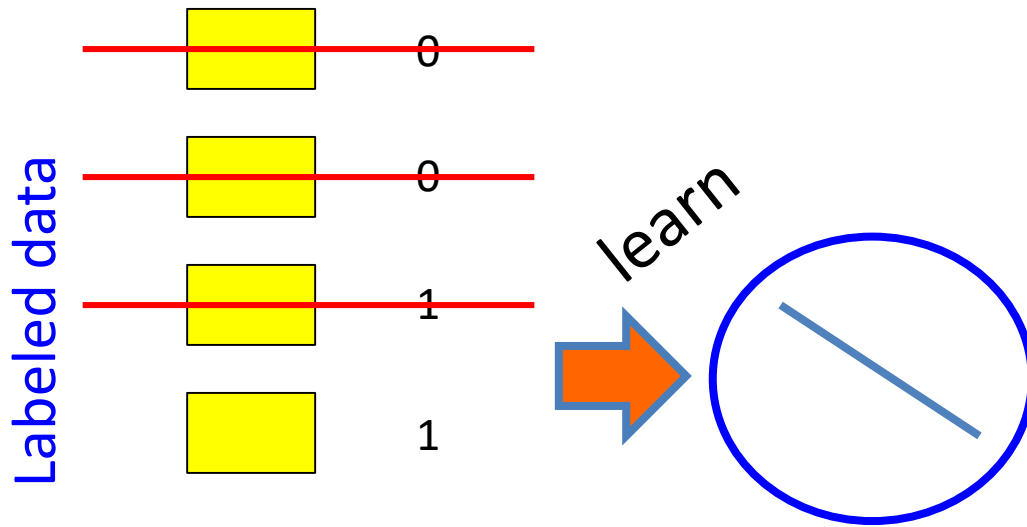
Only get to see one example at a time!

# Online learning algorithm



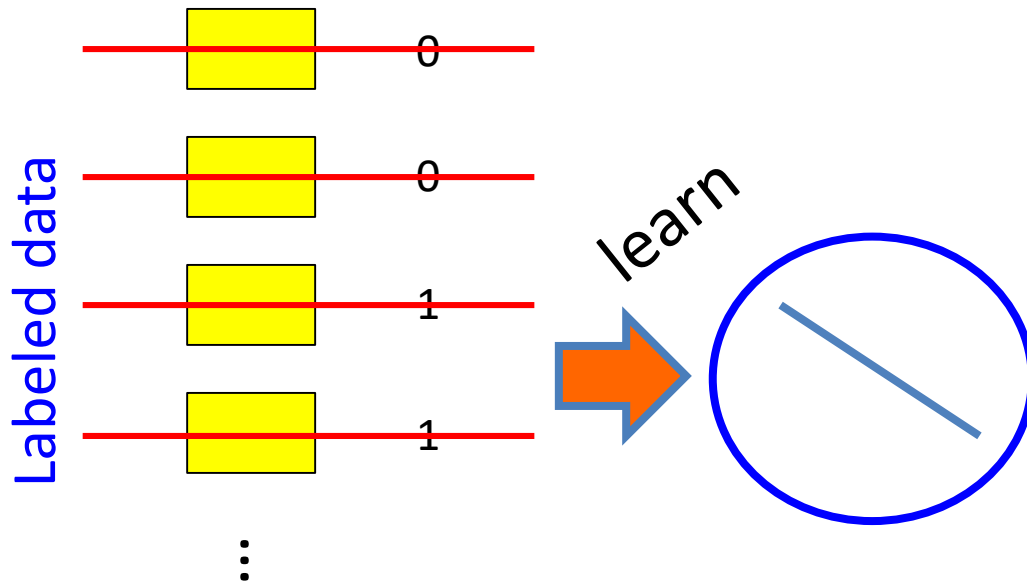
Only get to see one example at a time!

# Online learning algorithm



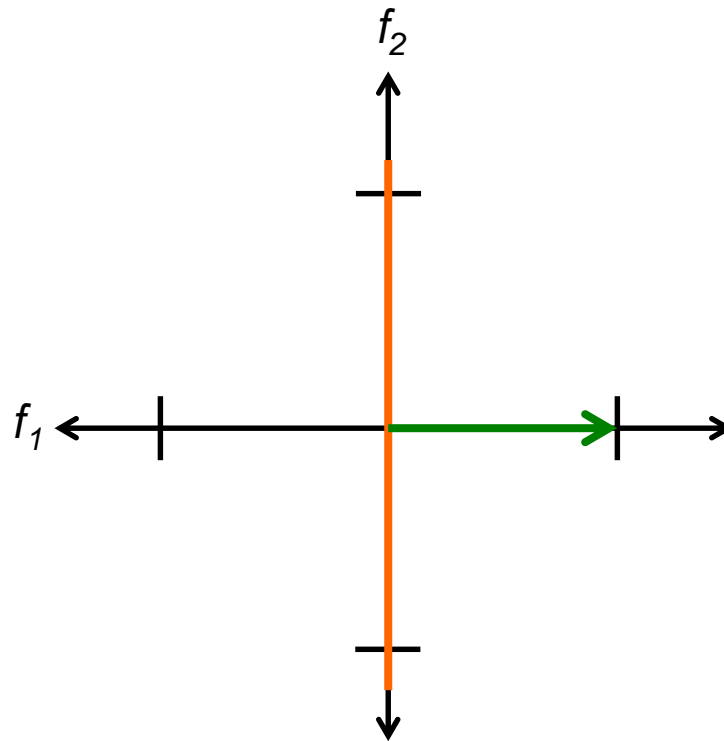
Only get to see one example at a time!

# Online learning algorithm



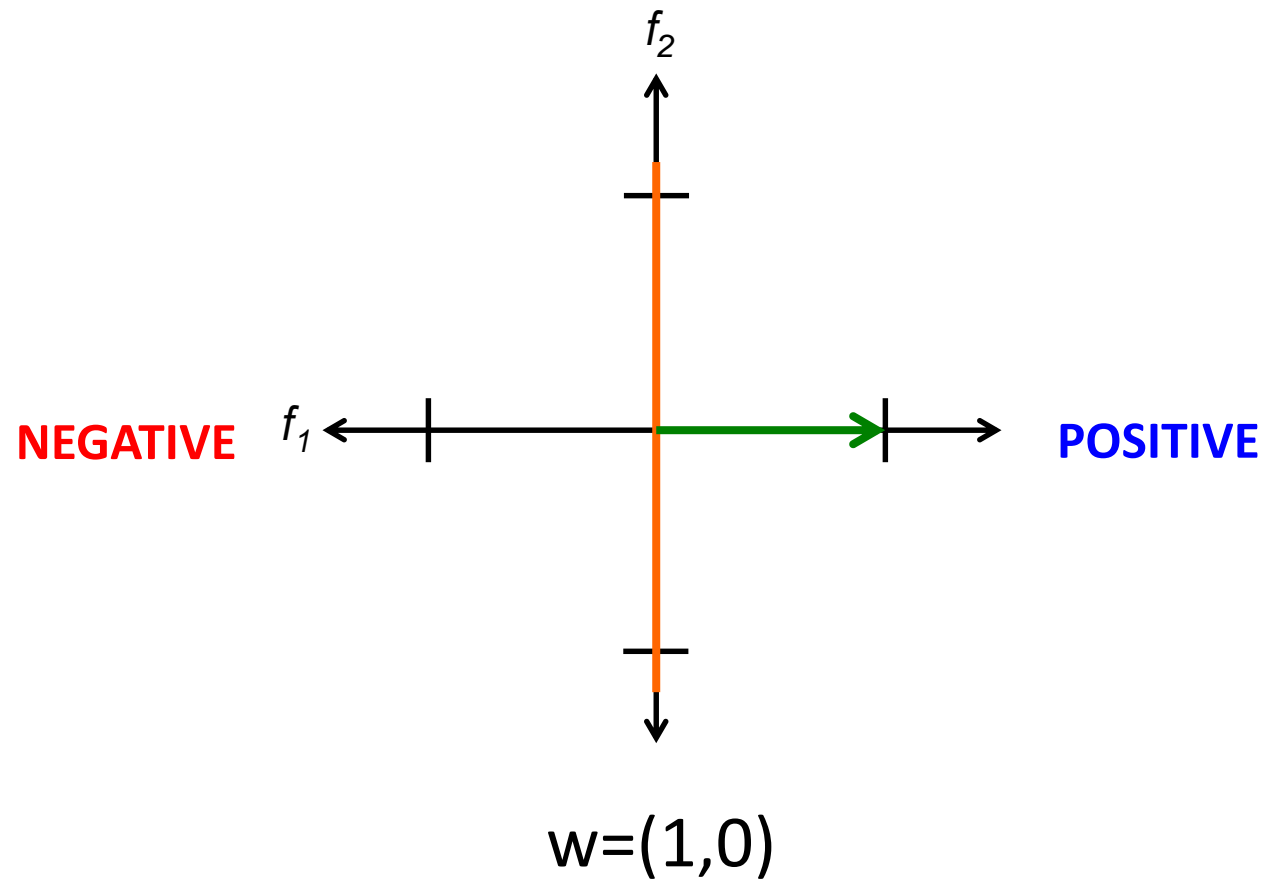
Only get to see one example at a time!

# Learning a linear classifier



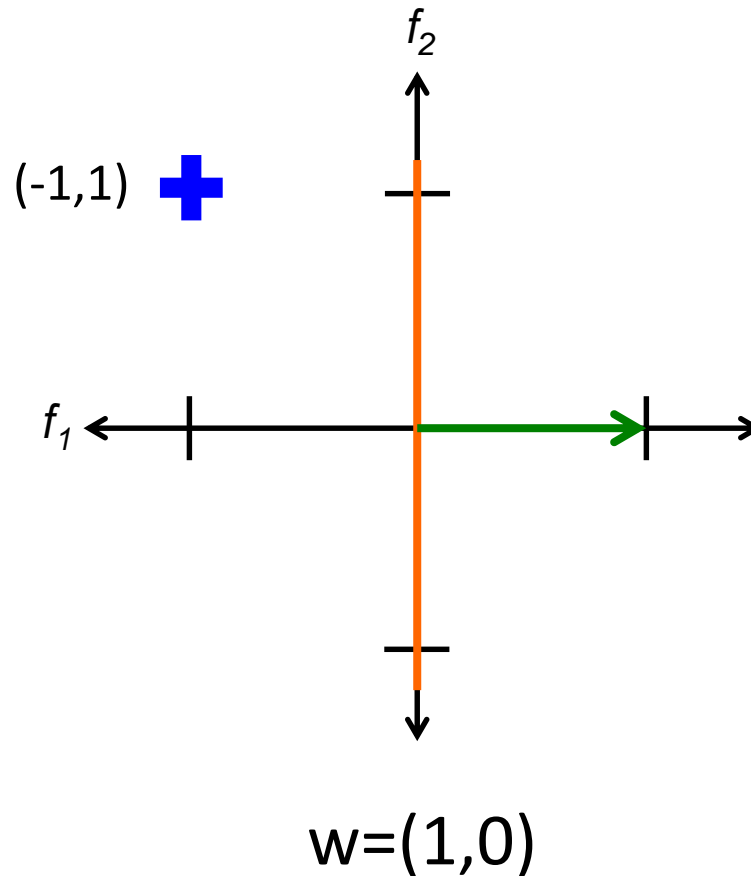
What does this model currently say?  $w=(1,0)$

# Learning a linear classifier



# Learning a linear classifier

$$0 = w_1 f_1 + w_2 f_2$$



Is our current guess:  
right or wrong?

# Learning a linear classifier

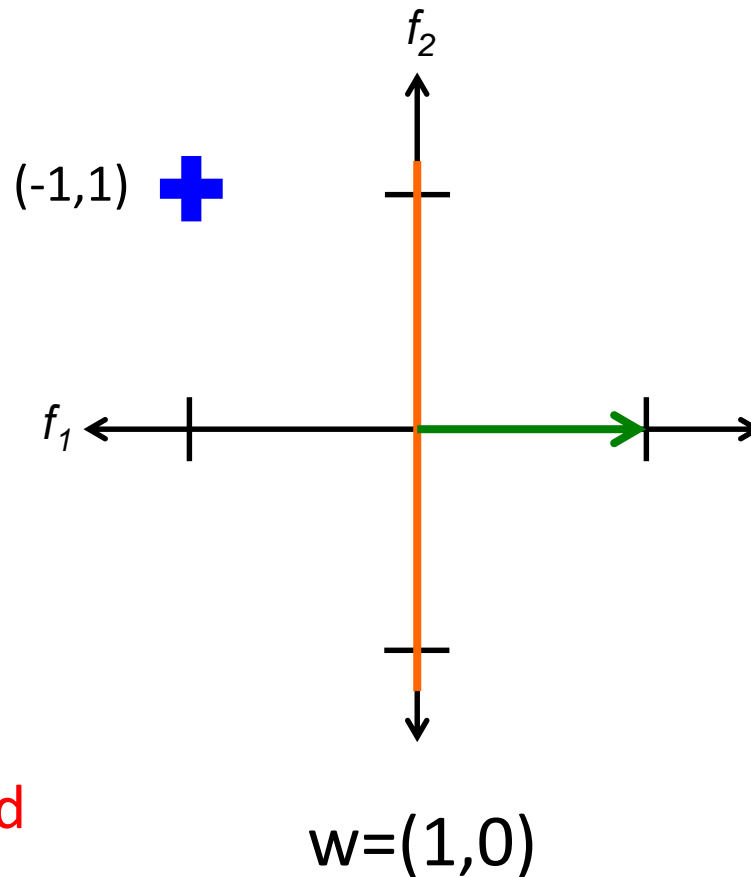
$$0 = w_1 f_1 + w_2 f_2$$

$$1 * f_1 + 0 * f_2 =$$

$$1 * -1 + 0 * 1 = -1$$

predicts negative, wrong

Geometrically, how should we update the model?

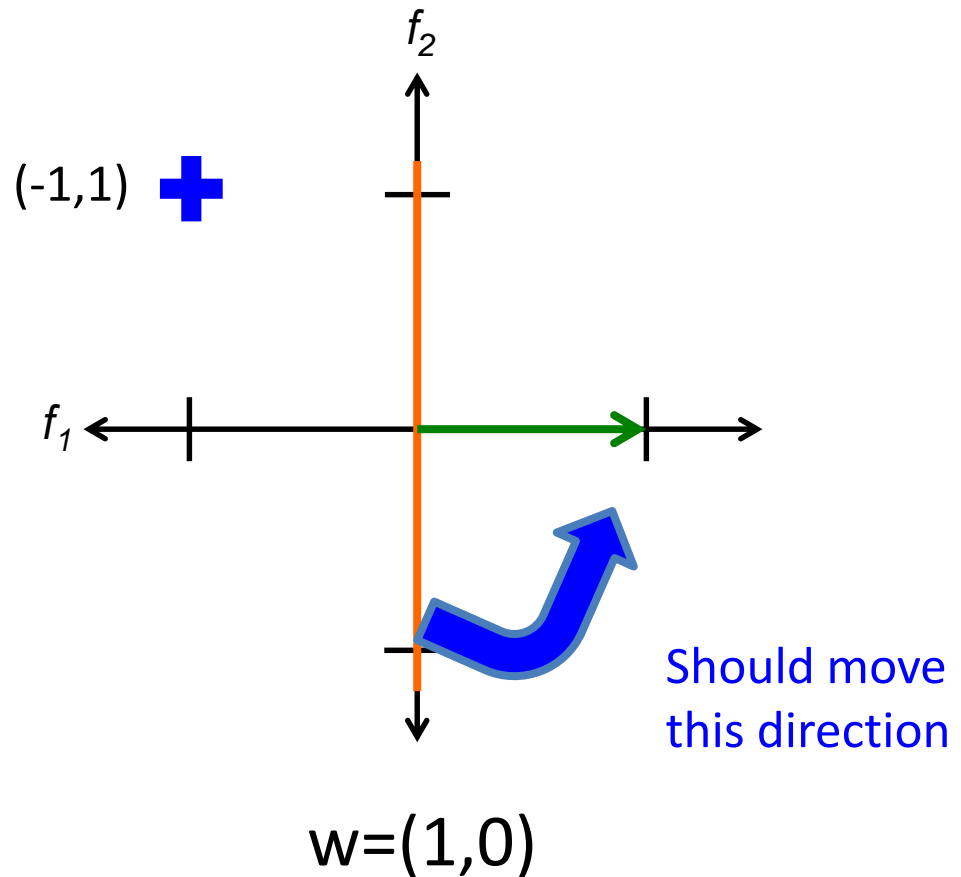


# Learning a linear classifier

$$0 = w_1 f_1 + w_2 f_2$$

$$1 * f_1 + 0 * f_2 =$$

$$1 * -1 + 0 * 1 = -1$$



# A closer look at why we got it wrong

$$w_1 \quad w_2$$

(-1, 1, positive)

$$1 * f_1 + 0 * f_2 =$$

$$1 * -1 + 0 * 1 = -1 \leftarrow$$

We'd like this value to be positive since it's a positive value



Which of these contributed to the mistake?

# A closer look at why we got it wrong

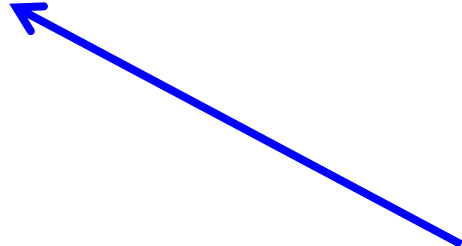
$w_1$        $w_2$

(-1, 1, positive)

$$1 * f_1 + 0 * f_2 =$$

$$1 * -1 + 0 * 1 = -1$$

We'd like this value to be positive since it's a positive value



contributed in the wrong direction

could have contributed (positive feature), but didn't

How should we change the weights?

# A closer look at why we got it wrong

$w_1$        $w_2$

(-1, 1, positive)

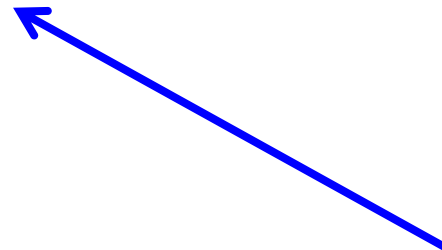
$$1 * f_1 + 0 * f_2 =$$

$$1 * -1 + 0 * 1 = -1$$

We'd like this value to be positive since it's a positive value



contributed in the wrong direction



could have contributed (positive feature), but didn't

decrease

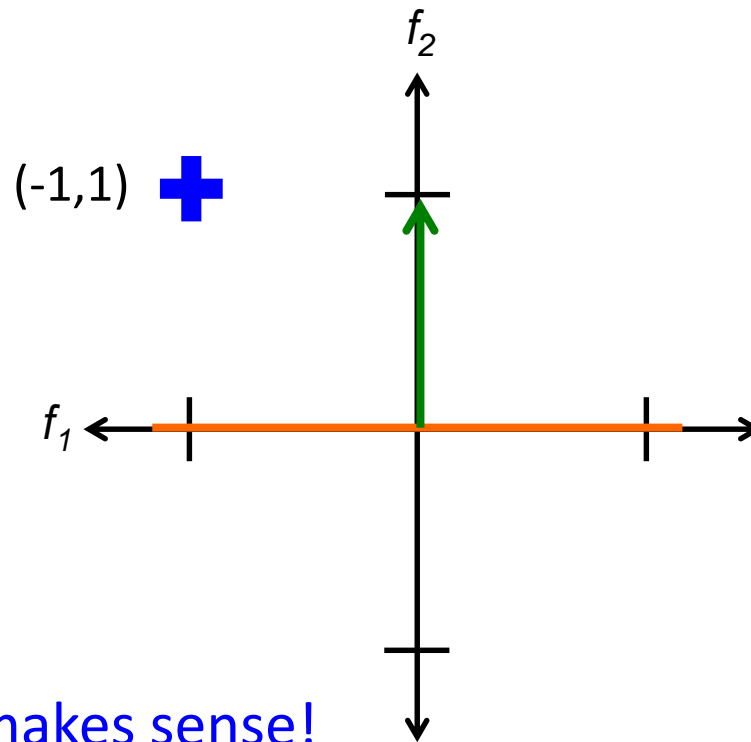
increase

1 -> 0

0 -> 1

# Learning a linear classifier

$$0 = w_1 f_1 + w_2 f_2$$

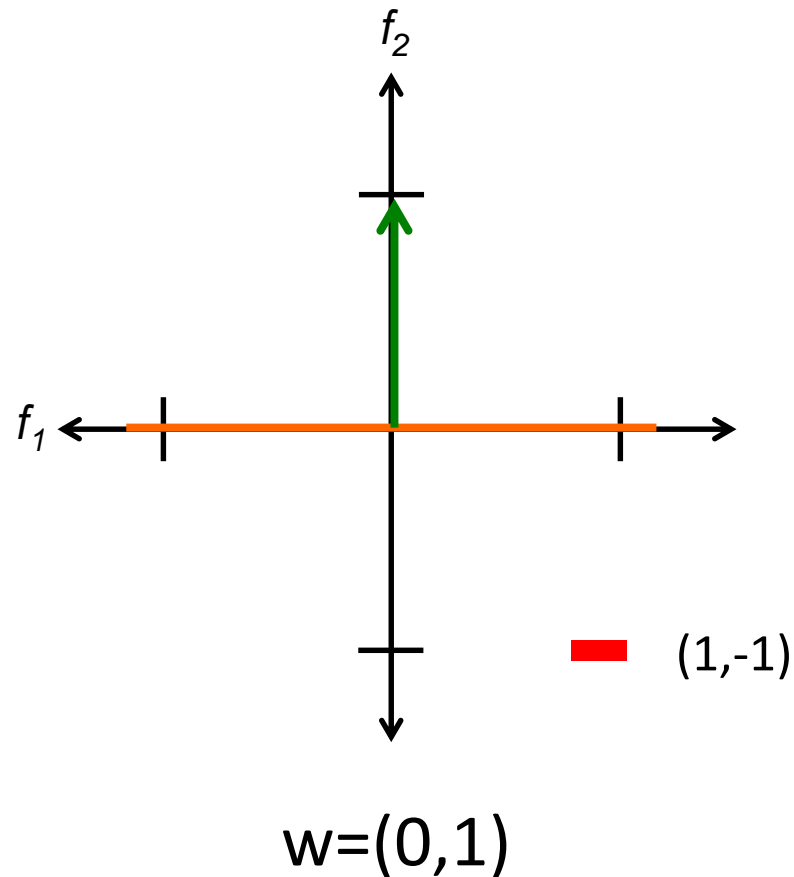


Geometrically, this also makes sense!

$$w = (0, 1)$$

# Learning a linear classifier

$$0 = w_1 f_1 + w_2 f_2$$



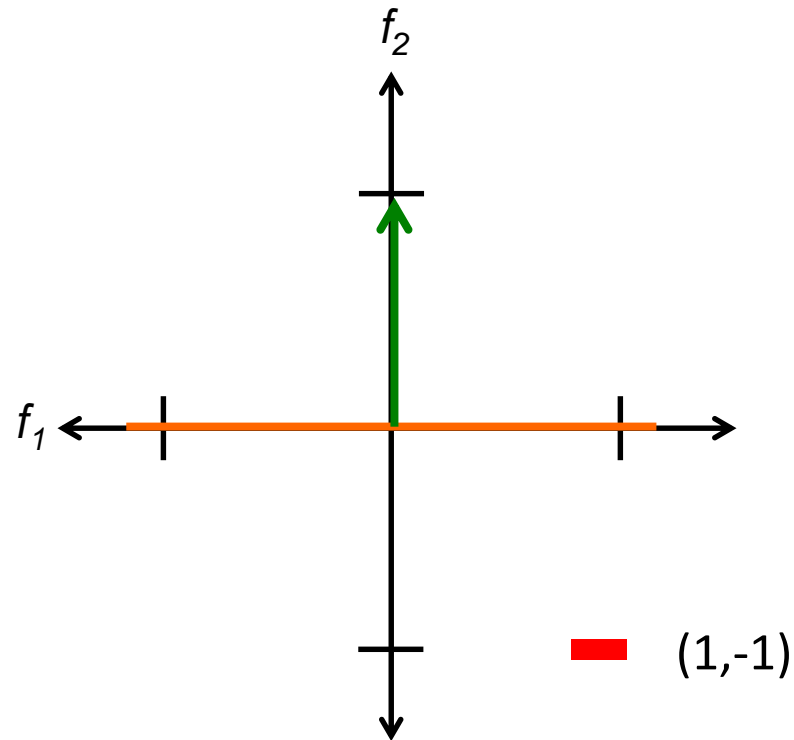
# Learning a linear classifier

$$0 = w_1 f_1 + w_2 f_2$$

$$0 * f_1 + 1 * f_2 =$$

$$0 * 1 + 1 * -1 = -1$$

predicts negative, correct



How should we update the model?

$$w = (0, 1)$$

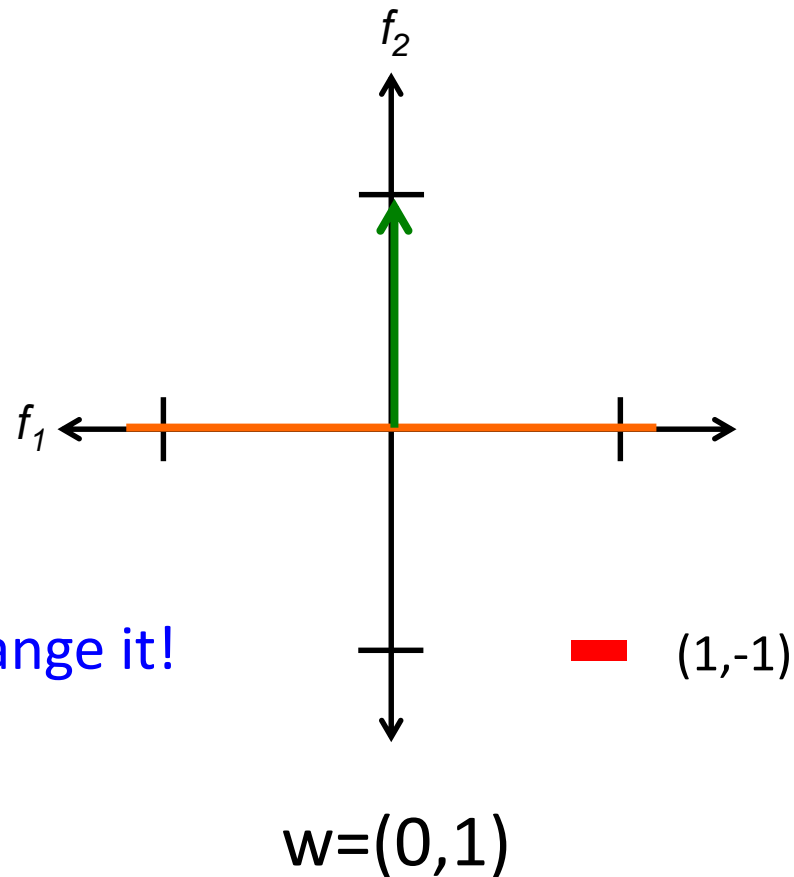
# Learning a linear classifier

$$0 = w_1 f_1 + w_2 f_2$$

$$0 * f_1 + 1 * f_2 =$$

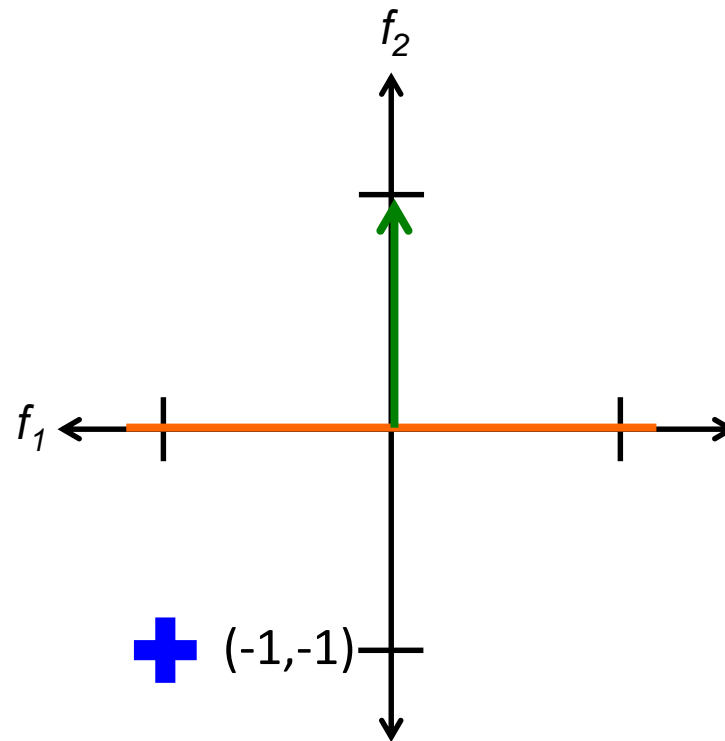
$$0 * 1 + 1 * -1 = -1$$

Already correct... don't change it!



# Learning a linear classifier

$$0 = w_1 f_1 + w_2 f_2$$



Is our current guess:  
right or wrong?

$$w = (0, 1)$$

# Learning a linear classifier

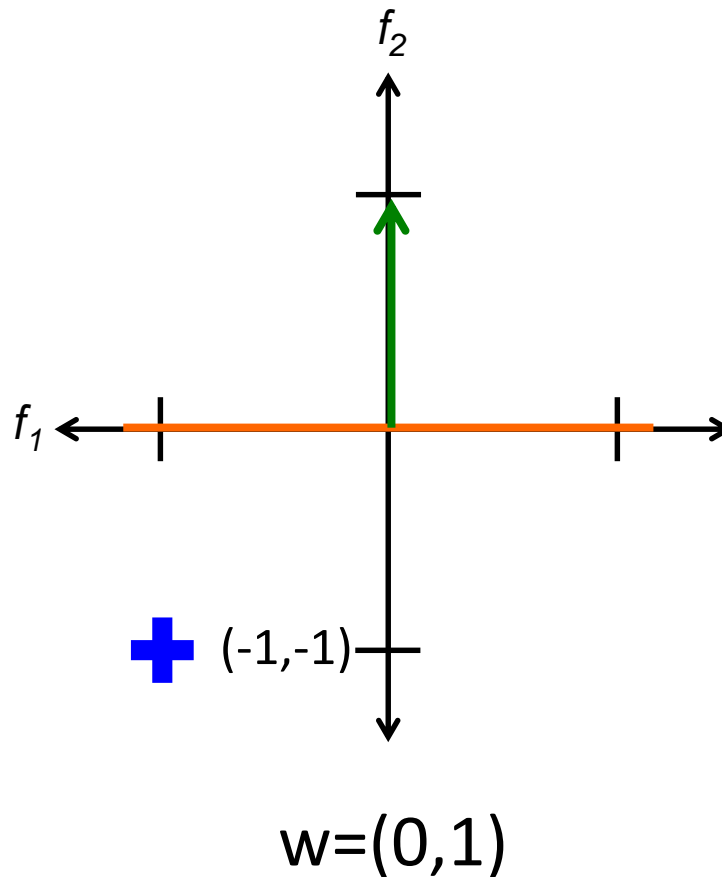
$$0 = w_1 f_1 + w_2 f_2$$

$$0 * f_1 + 1 * f_2 =$$

$$0 * -1 + 1 * -1 = -1$$

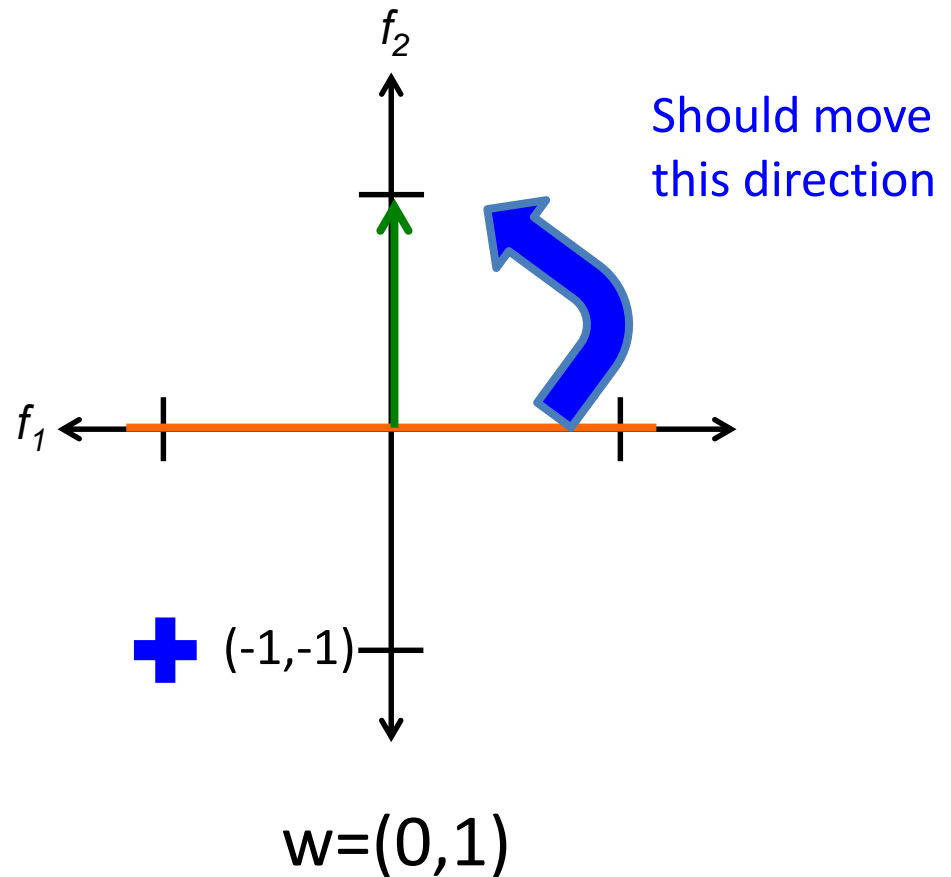
predicts negative, wrong

Geometrically, how should we update the model?



# Learning a linear classifier

$$0 = w_1 f_1 + w_2 f_2$$



# A closer look at why we got it wrong

$$w_1 \quad w_2$$

(-1, -1, positive)

$$0 * f_1 + 1 * f_2 =$$

$$0 * -1 + 1 * -1 = -1 \leftarrow$$

We'd like this value to be positive since it's a positive value



Which of these contributed to the mistake?

# A closer look at why we got it wrong

$$w_1 \quad w_2$$

(-1, -1, positive)

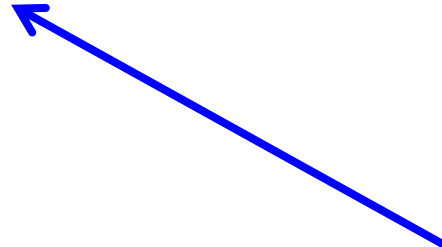
$$0 * f_1 + 1 * f_2 =$$

$$0 * -1 + 1 * -1 = -1$$

We'd like this value to be positive since it's a positive value



didn't contribute,  
but could have



contributed in the wrong  
direction

How should we change the weights?

# A closer look at why we got it wrong

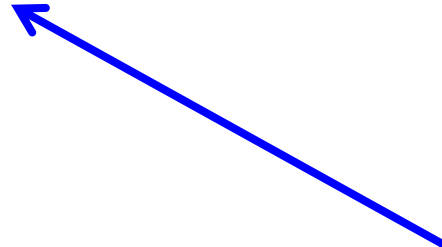
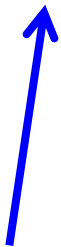
$w_1$        $w_2$

(-1, -1, positive)

$$0 * f_1 + 1 * f_2 =$$

$$0 * -1 + 1 * -1 = -1$$

We'd like this value to be positive since it's a positive value



didn't contribute,  
but could have

contributed in the wrong  
direction

decrease

0 -> -1

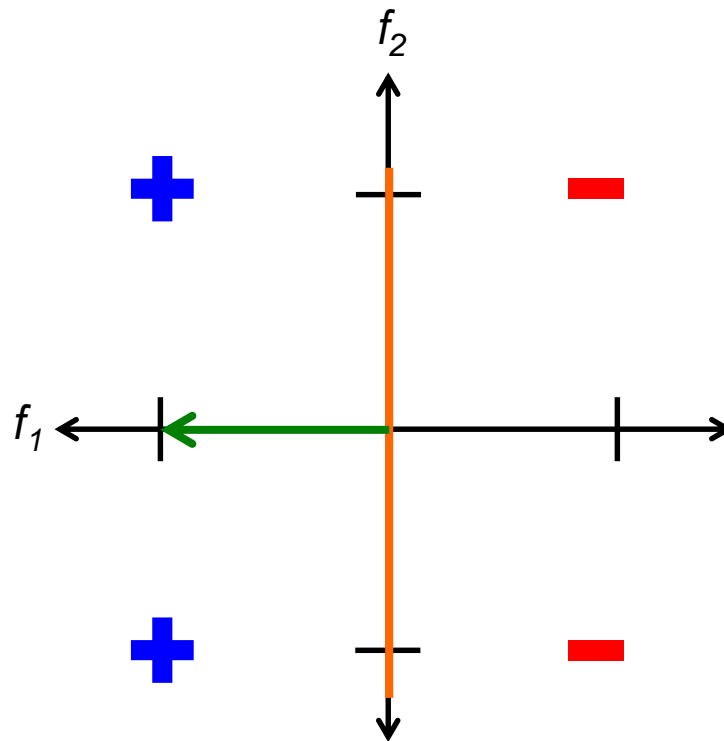
decrease

1 -> 0

# Learning a linear classifier

$f_1, f_2, label$

-1, -1, positive  
-1, 1, positive  
1, 1, negative  
1, -1, negative



$$w = (-1, 0)$$

# Perceptron learning algorithm

repeat until **convergence** (or for some # of iterations):

for each training example ( $f_1, f_2, \dots, f_n$ , label):

check if it's correct based on the current model

if **not correct**, update all the weights:

if label **positive** and feature **positive**:

**increase** weight (increase weight = predict more positive)

else if label **positive** and feature **negative**:

**decrease** weight (decrease weight = predict more positive)

else if label **negative** and feature **positive**:

**decrease** weight (decrease weight = predict more negative)

else if label **negative** and feature **negative**:

**increase** weight (increase weight = predict more negative)

# A trick...

if label positive and feature positive:

increase weight (increase weight = predict more positive)

else if label positive and feature negative:

decrease weight (decrease weight = predict more positive)

else if label negative and feature positive:

decrease weight (decrease weight = predict more negative)

else if label negative and negative weight:

increase weight (increase weight = predict more negative)

**label \*  $f_i$**

---

$$1 * 1 = 1$$

$$1 * -1 = -1$$

$$-1 * 1 = -1$$

$$-1 * -1 = 1$$

# A trick...

if label positive and feature positive:

**increase** weight (increase weight = predict more positive)

else if label positive and feature negative:

**decrease** weight (decrease weight = predict more positive)

else if label negative and feature positive:

**decrease** weight (decrease weight = predict more negative)

else if label negative and negative weight:

**increase** weight (increase weight = predict more negative)

label \*  $f_i$

$$1 * 1 = 1$$

$$1 * -1 = -1$$

$$-1 * 1 = -1$$

$$-1 * -1 = 1$$

# Perceptron learning algorithm

**repeat** until convergence (or for some # of iterations):  
**for** each training example  $(f_1, f_2, \dots, f_n, \text{label})$ :  
    check if it's correct based on the current model

**if**  $\text{prediction} * \text{label} \leq 0$ : // they don't agree  
    **for** each  $w_i$ :  
         $w_i = w_i + f_i * \text{label}$   
     $b = b + \text{label}$

How do we check if it's correct?

# Perceptron learning algorithm

**repeat** until convergence (or for some # of iterations):  
**for** each training example ( $f_1, f_2, \dots, f_n$ , label):

$$\text{prediction} = b + \sum_{i=1}^n w_i f_i$$

**if**  $\text{prediction} * \text{label} \leq 0$ : // they don't agree  
**for** each  $w_i$ :  
     $w_i = w_i + f_i * \text{label}$   
     $b = b + \text{label}$

# Perceptron learning algorithm

**repeat** until convergence (or for some # of iterations):  
**for** each training example  $(f_1, f_2, \dots, f_n, \text{label})$ :

$$\text{prediction} = b + \sum_{i=1}^n w_i f_i$$

**if**  $\text{prediction} * \text{label} \leq 0$ : // they don't agree

**for** each  $w_i$ :

$$w_i = w_i + f_i * \text{label}$$

$$b = b + \text{label}$$

Would this work for non-binary features, i.e. real-valued?

# Your turn

repeat until convergence (or for some # of iterations):

for each training example  $(f_1, f_2, \dots, f_n, \text{label})$ :

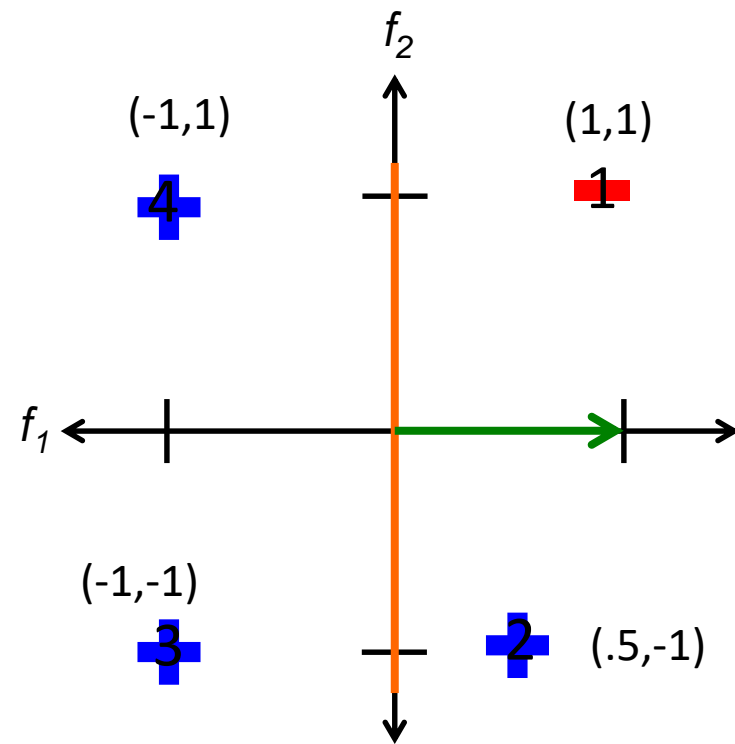
$$\text{prediction} = \hat{a}_{i=1}^n w_i f_i$$

if  $\text{prediction} * \text{label} \leq 0$ : // they don't agree

for each  $w_i$ :

$$w_i = w_i + f_i * \text{label}$$

- Repeat until convergence
- Keep track of  $w_1, w_2$  as they change
- Redraw the line after each step



$$w = (1, 0)$$

# Your turn

repeat until convergence (or for some # of iterations):

for each training example  $(f_1, f_2, \dots, f_n, \text{label})$ :

$$\text{prediction} = \hat{a}_{i=1}^n w_i f_i$$

if  $\text{prediction} * \text{label} \leq 0$ : // they don't agree

for each  $w_i$ :

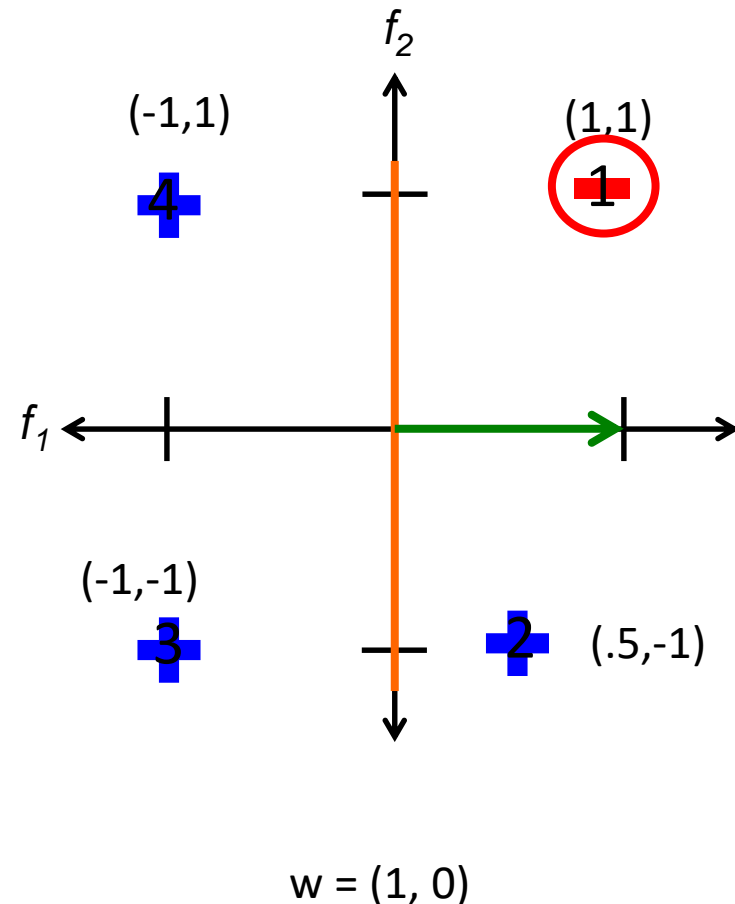
$$w_i = w_i + f_i * \text{label}$$

$$0 = 1 * f_1 - 0 * f_2$$

$$1 * 1 + 0 * 1 = 1 > 0 : \text{positive}$$

$$w_1 = 1 + 1 * (-1) = 0$$

$$w_2 = 0 + 1 * (-1) = -1$$



# Your turn

repeat until convergence (or for some # of iterations):

for each training example  $(f_1, f_2, \dots, f_n, \text{label})$ :

$$\text{prediction} = \hat{a}_{i=1}^n w_i f_i$$

if  $\text{prediction} * \text{label} \leq 0$ : // they don't agree

for each  $w_i$ :

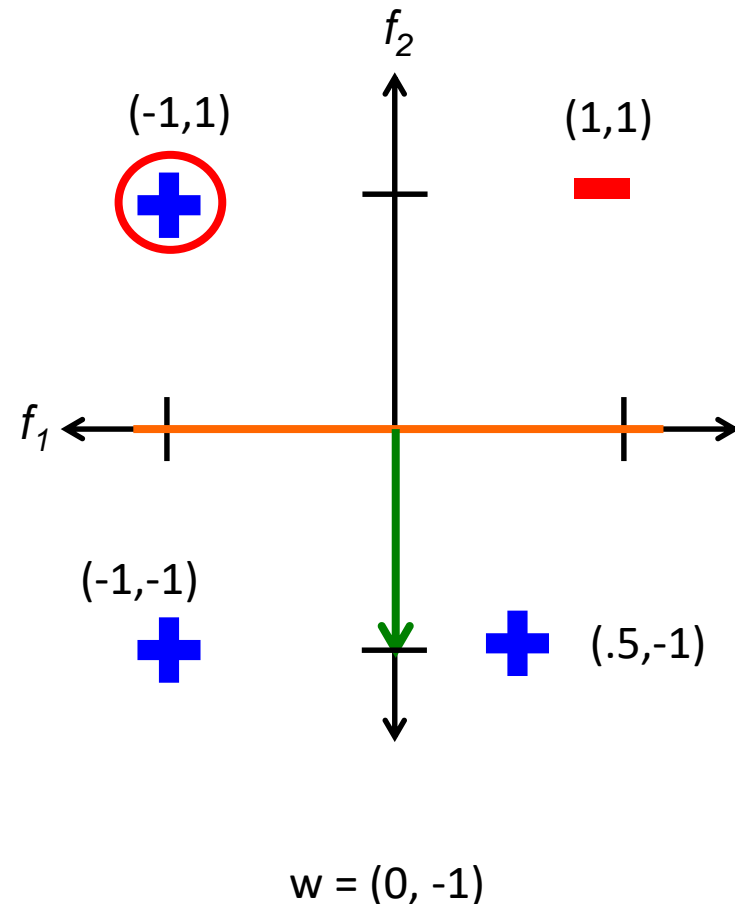
$$w_i = w_i + f_i * \text{label}$$

$$0 = 0 * f_1 - 1 * f_2$$

$$0 * (-1) + (-1) * 1 = -2 < 0 : \text{negative}$$

$$w_1 = 0 + (-1) * 1 = -1$$

$$w_2 = -1 + (1) * 1 = 0$$



# Your turn

repeat until convergence (or for some # of iterations):

for each training example  $(f_1, f_2, \dots, f_n, \text{label})$ :

$$\text{prediction} = \hat{a}_{i=1}^n w_i f_i$$

if  $\text{prediction} * \text{label} \leq 0$ : // they don't agree

for each  $w_i$ :

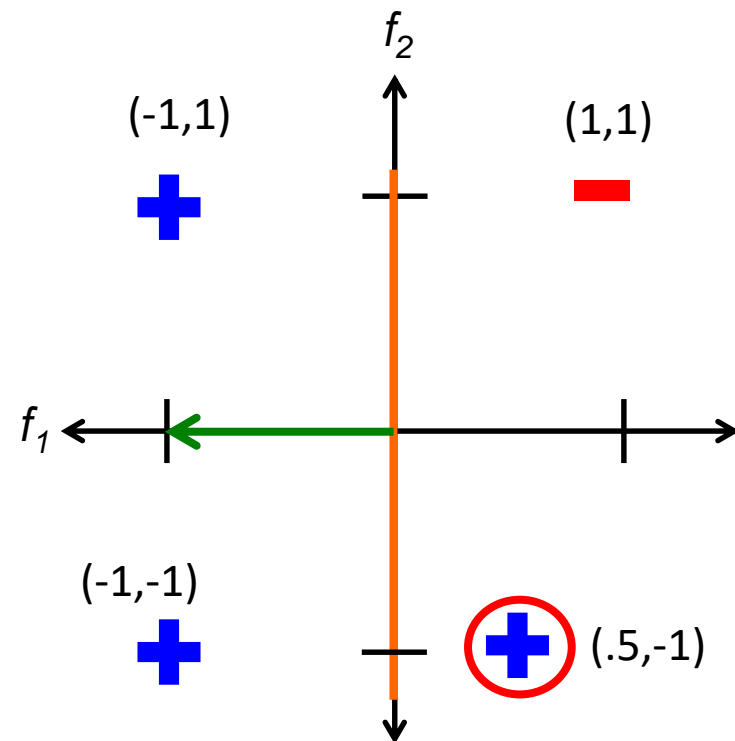
$$w_i = w_i + f_i * \text{label}$$

$$0 = -1 * f_1$$

$$(-1) * 0.5 = -0.5 < 0 : \text{negative}$$

$$w_1 = -1 + 0.5 * 1 = -0.5$$

$$w_2 = 0 + (-1) * 1 = -1$$



$$w = (-1, 0)$$

# Your turn

repeat until convergence (or for some # of iterations):

for each training example  $(f_1, f_2, \dots, f_n, \text{label})$ :

$$\text{prediction} = \sum_{i=1}^n w_i f_i$$

if  $\text{prediction} * \text{label} \leq 0$ : // they don't agree

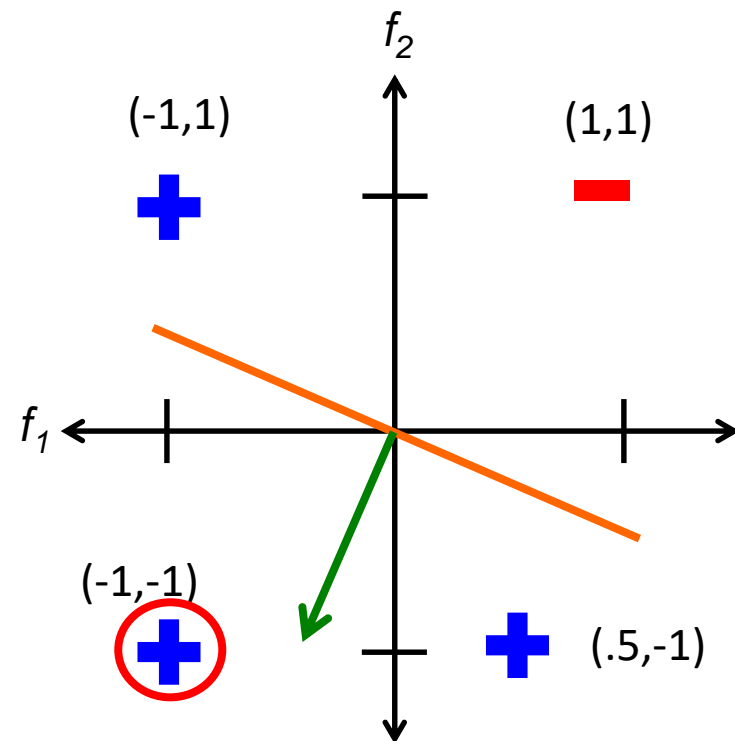
for each  $w_i$ :

$$w_i = w_i + f_i * \text{label}$$

$$0 = -0.5 * f_1 + (-1) * f_2$$

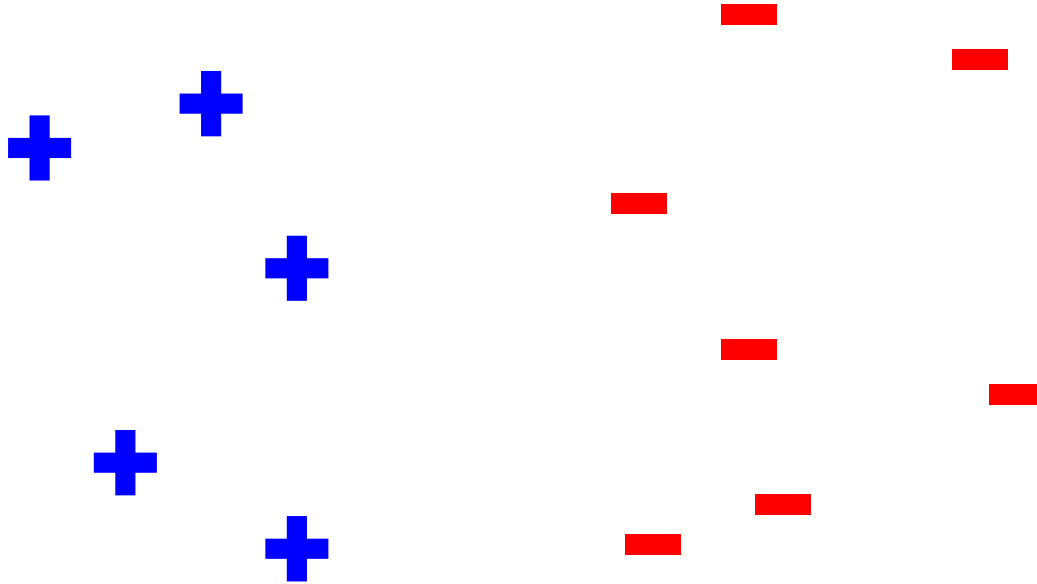
$$0.5 + 1 = 1.5 > 0 : \text{positive}$$

**No update**

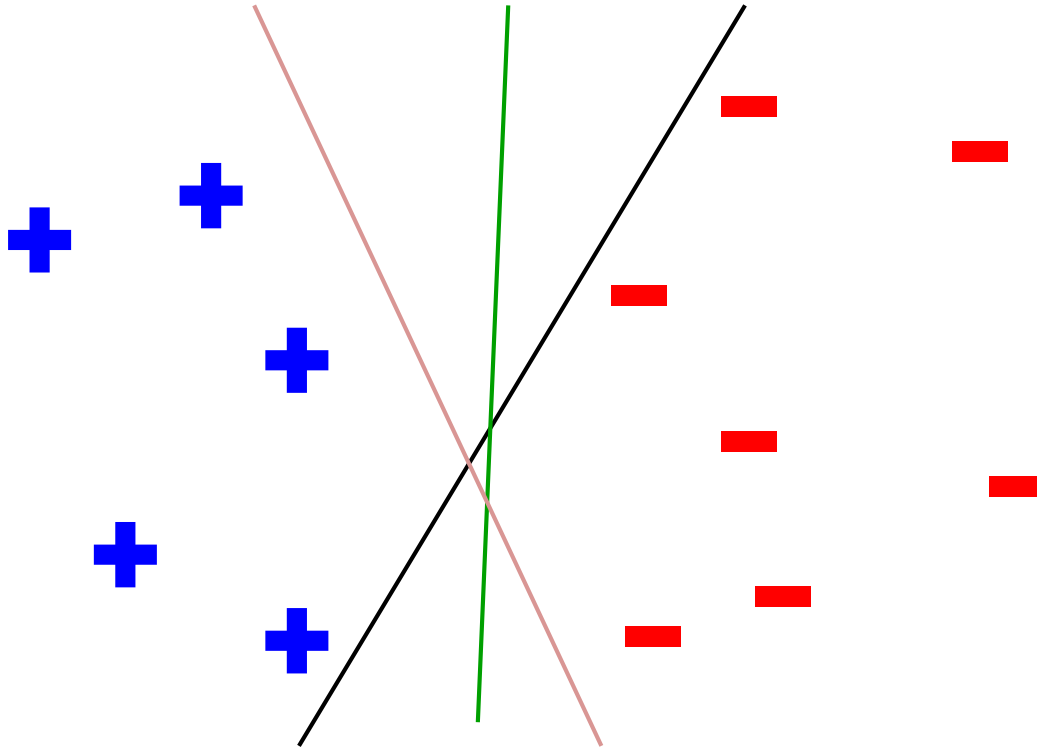


$$w = (-0.5, -1)$$

Which line will it find?



# Which line will it find?



Only guaranteed to find *some*  
line that separates the data

# Convergence

repeat until convergence (or for some # of iterations):

for each training example ( $f_1, f_2, \dots, f_n$ , label):

$$\text{prediction} = b + \sum_{i=1}^n w_i f_i$$

if  $\text{prediction} * \text{label} \leq 0$ : // they don't agree

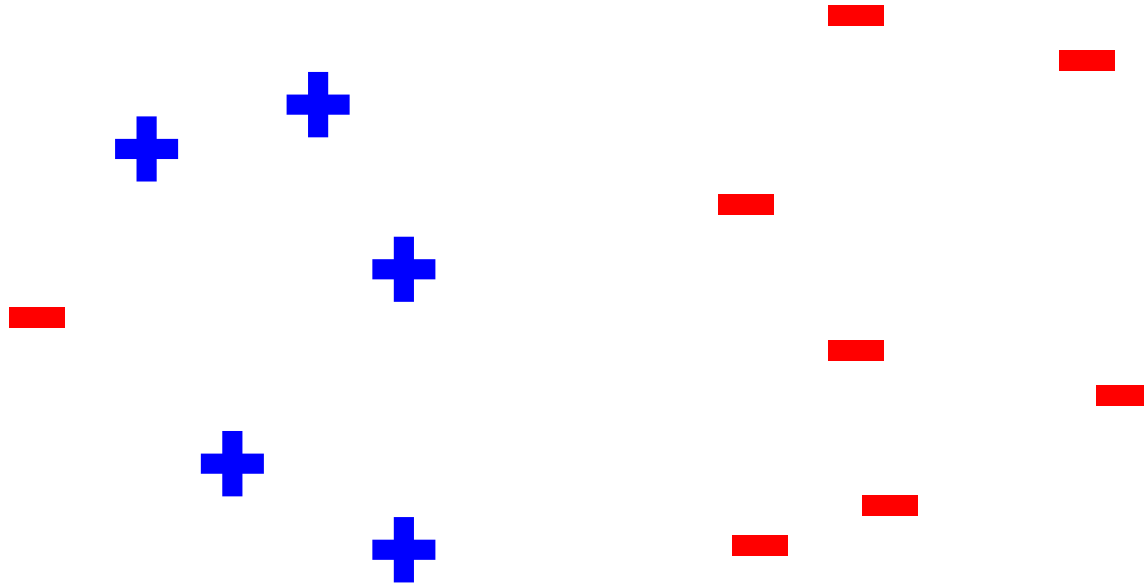
for each  $w_i$ :

$$w_i = w_i + f_i * \text{label}$$

$$b = b + \text{label}$$

Why do we also have the “some # iterations” check?

# Handling non-separable data



If we ran the algorithm on this it would never converge!

# Convergence

repeat until convergence (or **for some # of iterations**):  
for each training example ( $f_1, f_2, \dots, f_n$ , label):

$$\text{prediction} = b + \sum_{i=1}^n w_i f_i$$

if  $\text{prediction} * \text{label} \leq 0$ : // they don't agree  
for each  $w_i$ :

$$w_i = w_i + f_i * \text{label}$$

$$b = b + \text{label}$$

Also helps avoid overfitting!

(This is harder to see in 2-D examples)

# Ordering

repeat until convergence (or for some # of iterations):

for each training example  $(f_1, f_2, \dots, f_n, \text{label})$ :

$$\text{prediction} = b + \sum_{i=1}^n w_i f_i$$

if  $\text{prediction} * \text{label} \leq 0$ : // they don't agree

for each  $w_i$ :

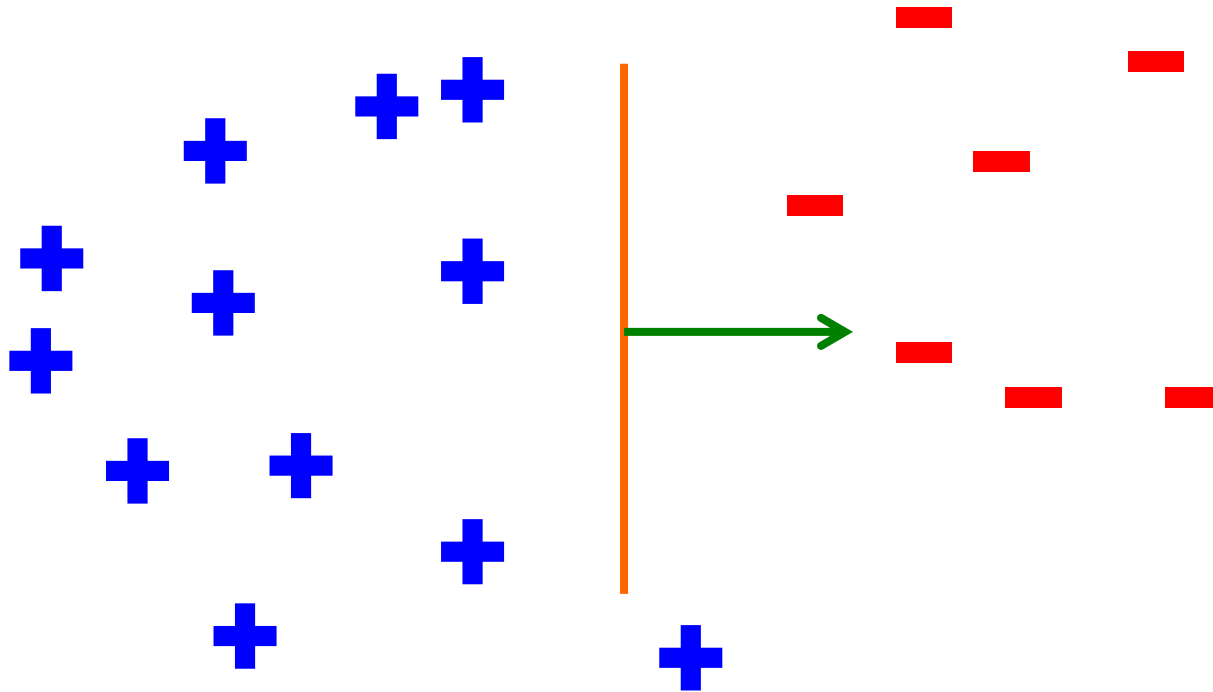
$$w_i = w_i + f_i * \text{label}$$

$$b = b + \text{label}$$

What order should we traverse the examples?

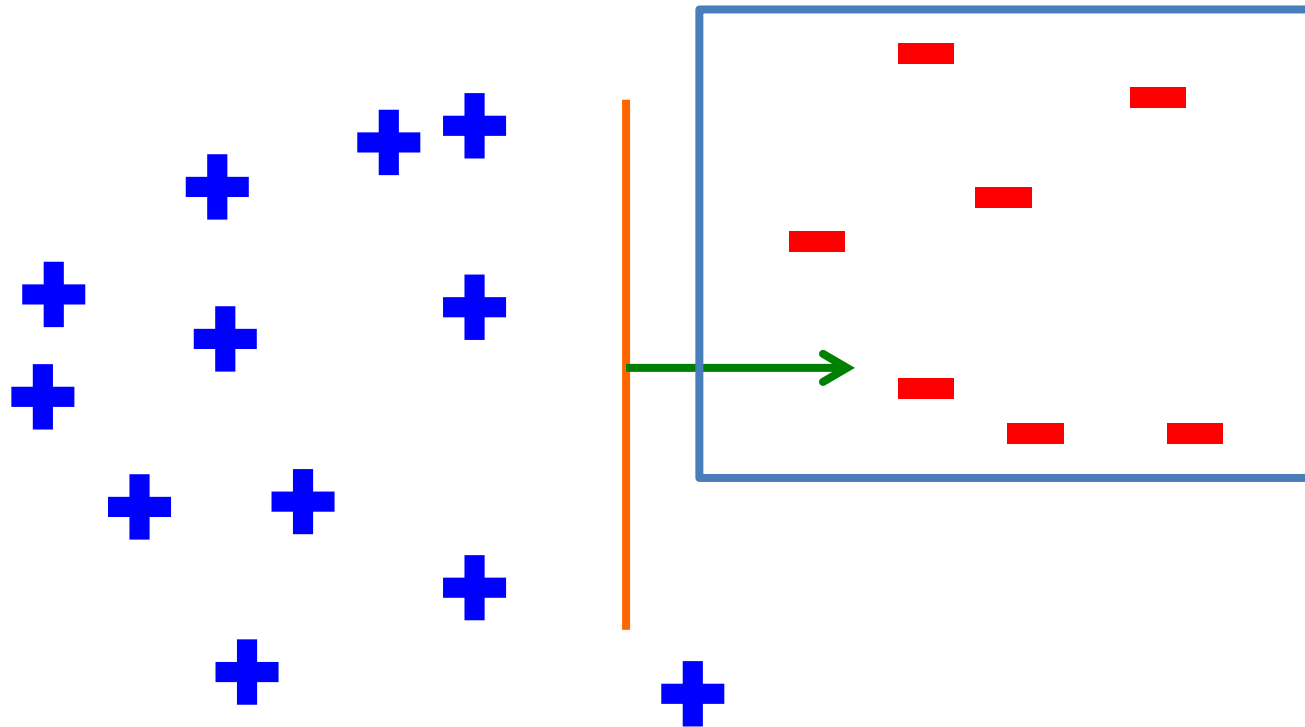
Does it matter?

# Order matters

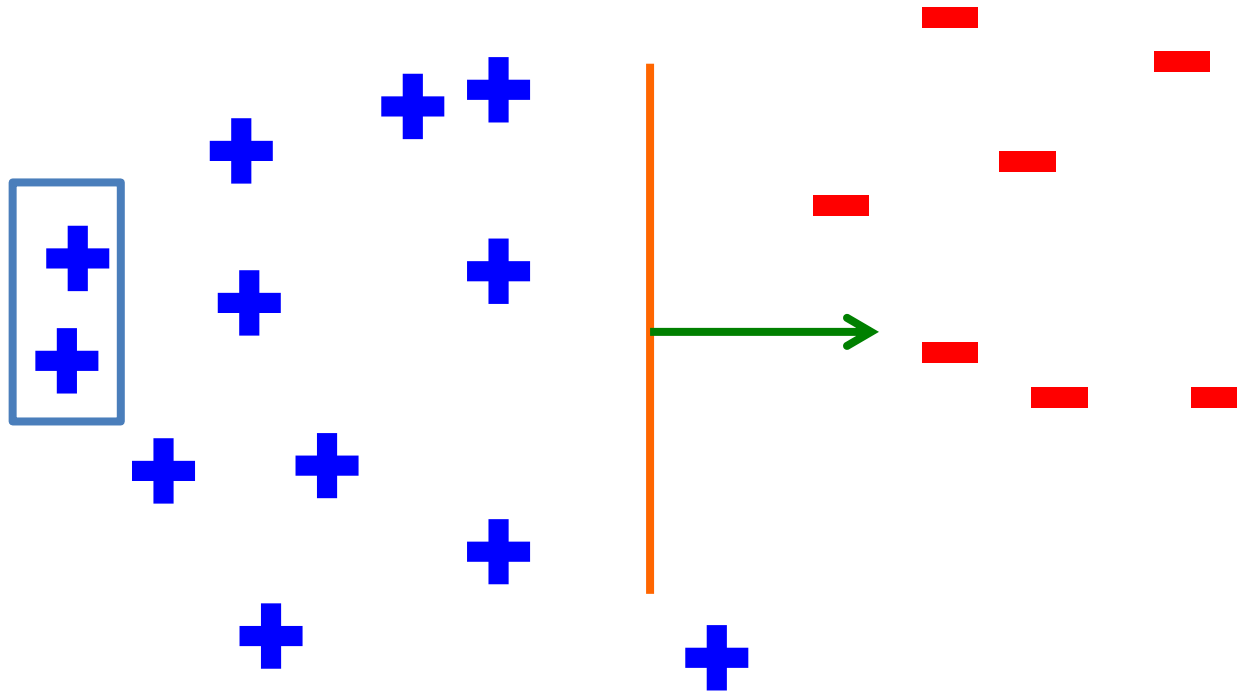


What would be a good/bad order?

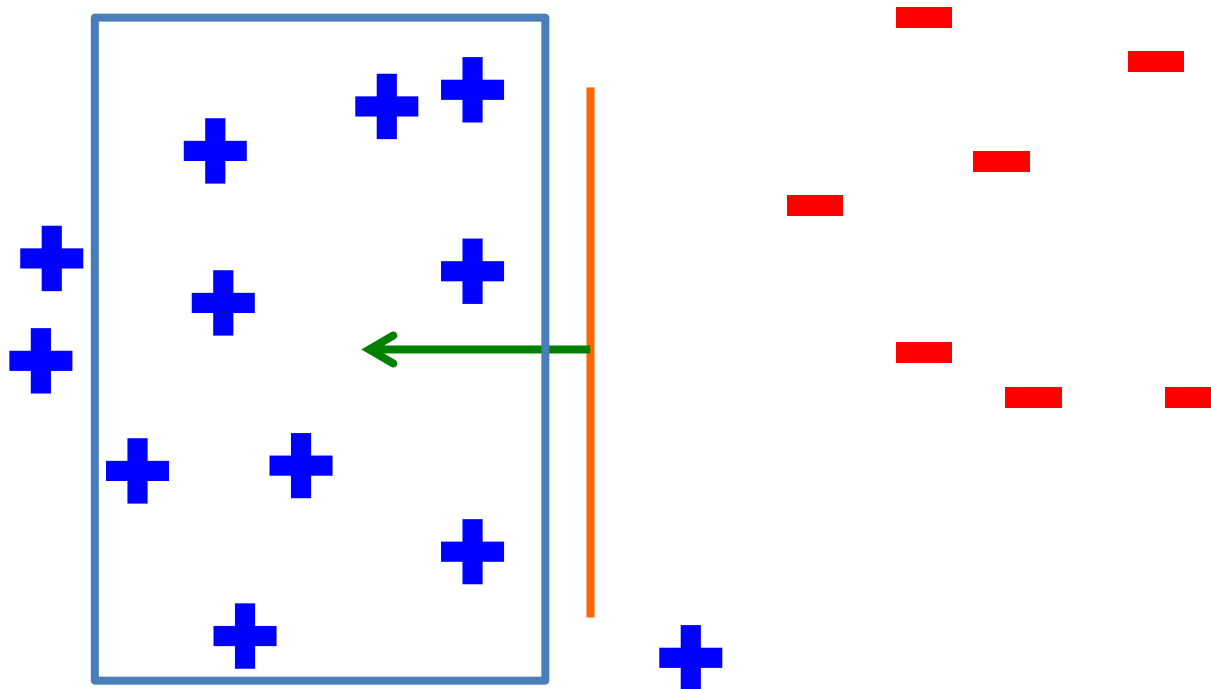
# Order matters: a bad order



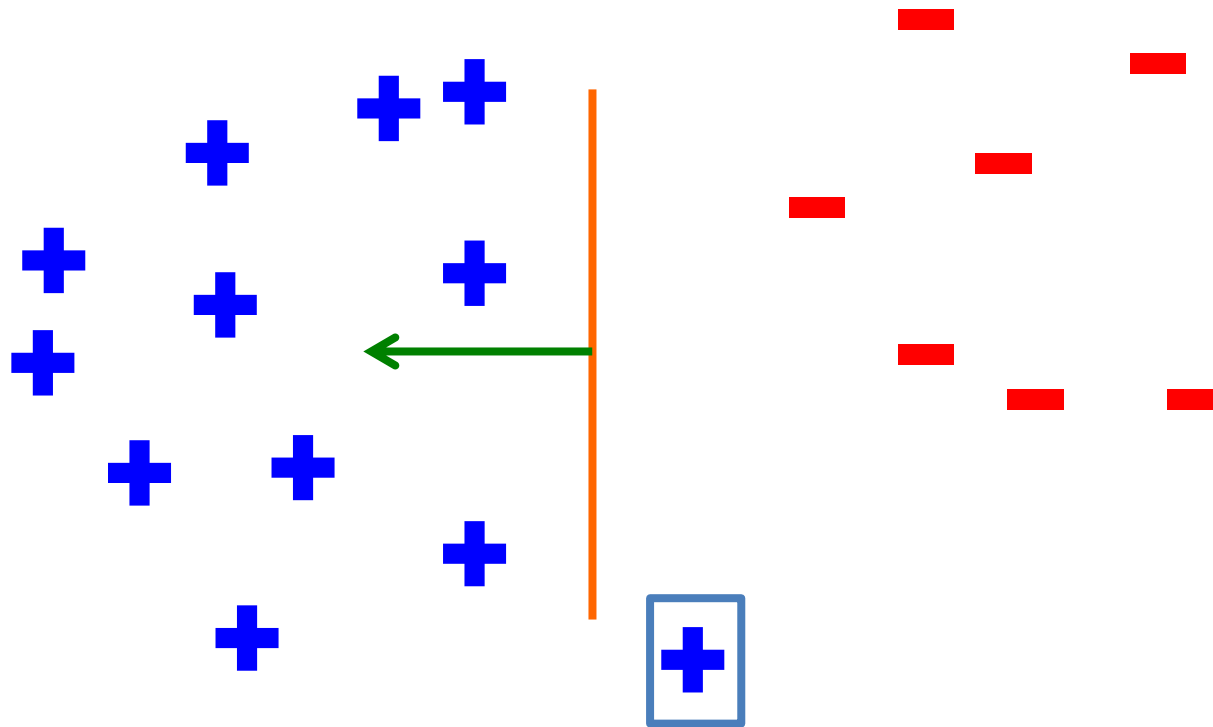
# Order matters: a bad order



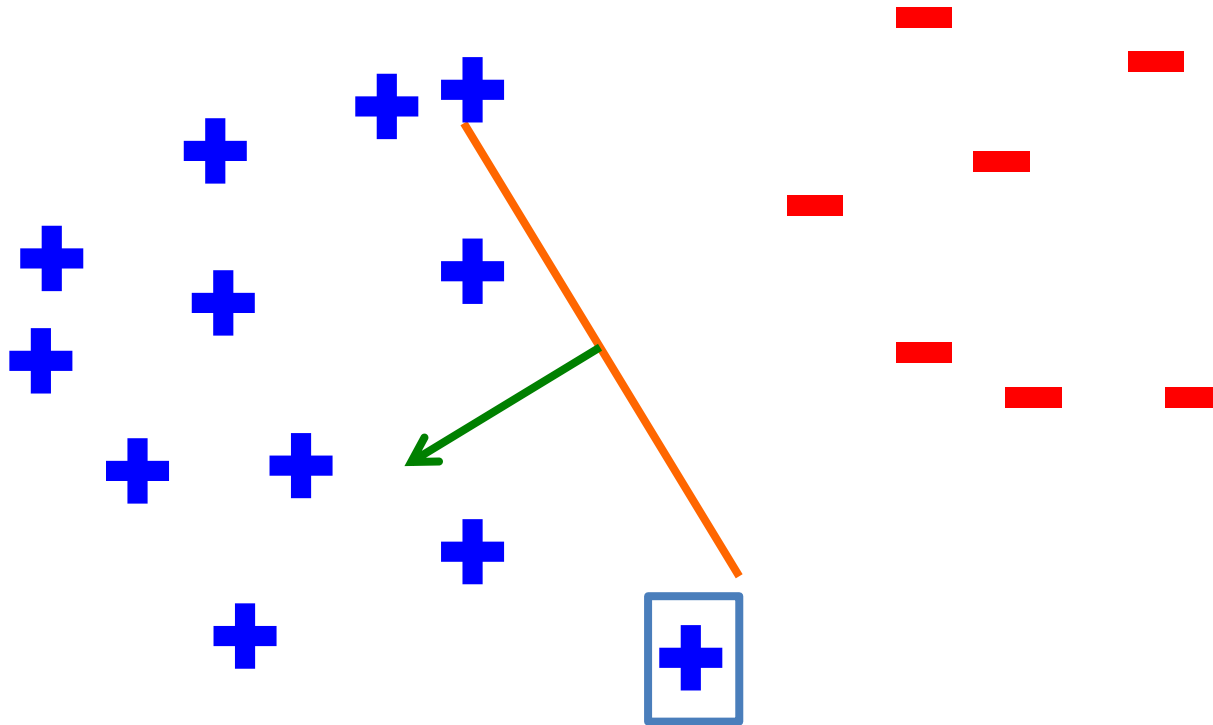
# Order matters: a bad order



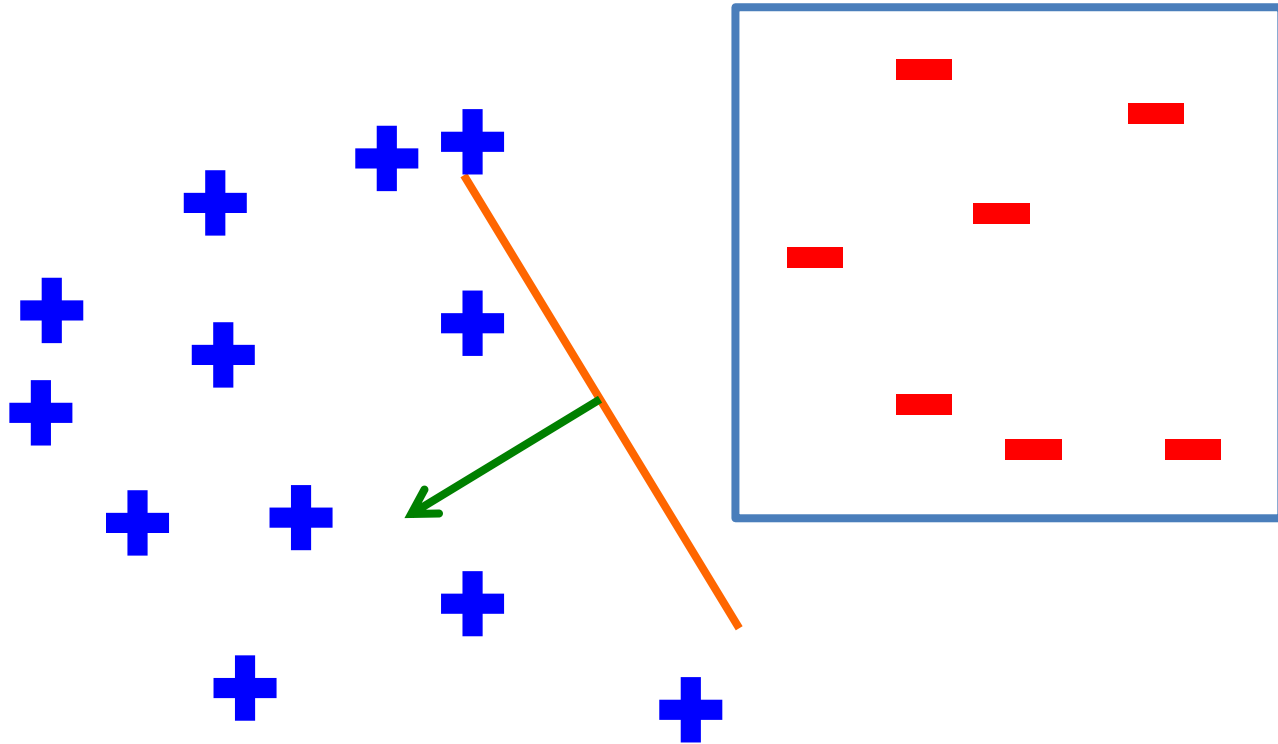
# Order matters: a bad order



# Order matters: a bad order



# Order matters: a bad order



Solution?

# Ordering

repeat until convergence (or for some # of iterations):

randomize order of training examples

for each training example  $(f_1, f_2, \dots, f_n, \text{label})$ :

$$\text{prediction} = b + \sum_{i=1}^n w_i f_i$$

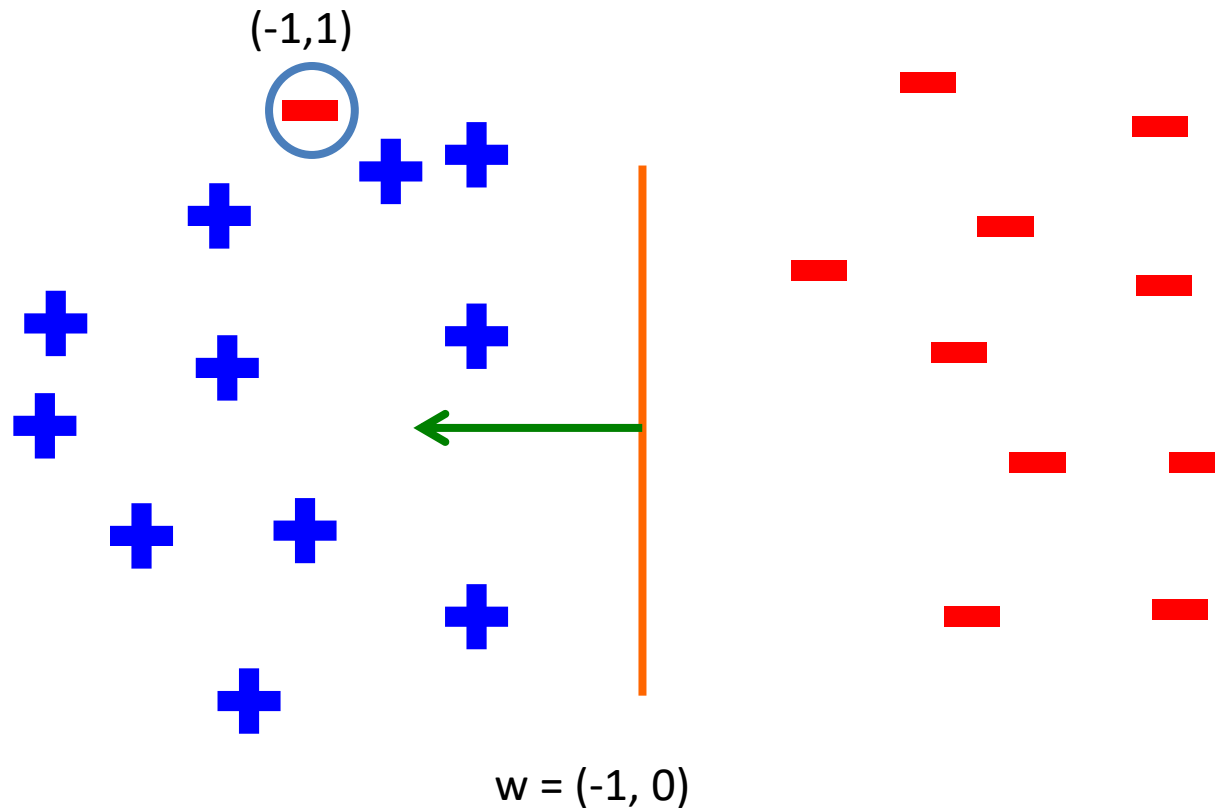
if  $\text{prediction} * \text{label} \leq 0$ : // they don't agree

for each  $w_i$ :

$$w_i = w_i + f_i * \text{label}$$

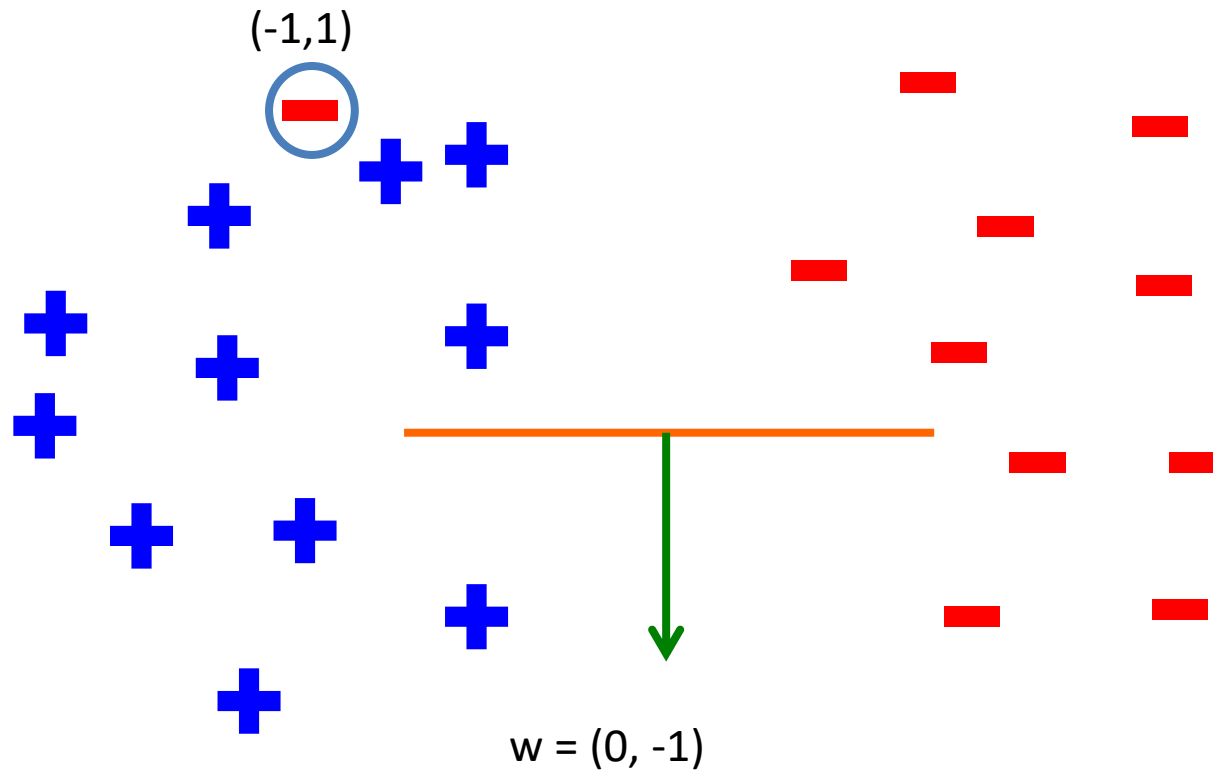
$$b = b + \text{label}$$

# Improvements



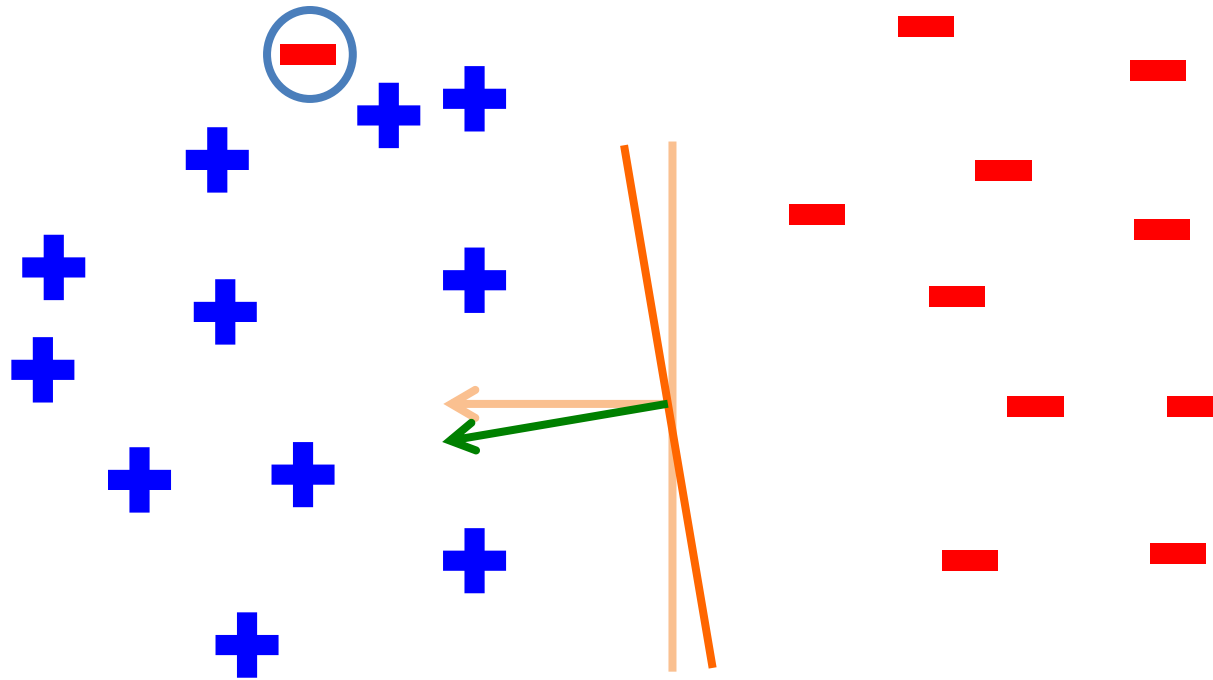
What will happen when we examine this example?

# Improvements



Does this make sense? What if we had previously gone through ALL of the other examples correctly?

# Improvements



Maybe just move it slightly in the direction of correction

# Voted perceptron learning

## Training

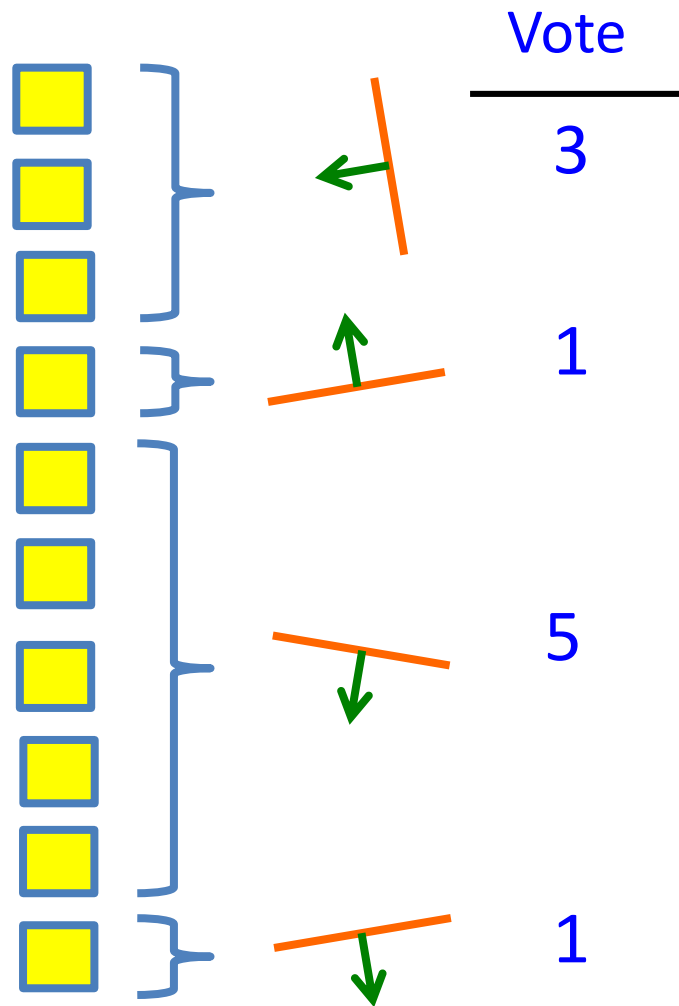
### -every time a mistake is made on an example:

- store the weights (i.e. before changing for current example)
- store the number of examples that set of weights got correct

## Classify

- calculate the prediction from **ALL saved weights**
- multiply each prediction by the number it got correct (i.e., a weighted vote) and take the sum over all predictions
- said another way: pick whichever prediction has the most votes

# Voted perceptron learning

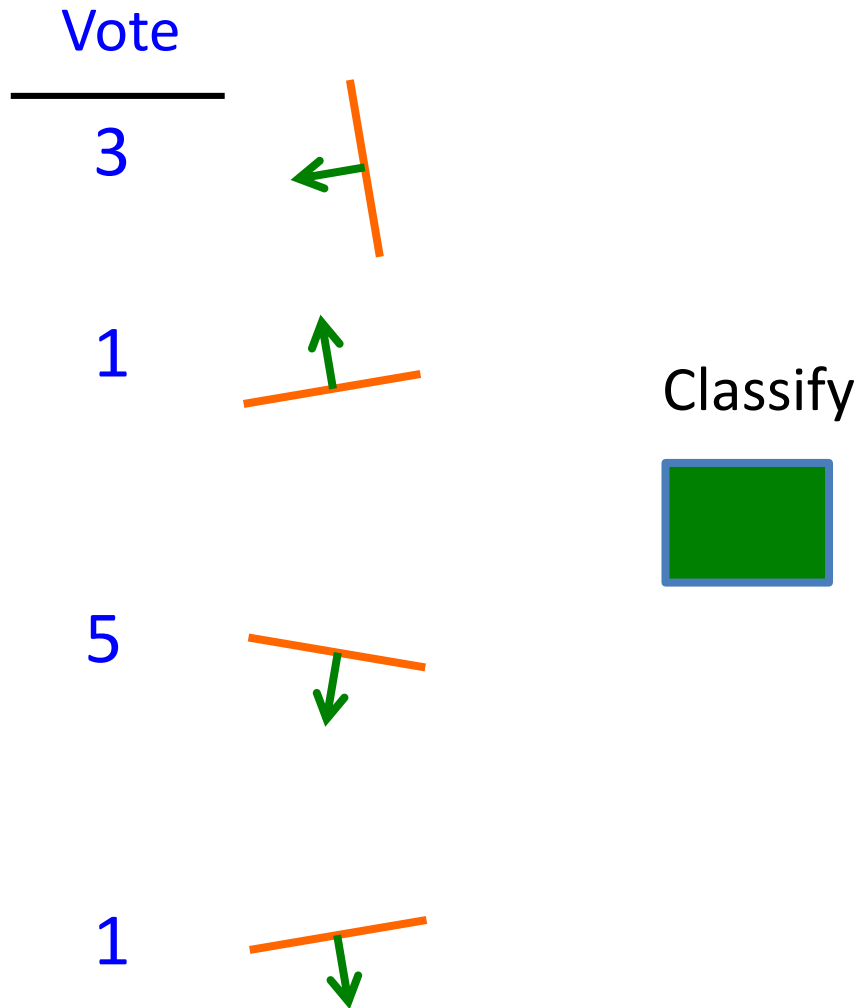


## Training

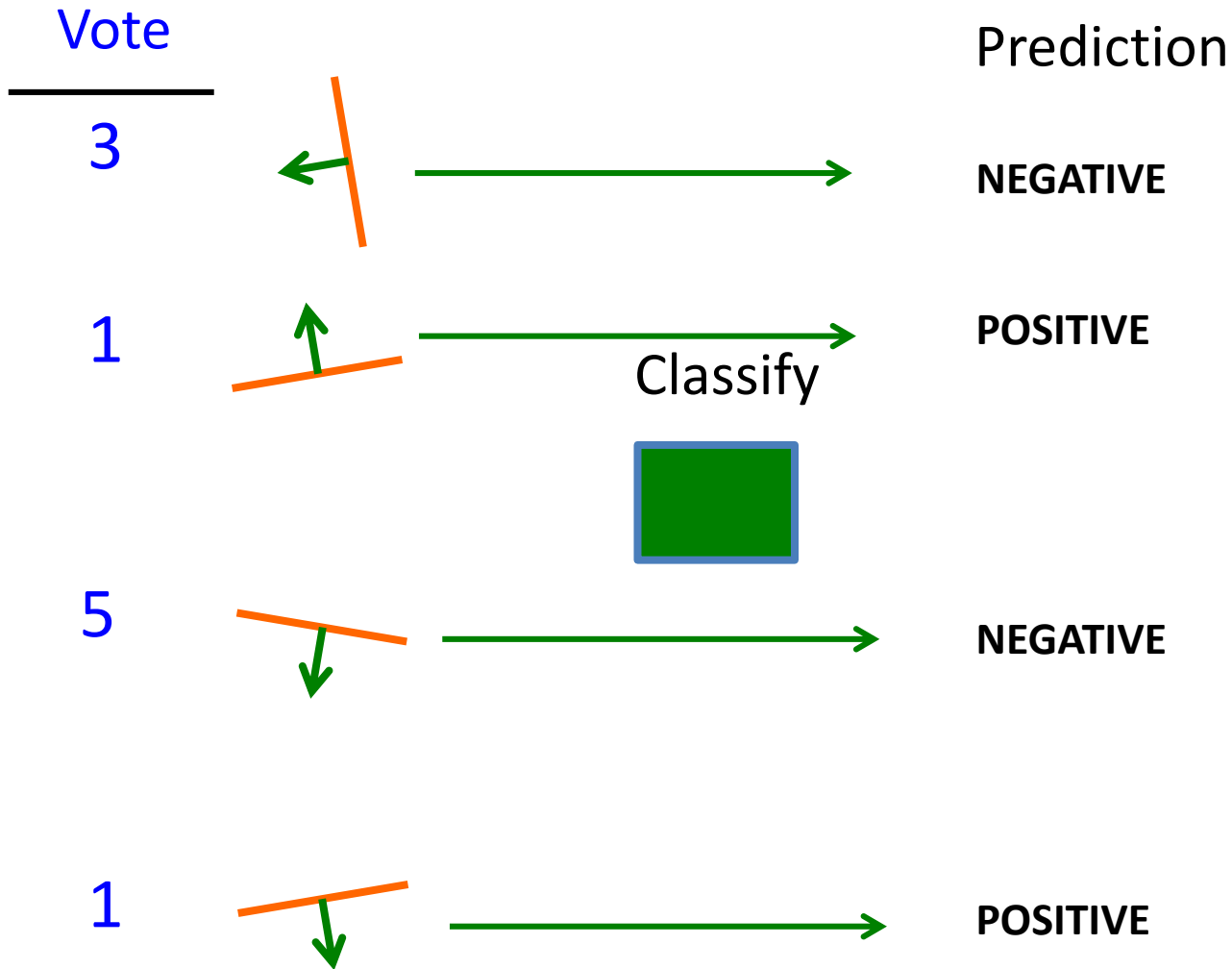
every time a mistake is made on an example:

- store the weights
- store the number of examples that set of weights got correct

# Voted perceptron learning

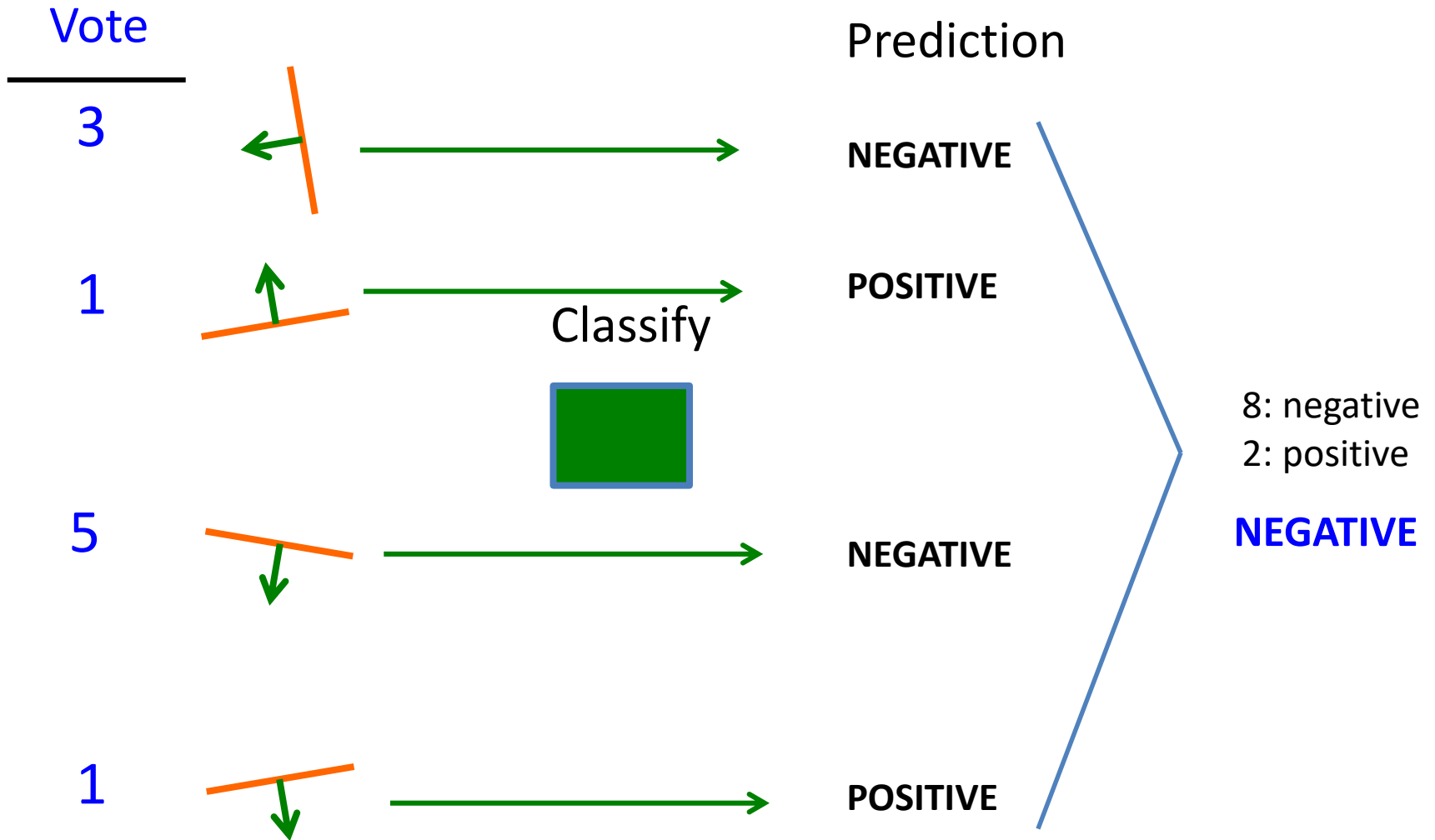


# Voted perceptron learning



Decision?

# Voted perceptron learning



# Voted perceptron learning

Works much better in practice

Avoids overfitting, though it can still happen

Avoids big changes in the result by examples examined at the end of training

# Voted perceptron learning

## Training

- every time a mistake is made on an example:
  - store the weights (i.e. before changing for current example)
  - store the number of examples that set of weights got correct

## Classify

- calculate the prediction from ALL saved weights
- multiply each prediction by the number it got correct (i.e a weighted vote) and take the sum over all predictions
- said another way: pick whichever prediction has the most votes

Any issues/concerns?

# Voted perceptron learning

## Training

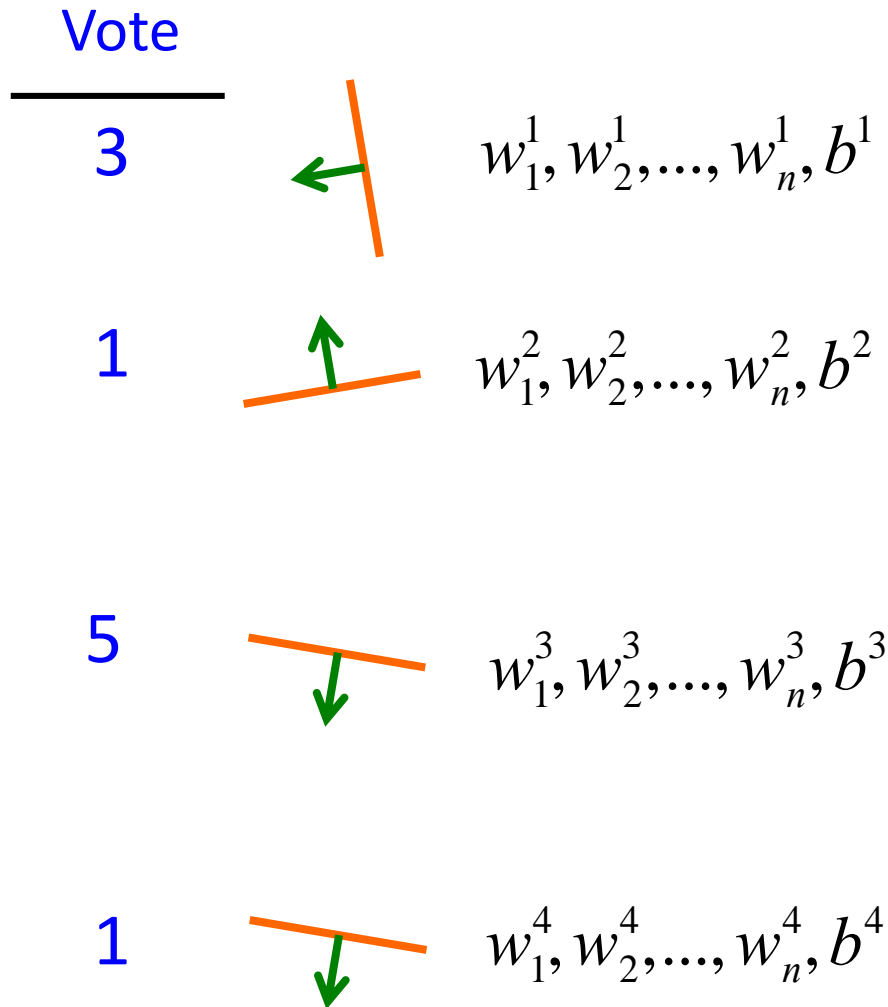
- every time a mistake is made on an example:
  - store the weights (i.e. before changing for current example)
  - store the number of examples that set of weights got correct

## Classify

- calculate the prediction from ALL saved weights
- multiply each prediction by the number it got correct (i.e a weighted vote) and take the sum over all predictions
- said another way: pick whichever prediction has the most votes

1. Can require a lot of storage
2. Classifying becomes very, very expensive

# Average perceptron



$$\bar{w}_i = \frac{3w_i^1 + 1w_i^2 + 5w_i^3 + 1w_i^4}{10}$$

The final weights are the *weighted average* of the previous weights

Can just keep a running average!

# Perceptron learning algorithm

repeat until convergence (or for some # of iterations):

for each training example  $(f_1, f_2, \dots, f_n, \text{label})$ :

$$\text{prediction} = b + \sum_{i=1}^n w_i f_i$$

if  $\text{prediction} * \text{label} \leq 0$ : // they don't agree

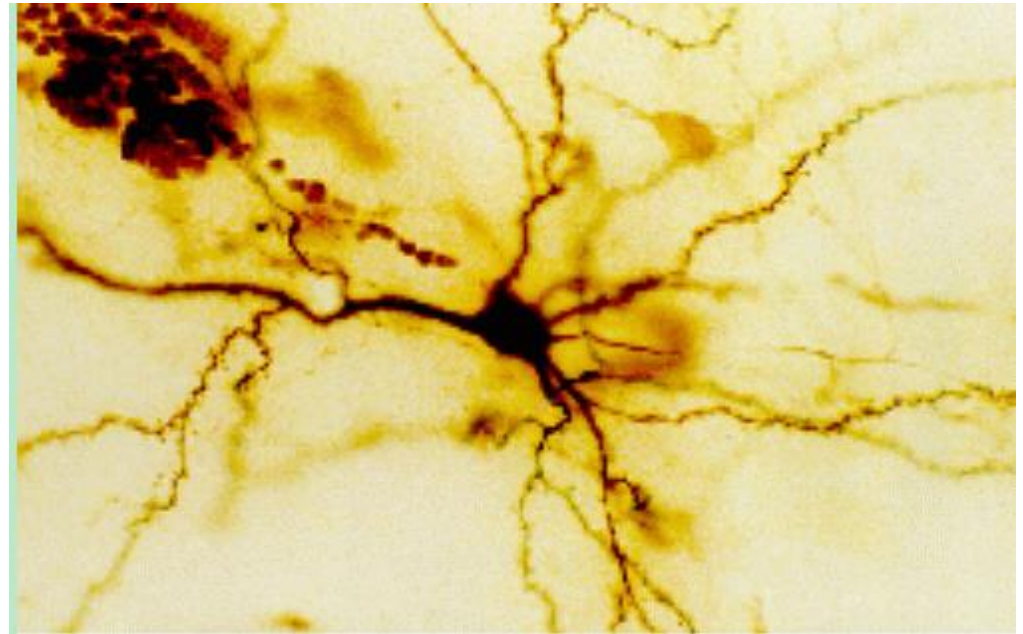
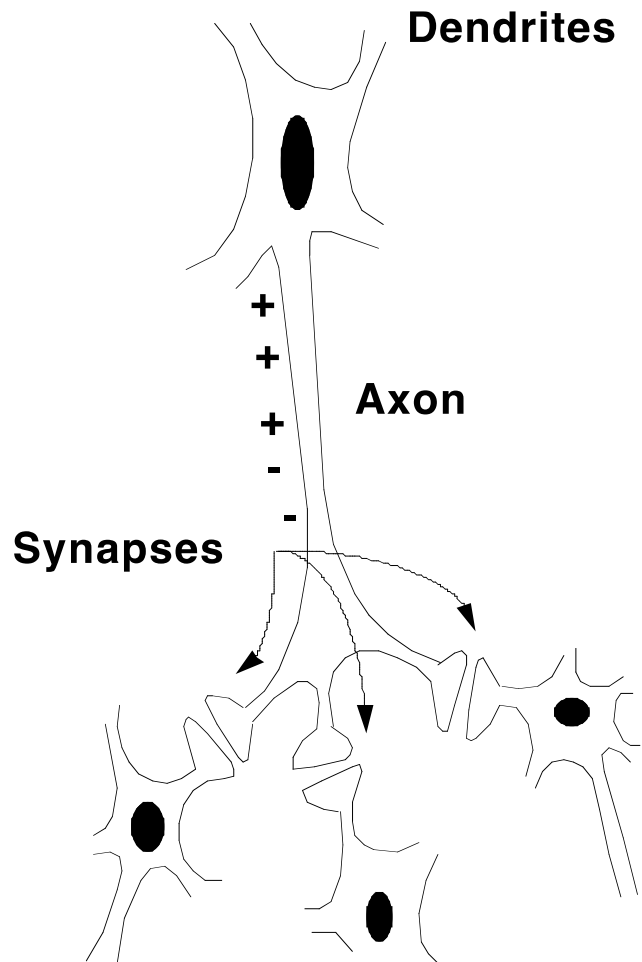
for each  $w_i$ :

$$w_i = w_i + f_i * \text{label}$$

$$b = b + \text{label}$$

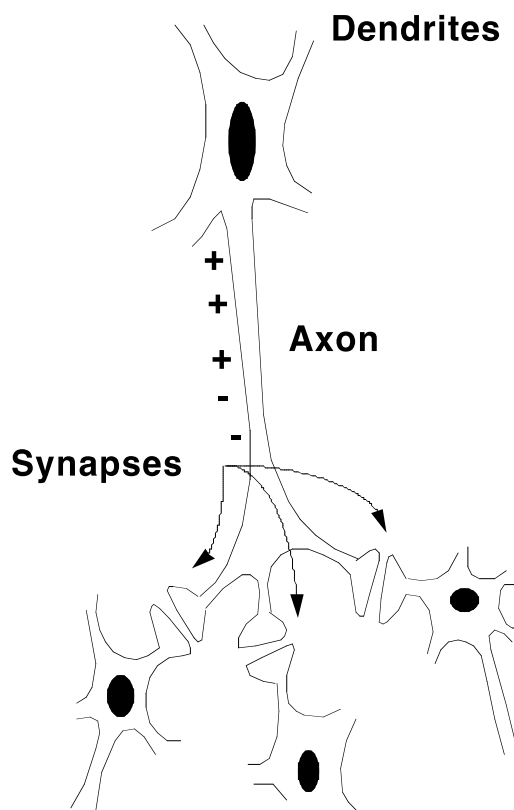
Why is it called the “perceptron” learning algorithm if what it learns is a line? Why not “line learning” algorithm?

# Our Nervous System



Neuron

# Our nervous system: *the computer science view*



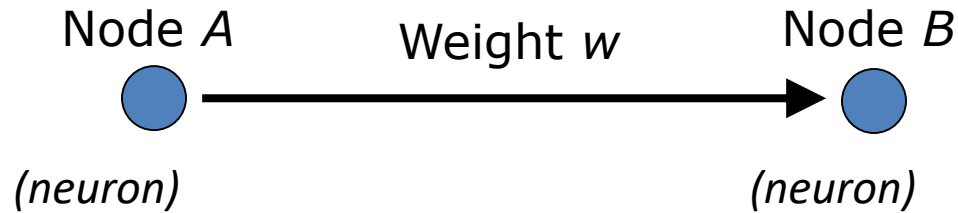
the human brain is a large collection of interconnected neurons

a **NEURON** is a brain cell

collect, process, and disseminate electrical signals

Neurons are connected via synapses

They **FIRE** depending on the conditions of the neighboring neurons



$w$  is the strength of signal sent between A and B.

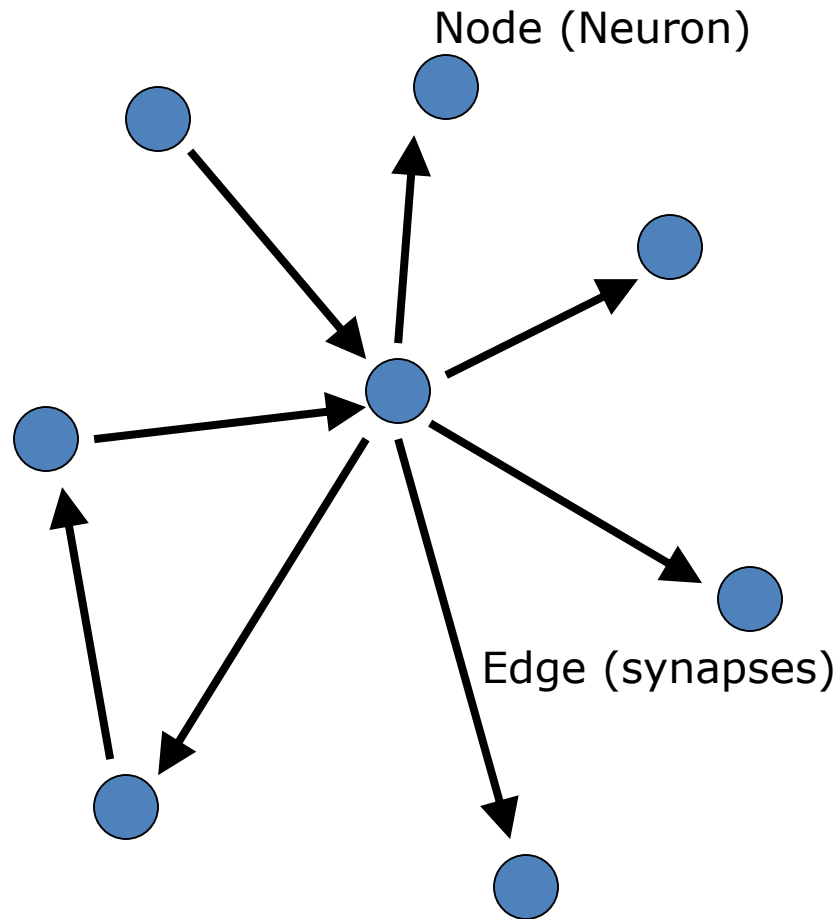
If A fires and  **$w$  is positive**, then A **stimulates** B.

If A fires and  **$w$  is negative**, then A **inhibits** B.

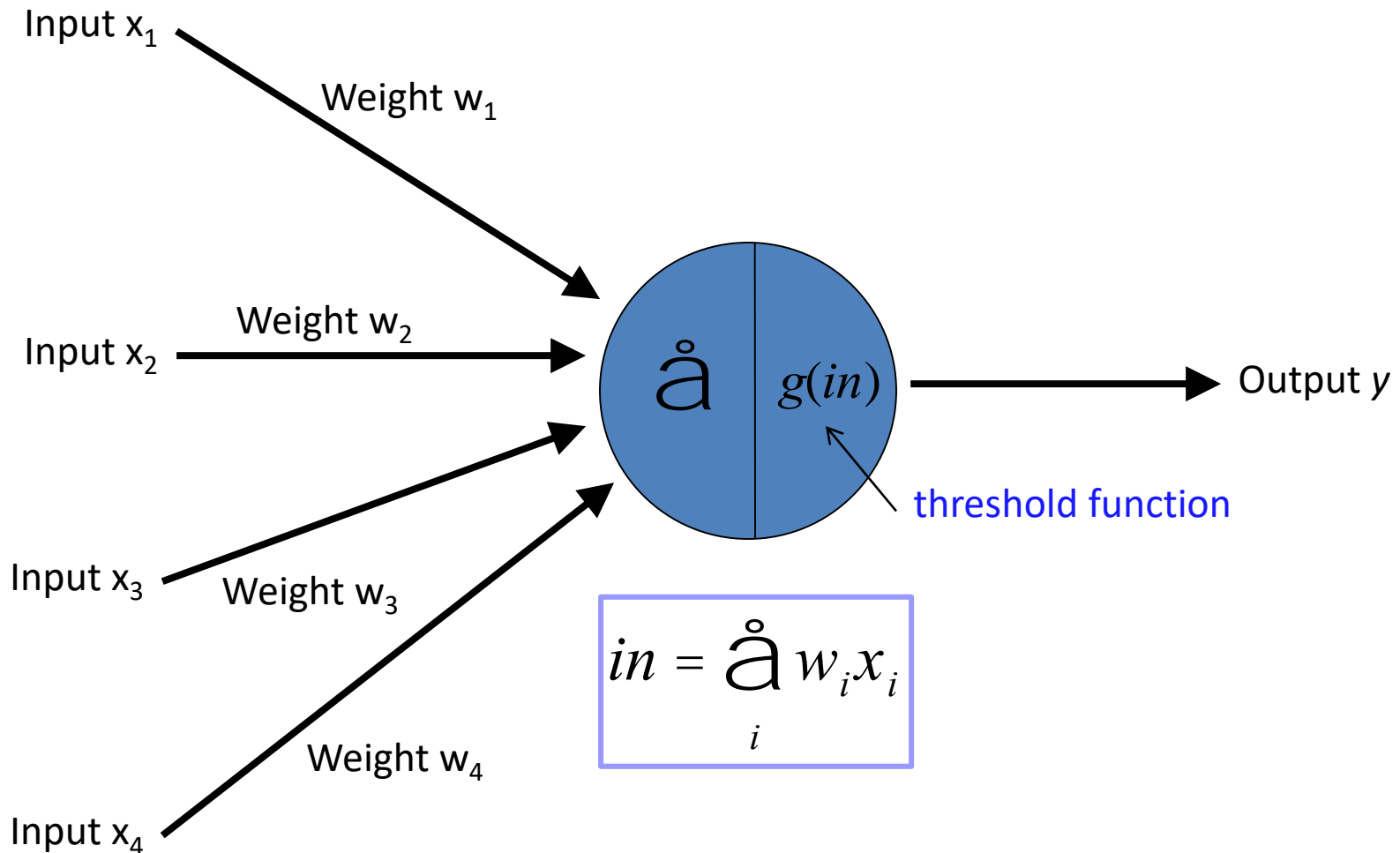
If a node is stimulated enough, then it also fires.

How much stimulation is required is determined by its **threshold**.

# Neural Networks



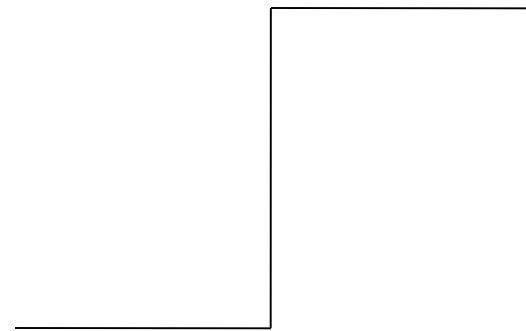
# A Single Neuron/Perceptron



# Possible threshold functions

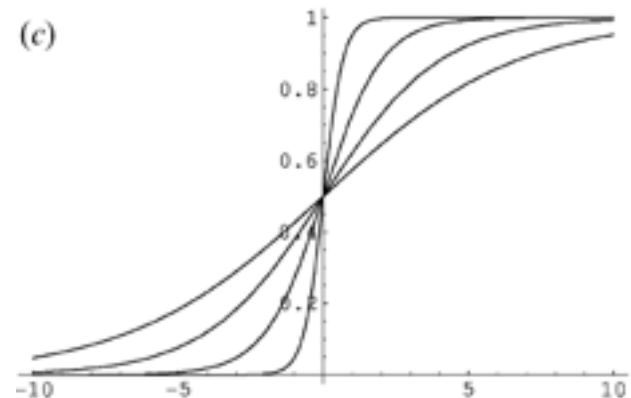
hard threshold:

if  $in$  (the sum of weights)  $\geq$  *threshold* 1  
else 0 otherwise

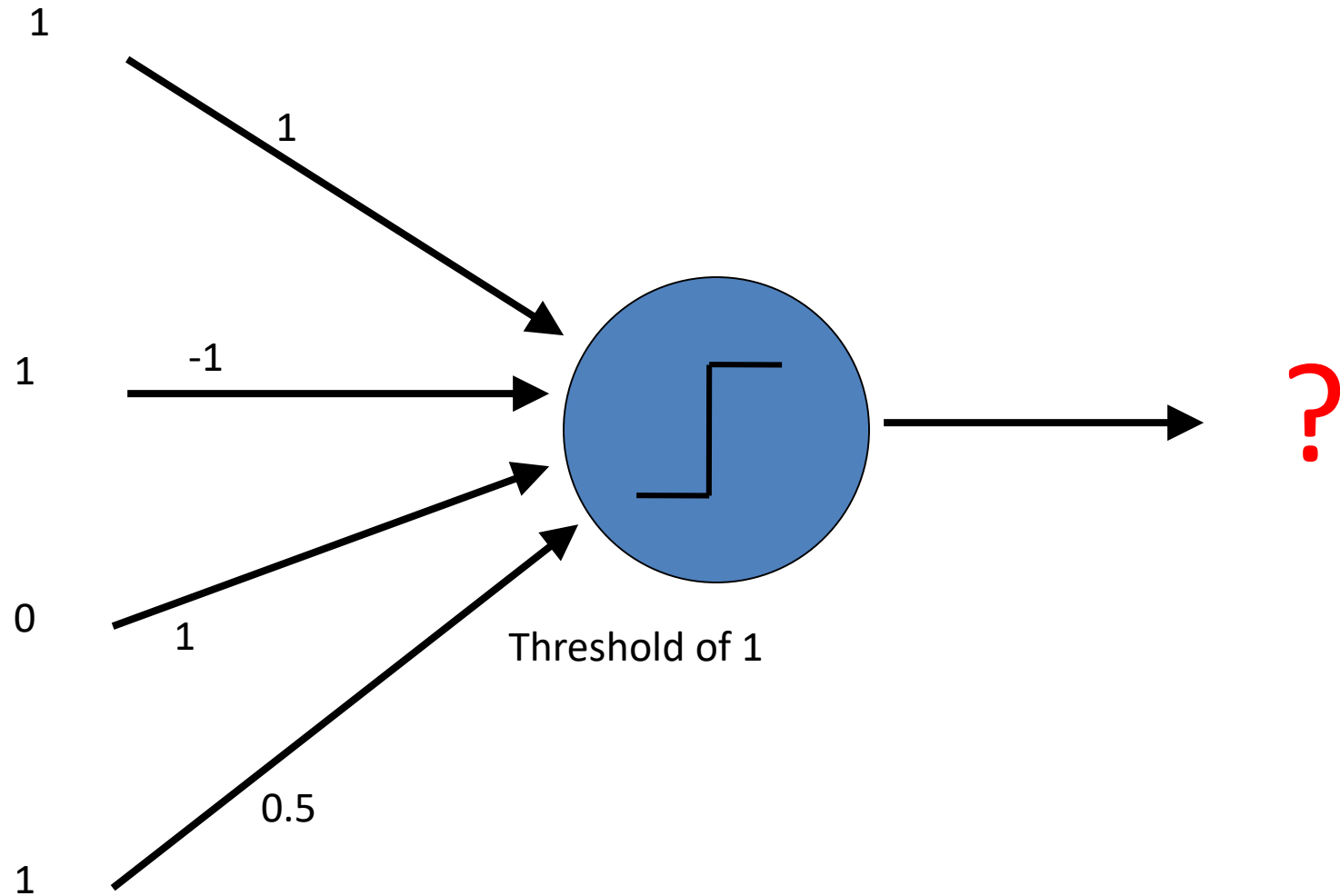


Sigmoid

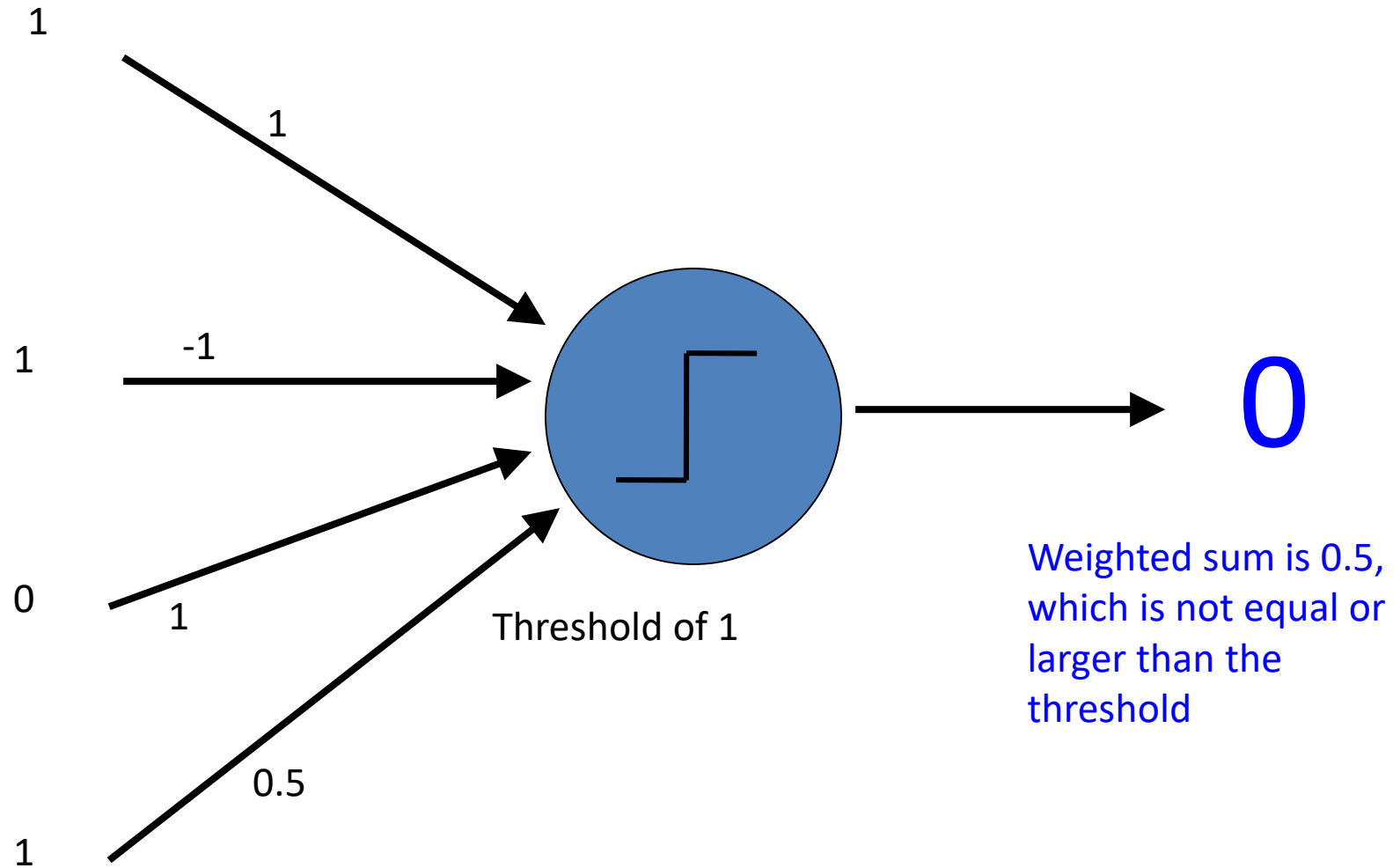
$$g(x) = \frac{1}{1 + e^{-ax}}$$



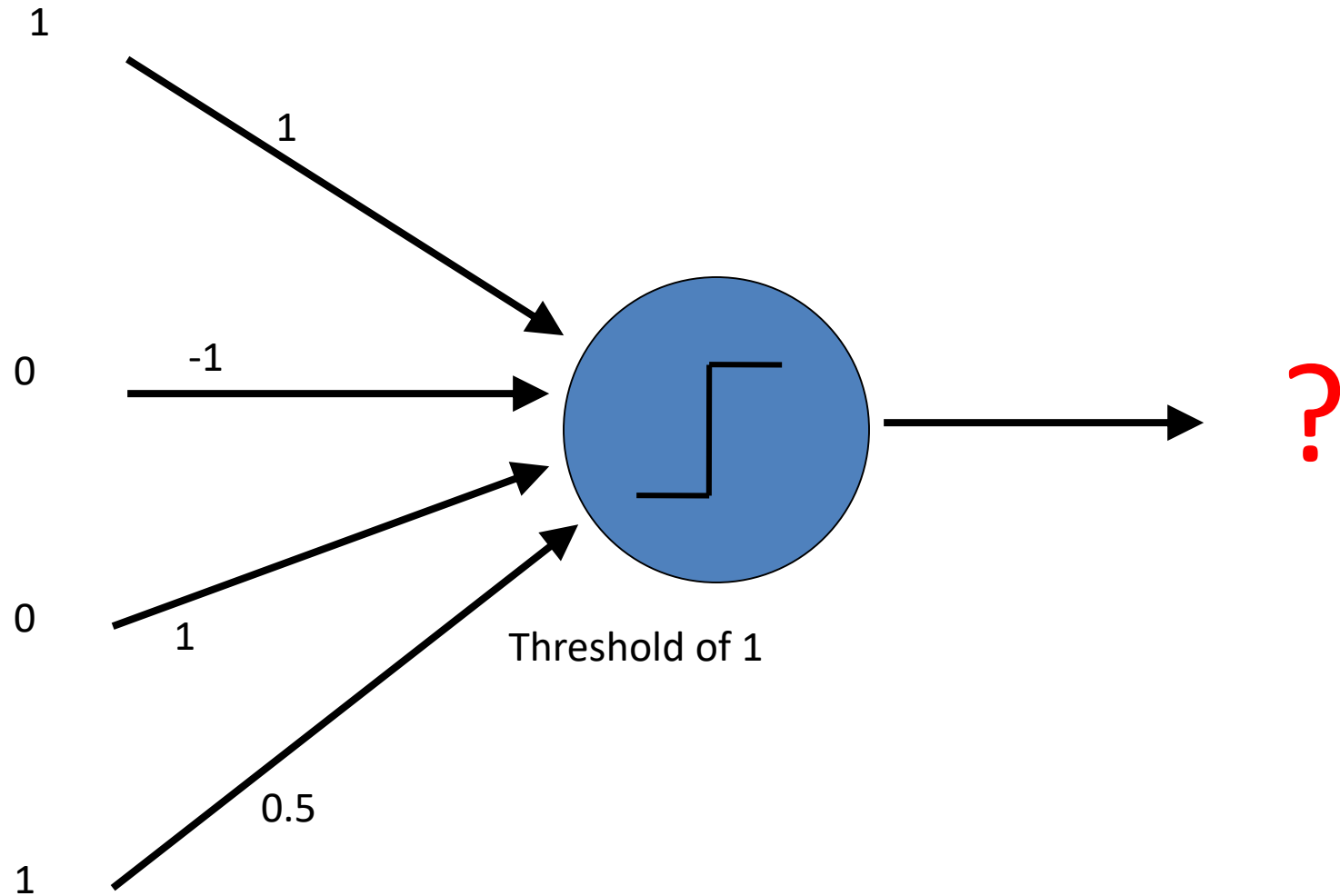
# A Single Neuron/Perceptron



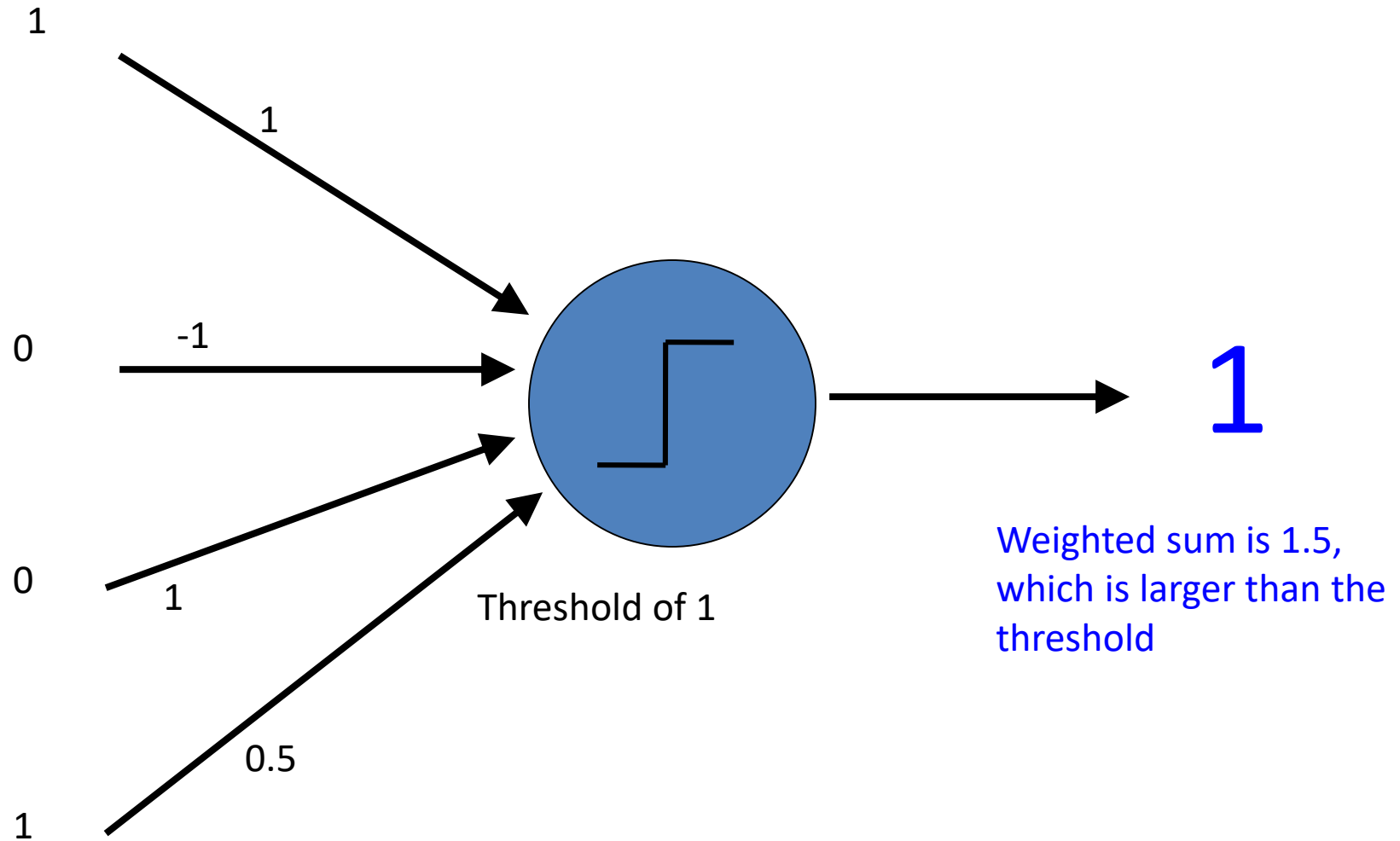
# A Single Neuron/Perceptron



# A Single Neuron/Perceptron



# A Single Neuron/Perceptron



# A Single Neuron/Perceptron

