

# Perceptron

# Agenda

- Review of models
- Parts of a perceptron
- Making decisions with a learned model
- Learning a perceptron representation
- Batch and Online learning
- Sparse and Dense representations

# Learning from examples

- Decision Trees
  - learn nested **if-then-else** rules
- kNN
  - memorize examples
- Linear classifiers
  - find simple boundaries in space

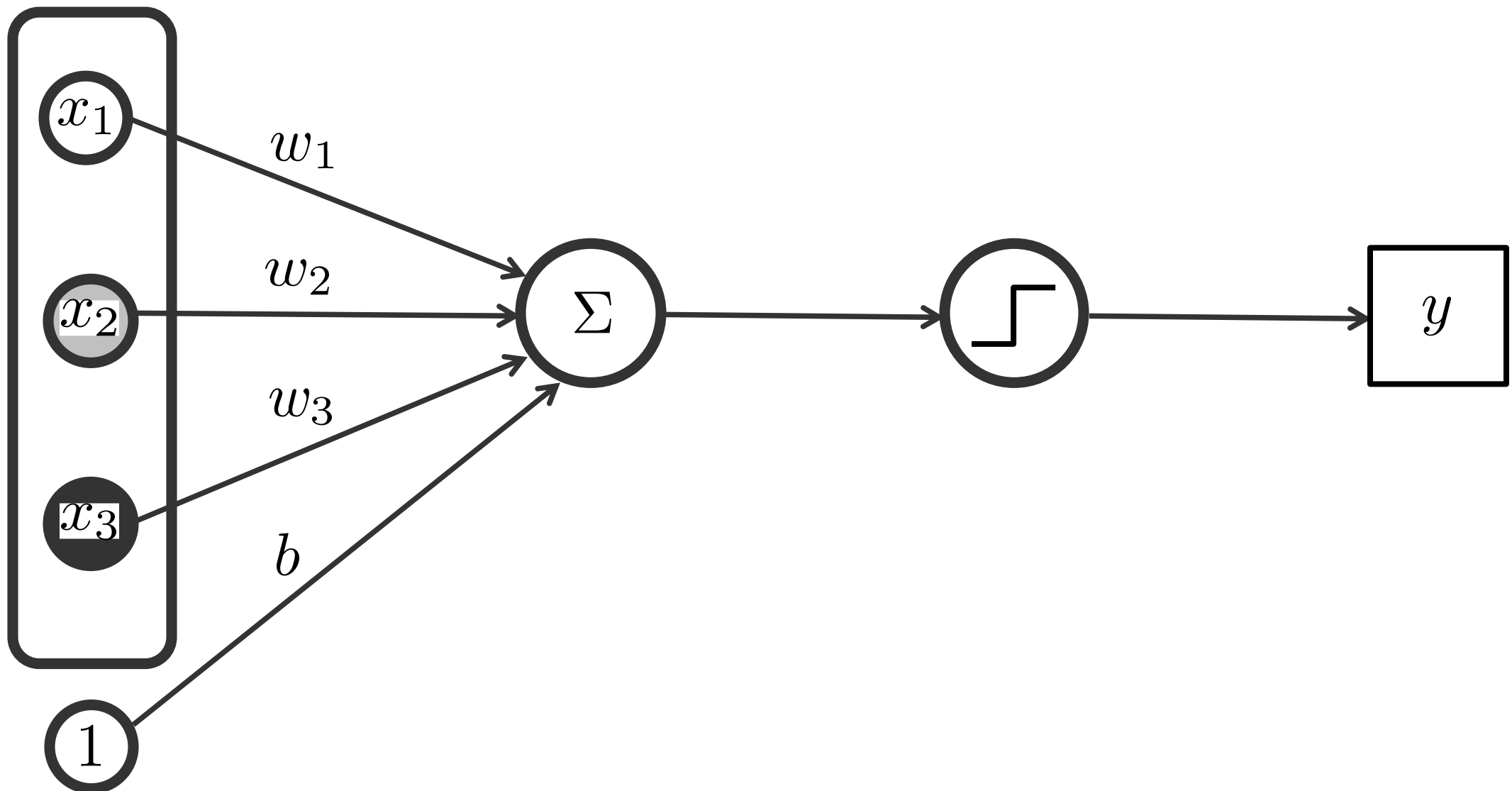


Frank Rosenblatt 1928 – 1971. Psychologist, inventor of the perceptron algorithm.

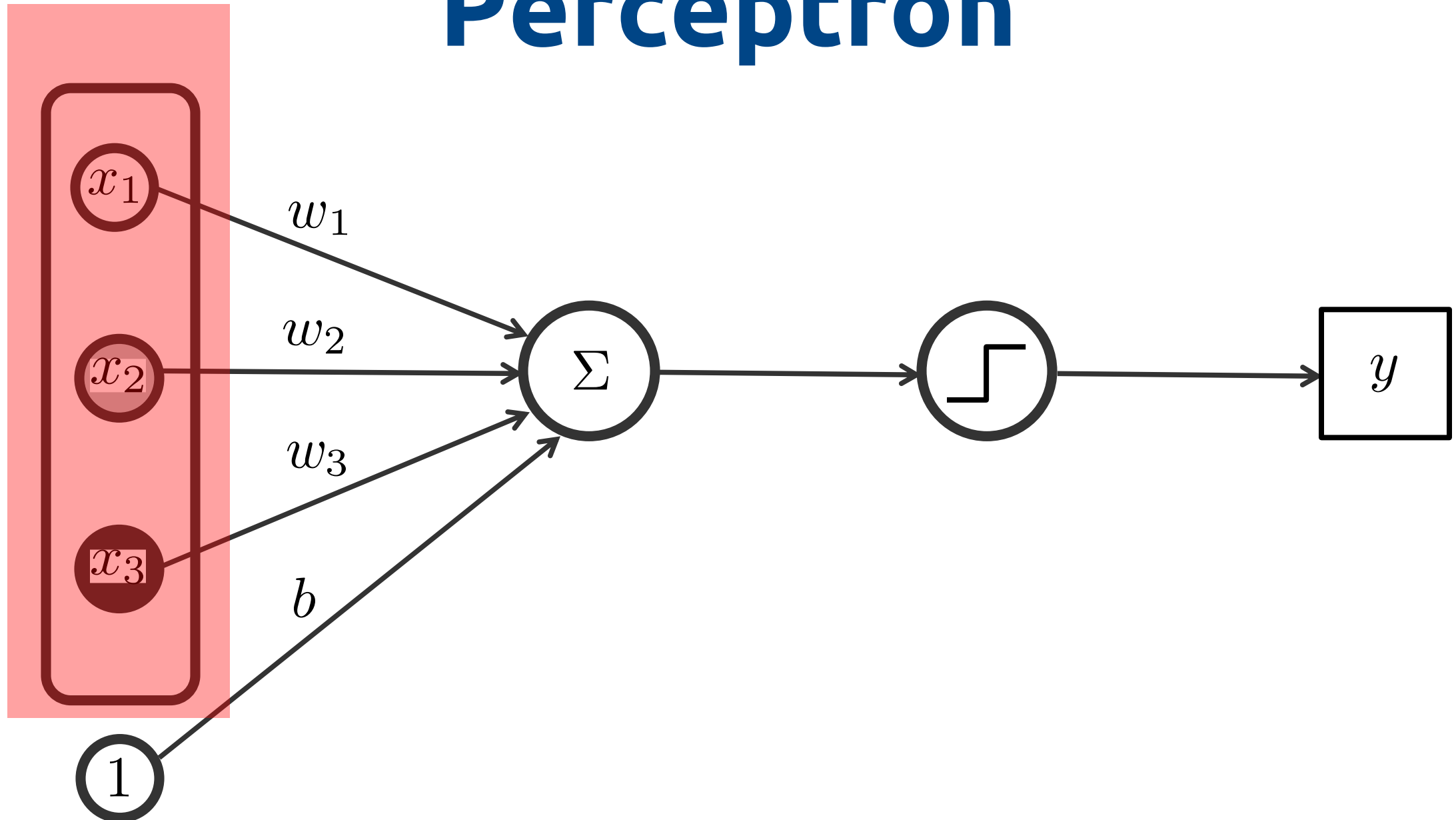
# Agenda

- Review of models
- Parts of a perceptron
- Making decisions with a learned model
- Learning a perceptron representation
- Batch and Online learning
- Sparse and Dense representations

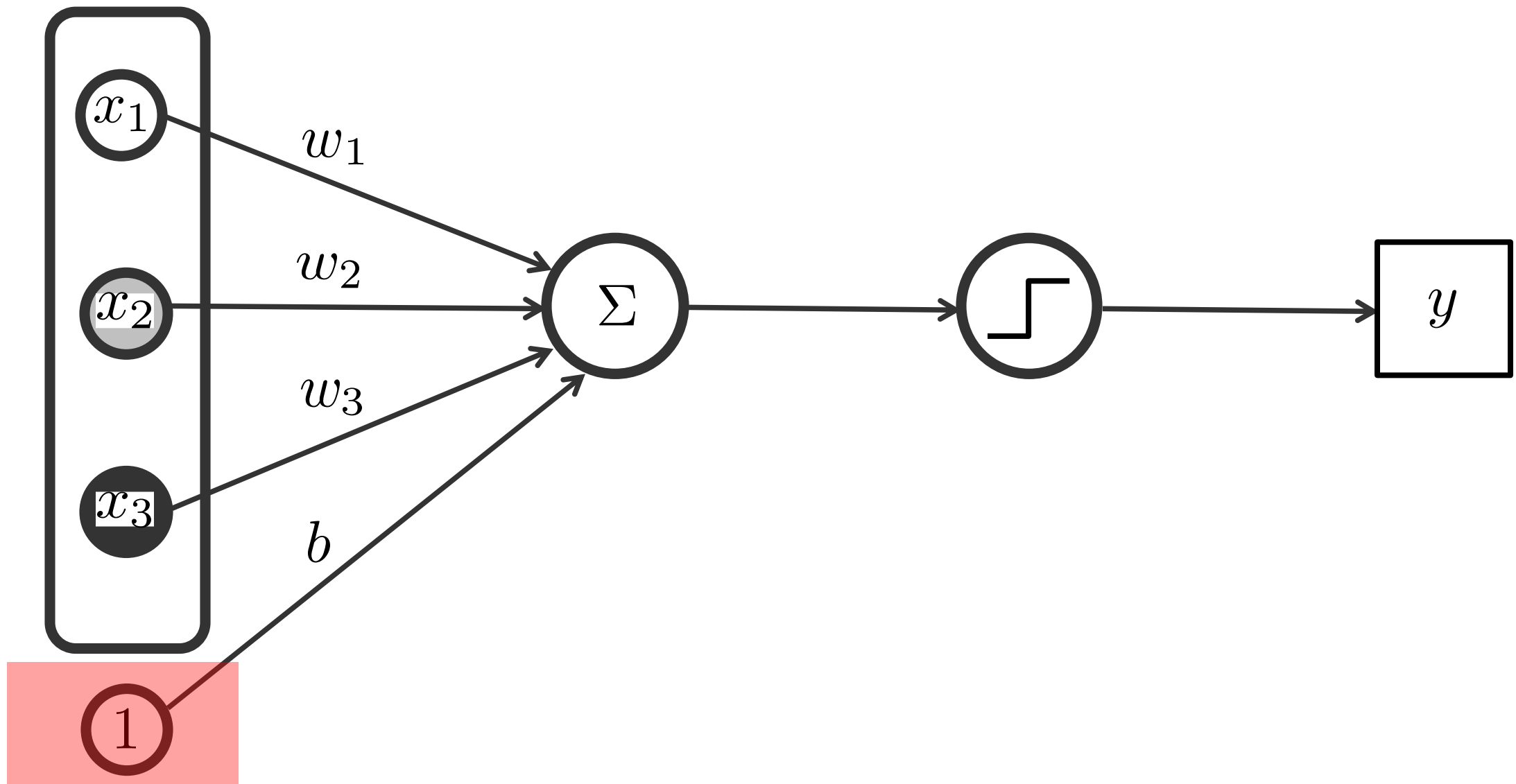
# Perceptron



# Perceptron

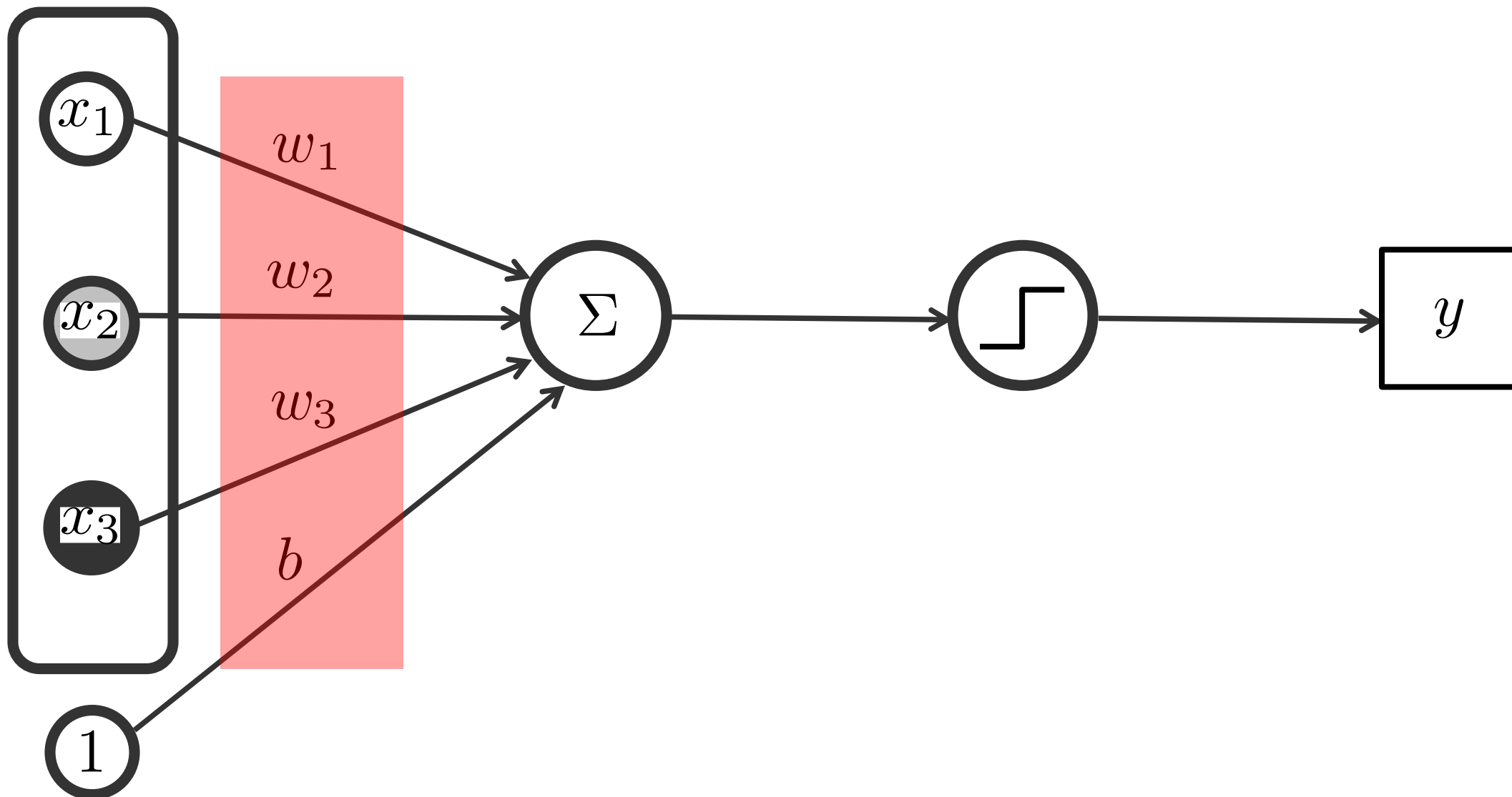


# Perceptron

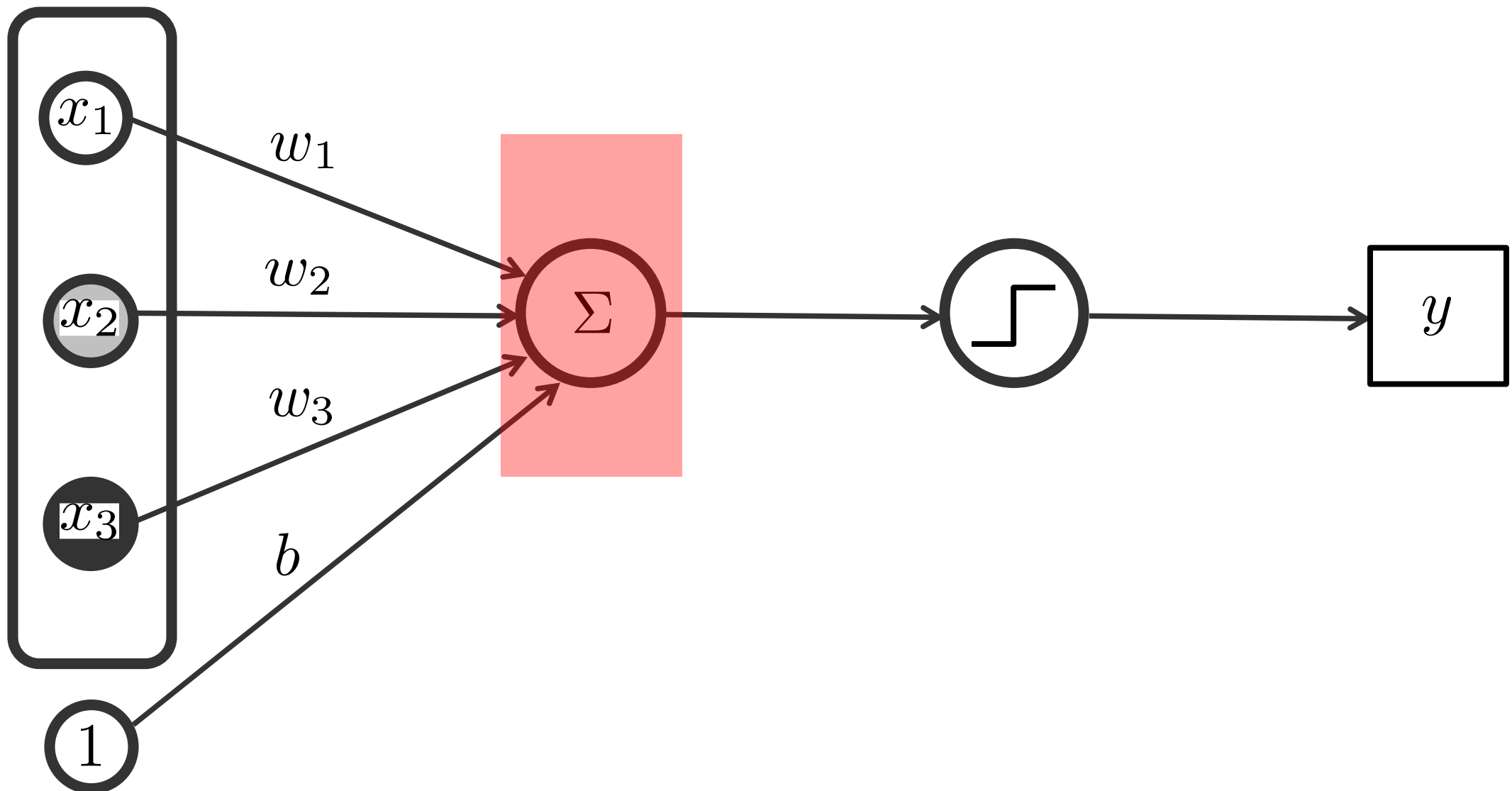




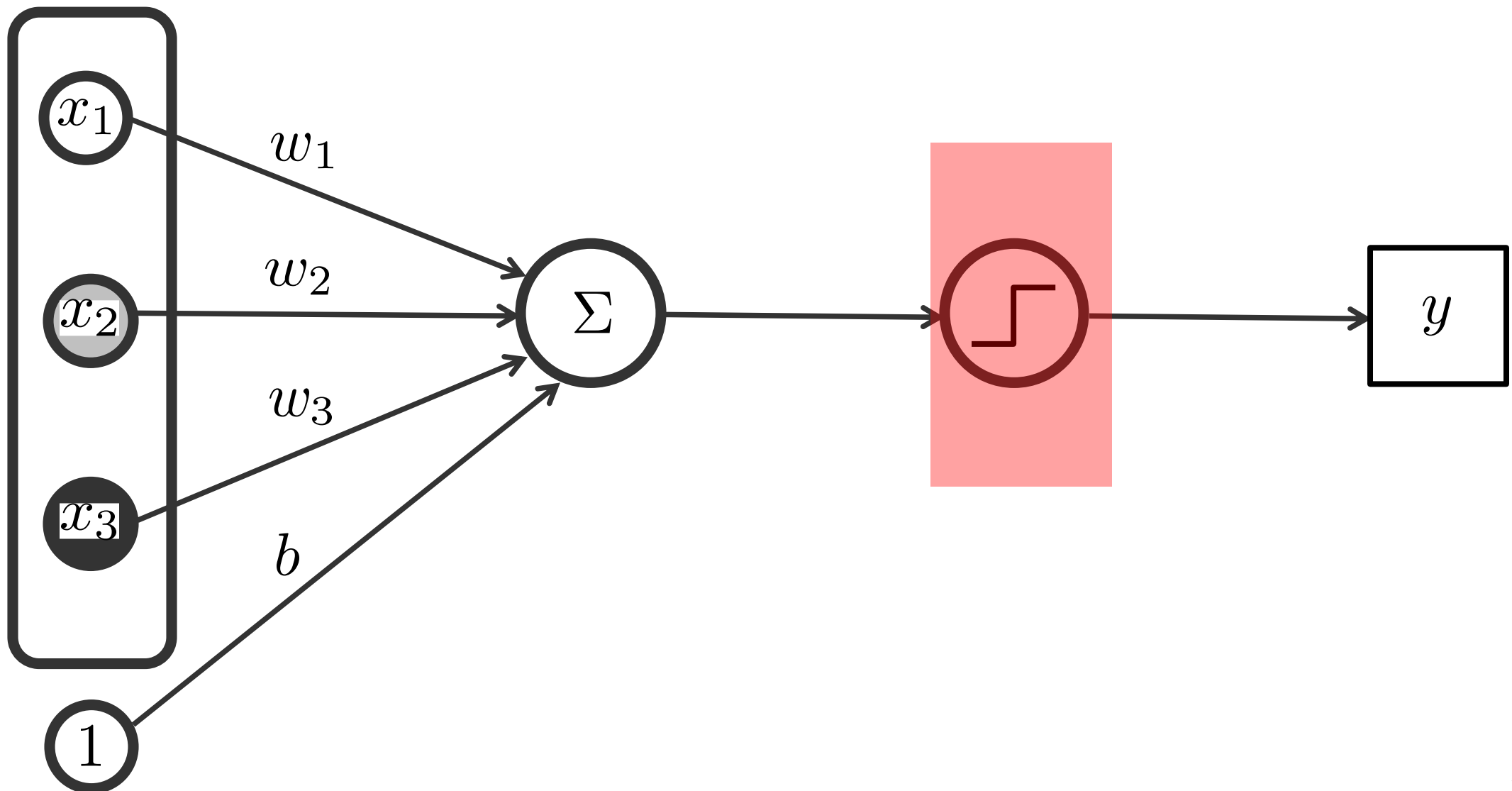
# Perceptron



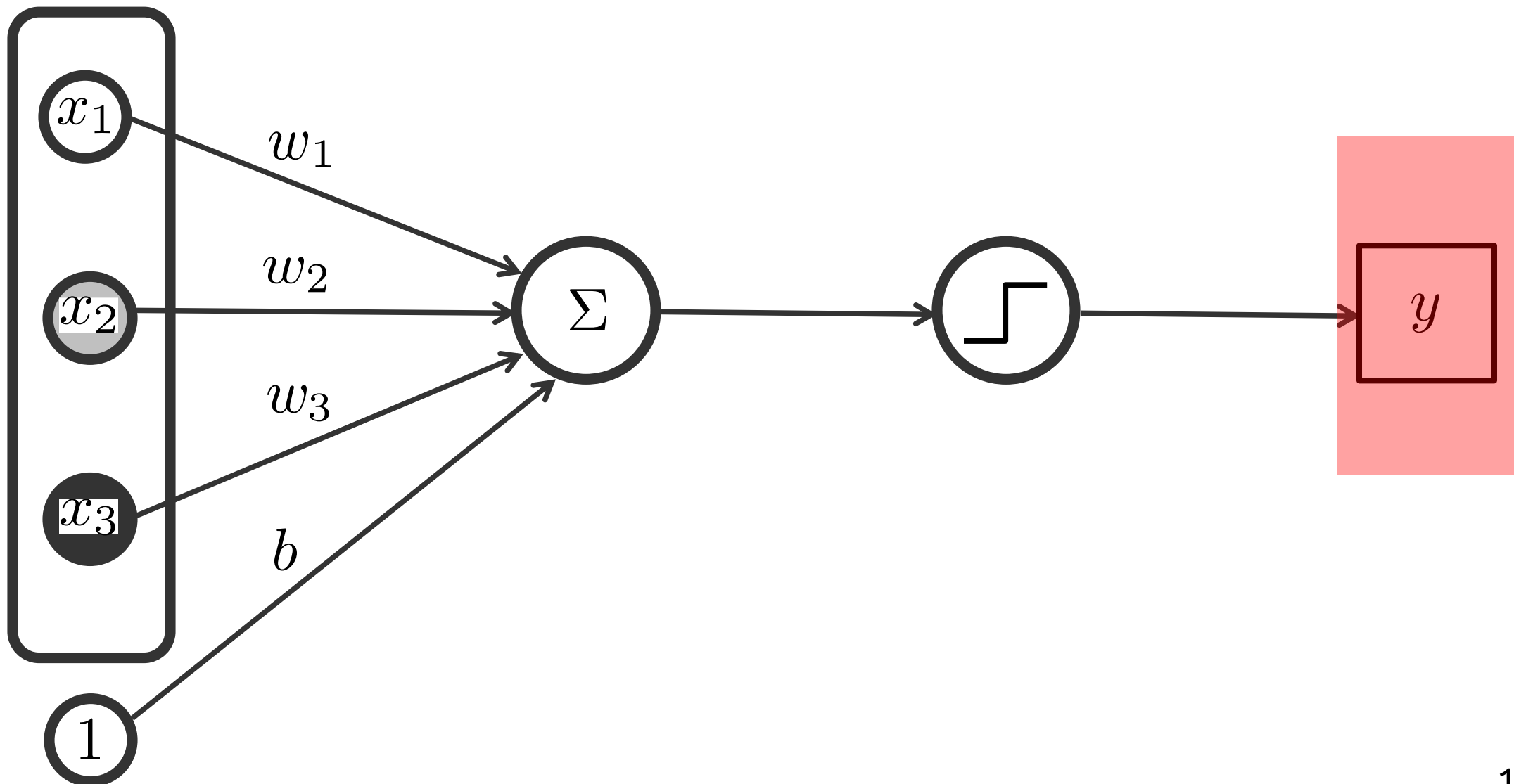
# Perceptron



# Perceptron



# Perceptron



# Agenda

- Review of models
- Parts of a perceptron
- Making decisions with a learned model
- Learning a perceptron representation
- Batch and Online learning
- Sparse and Dense representations

# Perceptron classification rule

- Perceptron uses a simple rule to classify objects
- It computes the weighted sum of the input features (plus **bias**)
- If this sum is greater than or equal to **0**, it outputs positive class **+1**
- Otherwise it outputs negative class **-1**

# Discriminant function

$$f(\mathbf{x}) = \left( \sum_{i=1}^N w_i x_i \right) + b$$

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

# Example: movie reviews

	#good	#dark	#mediocre	#the
$\mathbf{x}^1 =$	( 2,	0,	0,	5 )

$\mathbf{w} =$	( 2.5,	0.5,	-4.0,	0.0 )
----------------	--------	------	-------	-------

$$b = 0.5$$

$$\text{score } f(\mathbf{x}^1) =$$



# Role of bias

- When  $w \cdot x \approx 0$ ,  
bias decides which class to predict
- Makes the default decision
- Biases the classifier towards positive or negative class

# Geometric interpretation in 2D

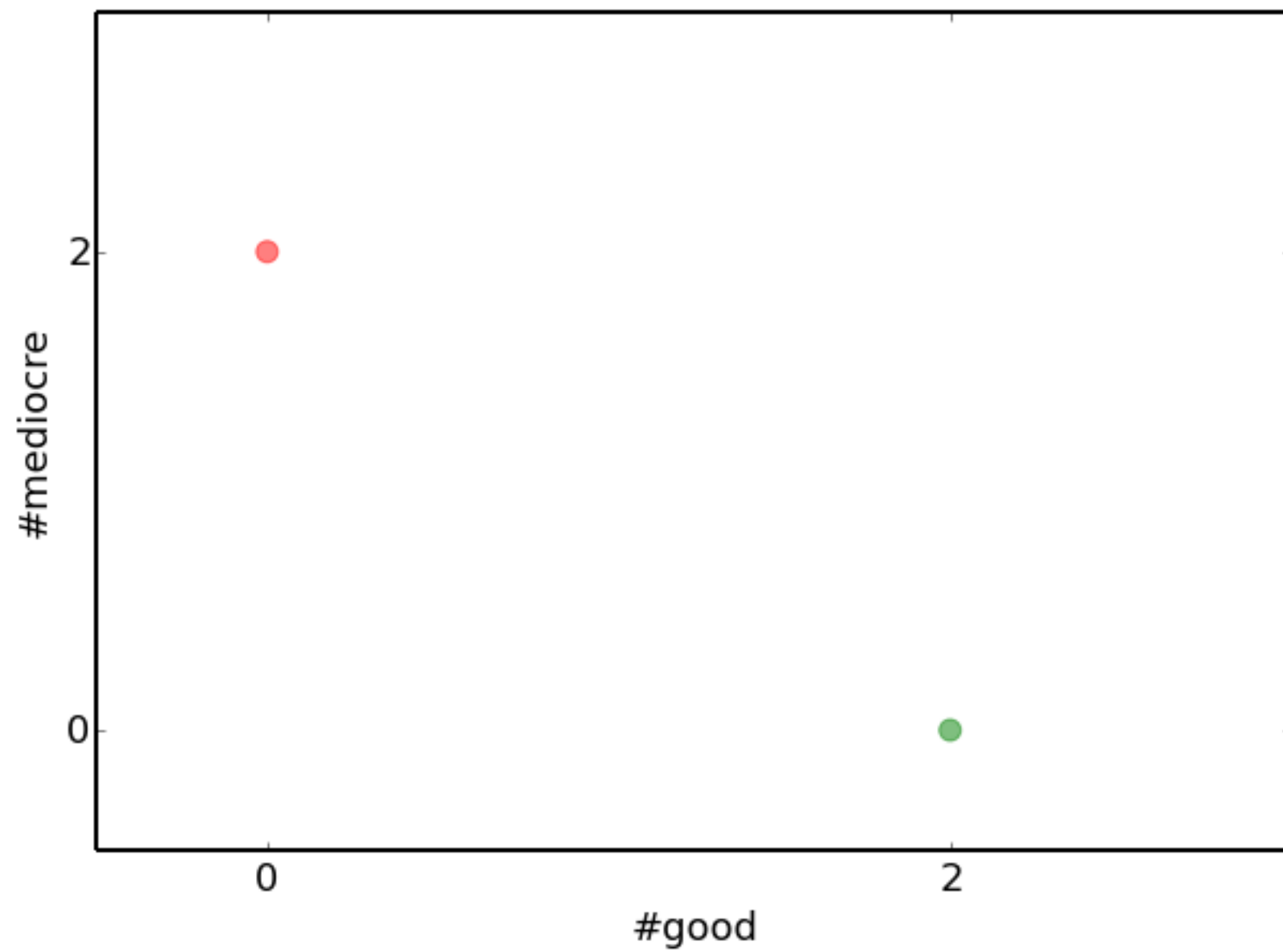
#good      #mediocre

$$\mathbf{x}^1 = (2, 0)$$

$$\mathbf{x}^2 = (0, 2)$$

$$\mathbf{w} = (2.5, -4.0)$$

$$b = 0.5$$



# Decision boundary

$$w_1x_1 + w_2x_2 + b = 0$$

Solve for  $x_2$

$$x_2 = -\frac{w_1}{w_2}x_1 - \frac{b}{w_2}$$

slope

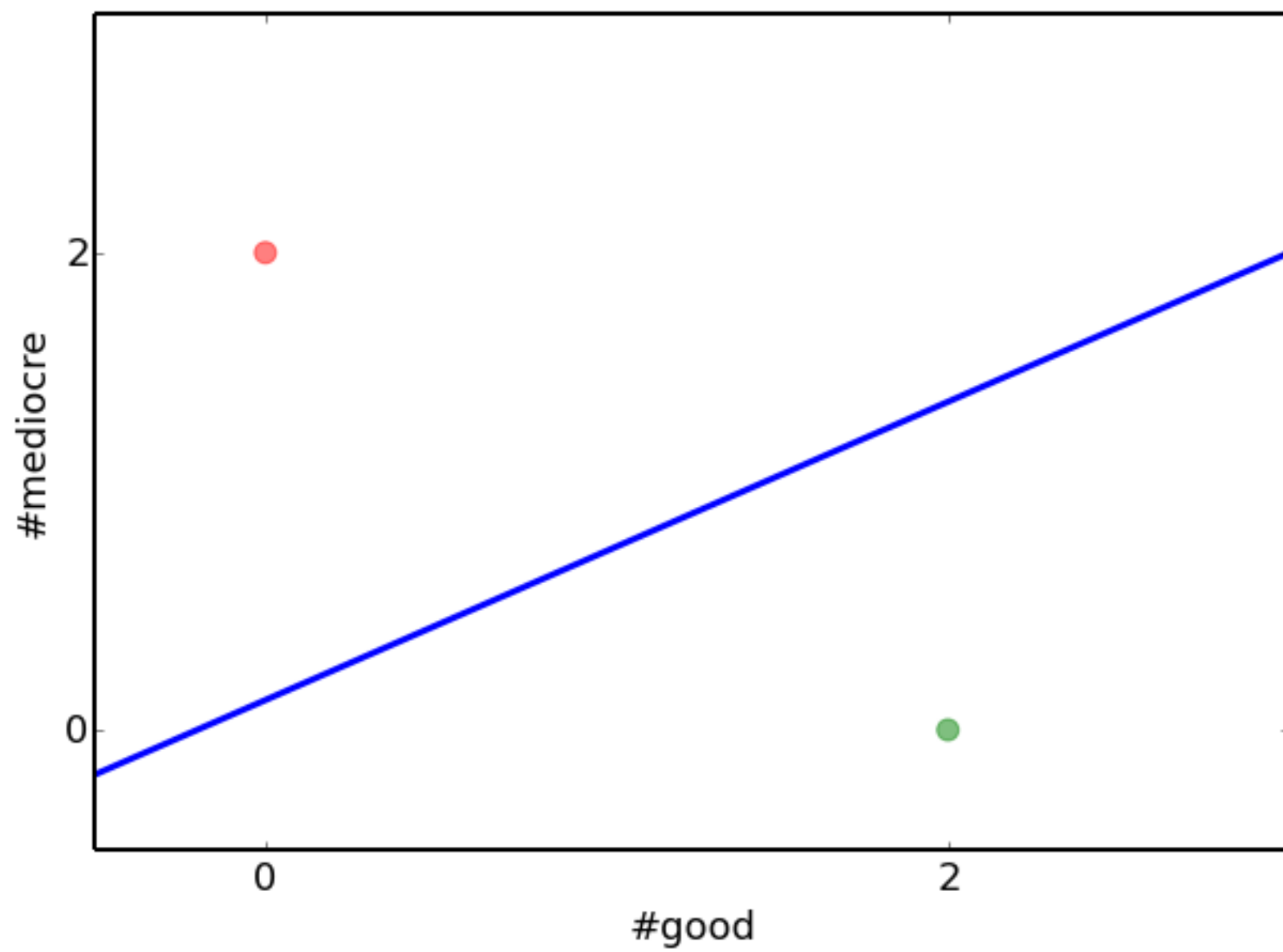
intercept

# Aside: Multiple dimensions

- Examples with only 2 dimensions are not very realistic
  - What about domains with more than two features?
- 2 dimensions are “split” by a 1 dimensional shape: a line
- 3 dimensional space is “split” by a 2 dimensional shape: a plane
- $N$  dimensional space is “split” by an  $(N-1)$  dimensional shape: a hyperplane

# Agenda

- Review of models
- Parts of a perceptron
- Making decisions with a learned model
- Learning a perceptron representation
- Batch and Online learning
- Sparse and Dense representations



# How can we find good $(w,b)$ ?

- Go through examples one by one
- Try classifying current example with current  $(w,b)$
- If correct, keep going
- If not correct, adjust  $(w,b)$



- New item:  $(\mathbf{x}, +1)$
- With current  $(\mathbf{w}, b)$ , the score  $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$  is less than 0

# How to adjust $(w, b)$ ?

- New item:  $(x, +1)$
- With current  $(w, b)$ , the score  $f(x) = w \cdot x + b$  is less than 0

# How to adjust $(w, b)$ ?

- New item:  $(x, +1)$
- With current  $(w, b)$ , the score  $f(x) = w \cdot x + b$  is less than 0
- How do we change  $b$  to make it higher?
- How do we change  $w$  to make it higher?

# A Concrete Example

$$\mathbf{x}^1 = ( \quad 2, \quad \quad 0, \quad \quad 0, \quad \quad 5 \quad )$$

$$\mathbf{w} = ( -0.5, \quad 1.0, \quad -2.0, \quad \quad 0.0 \quad )$$

$$\mathbf{b} = 0$$

$$f(\mathbf{x}^1) = \mathbf{w} \cdot \mathbf{x}^1 + \mathbf{b} = -1.0$$

Change  $\mathbf{b}$  to increase  $f(\mathbf{x}^1)$

Change  $\mathbf{w}$  to increase  $f(\mathbf{x}^1)$

# Vector addition and subtraction

$$\mathbf{c} = \mathbf{a} + \mathbf{b}$$

for all  $i$ ,  $c_i = a_i + b_i$

$$\begin{array}{lcl} \mathbf{x}^1 & = & ( \quad 2, \quad \quad 0, \quad \quad 0, \quad \quad 5 \ ) \\ \mathbf{w} & = & ( -0.5, \quad 1.0, \quad -2.0, \quad 0.0 \ ) \\ \mathbf{w} + \mathbf{x}^1 & = & ( \quad 1.5, \quad 1.0, \quad -2.0, \quad 5.0 \ ) \end{array}$$

# Reminder:

# Discriminant function

$$f(\mathbf{x}) = \left( \sum_{i=1}^N w_i x_i \right) + b$$

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

# Update rule: example $(\mathbf{x}, y)$ , model $(\mathbf{w}, b)$

$$1: y_{\text{pred}} = \text{predict}((\mathbf{w}, b), \mathbf{x})$$

# Update rule: example $(\mathbf{x}, y)$ , model $(\mathbf{w}, b)$

```
1:  $y_{\text{pred}} = \text{predict}((\mathbf{w}, b), \mathbf{x})$   
2: if  $y = +1$  and  $y_{\text{pred}} = -1$  then  
3:    $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}$   
4:    $b \leftarrow b + 1$   
5: else if  $y = -1$  and  $y_{\text{pred}} = +1$  then  
6:    $\mathbf{w} \leftarrow \mathbf{w} - \mathbf{x}$   
7:    $b \leftarrow b - 1$ 
```



# Aside: Dot product notation

$$\mathbf{w} \cdot \mathbf{x} = \sum_{i=1}^N w_i x_i$$

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$$

# Update rule: example $(\mathbf{x}, y)$ , model $(\mathbf{w}, b)$

```
1:  $y_{\text{pred}} = \text{predict}((\mathbf{w}, b), \mathbf{x})$   
2: if  $y = +1$  and  $y_{\text{pred}} = -1$  then  
3:    $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}$   
4:    $b \leftarrow b + 1$   
5: else if  $y = -1$  and  $y_{\text{pred}} = +1$  then  
6:    $\mathbf{w} \leftarrow \mathbf{w} - \mathbf{x}$   
7:    $b \leftarrow b - 1$ 
```

# Proof Step 3 increases predicted score

New weight vector

$$(\mathbf{w} + \mathbf{x}) \cdot \mathbf{x} = \mathbf{w} \cdot \mathbf{x} + \mathbf{x} \cdot \mathbf{x}$$

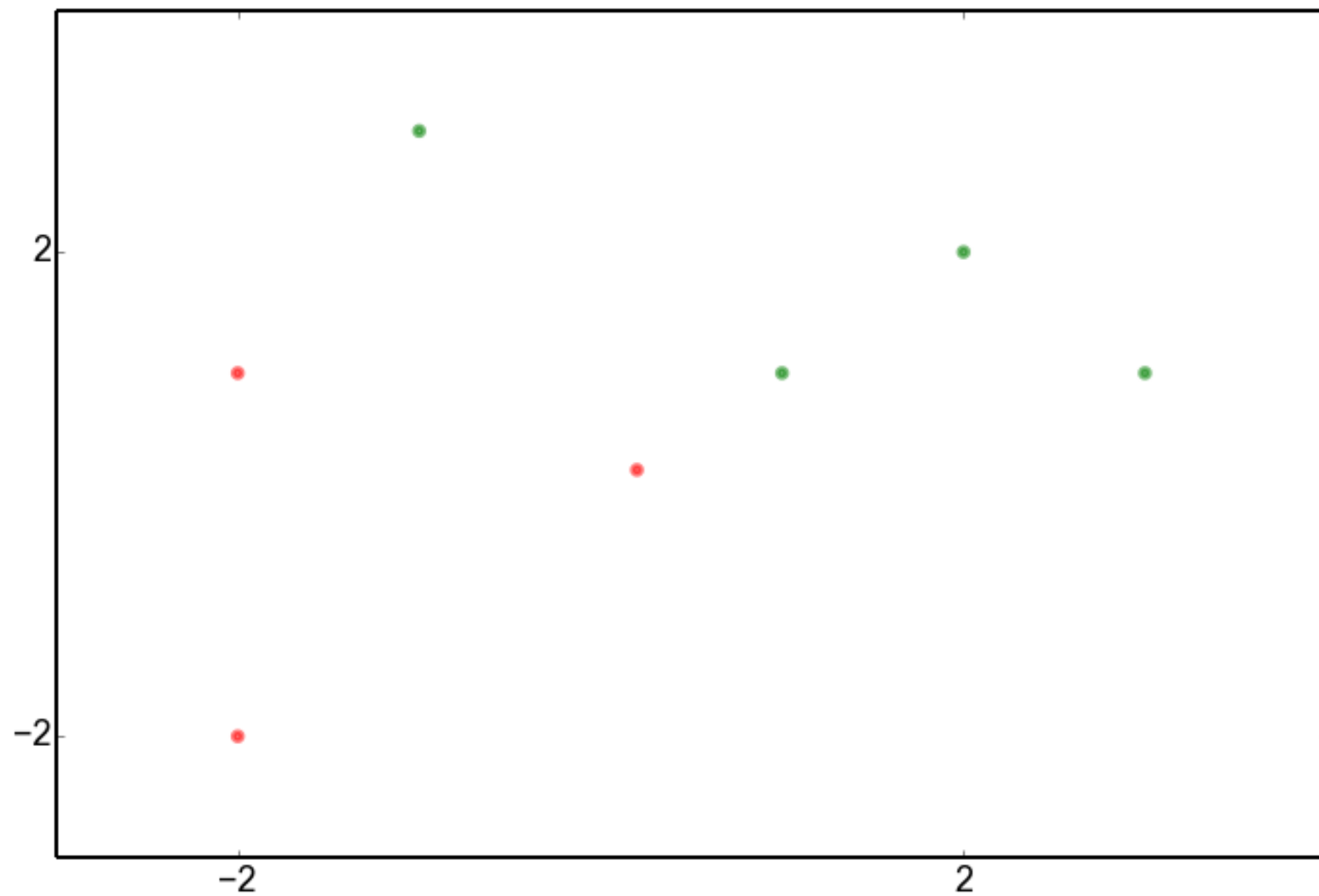
Since  $\mathbf{x} \cdot \mathbf{x} \geq 0$

Thus  $\mathbf{w} \cdot \mathbf{x} + \mathbf{x} \cdot \mathbf{x} \geq \mathbf{w} \cdot \mathbf{x}$

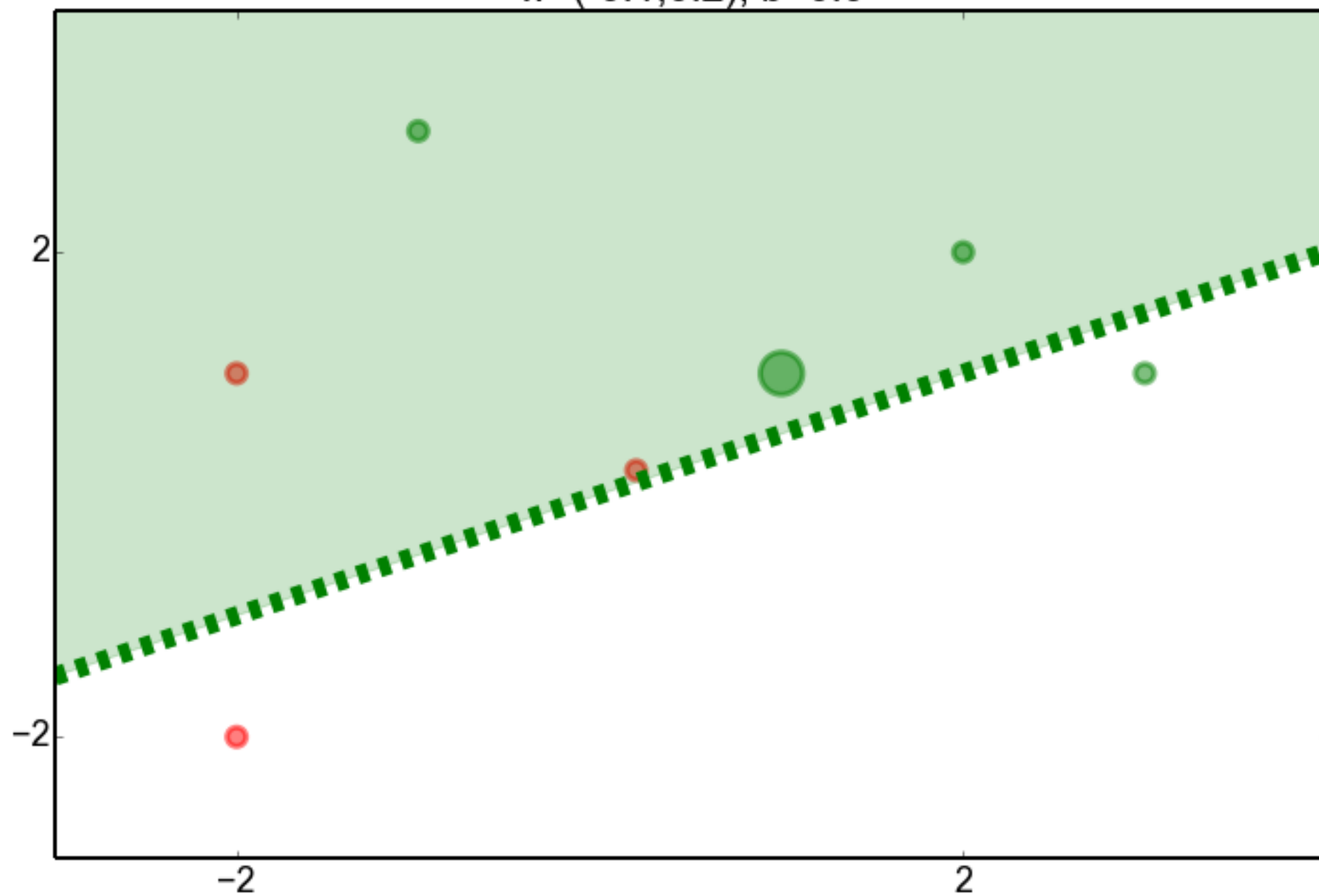
# One learning iteration over $N$ examples

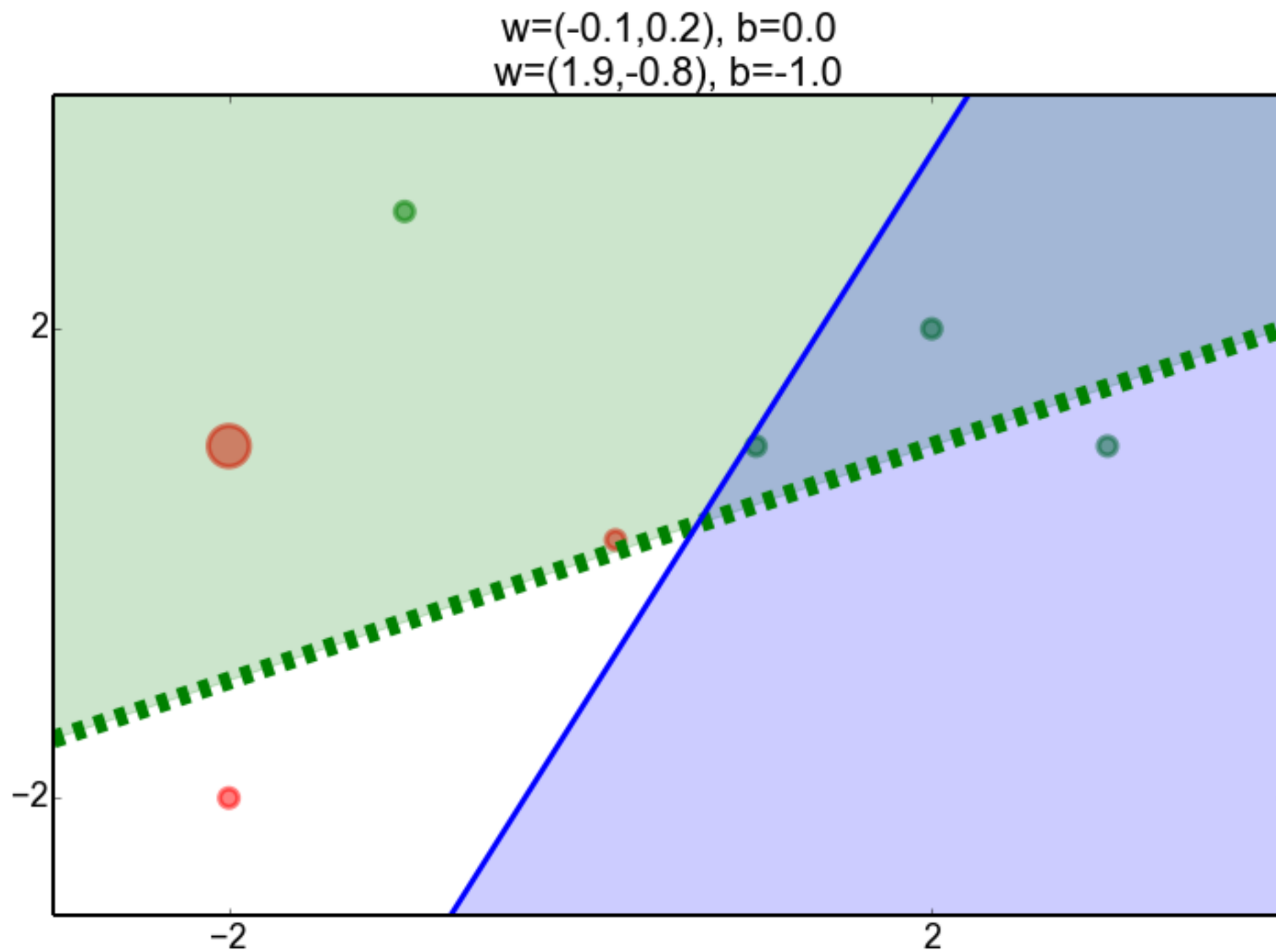
```
1:  $\mathbf{w} \leftarrow \mathbf{0}$   
2:  $b \leftarrow 0$   
3: for  $n = 1..N$  do  
4:    $y_{\text{pred}}^n = \text{predict}((\mathbf{w}, b), \mathbf{x}^n)$   
5:   if  $y^n = +1$  and  $y_{\text{pred}}^n = -1$  then  
6:      $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}^n$   
7:      $b \leftarrow b + 1$   
8:   else if  $y^n = -1$  and  $y_{\text{pred}}^n = +1$  then  
9:      $\mathbf{w} \leftarrow \mathbf{w} - \mathbf{x}^n$   
10:     $b \leftarrow b - 1$ 
```

# Example



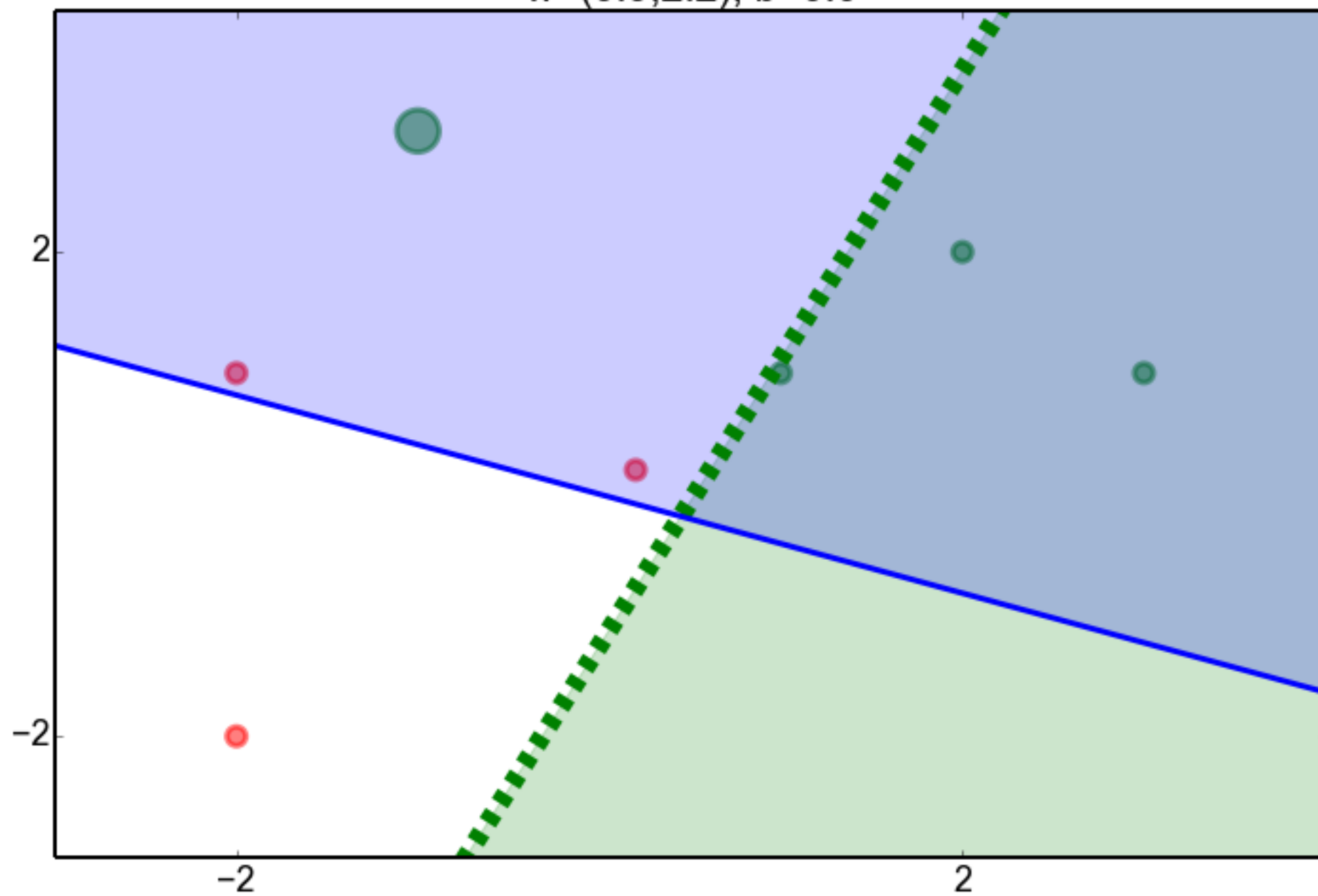
$w=(-0.1,0.2), b=0.0$   
 $w=(-0.1,0.2), b=0.0$



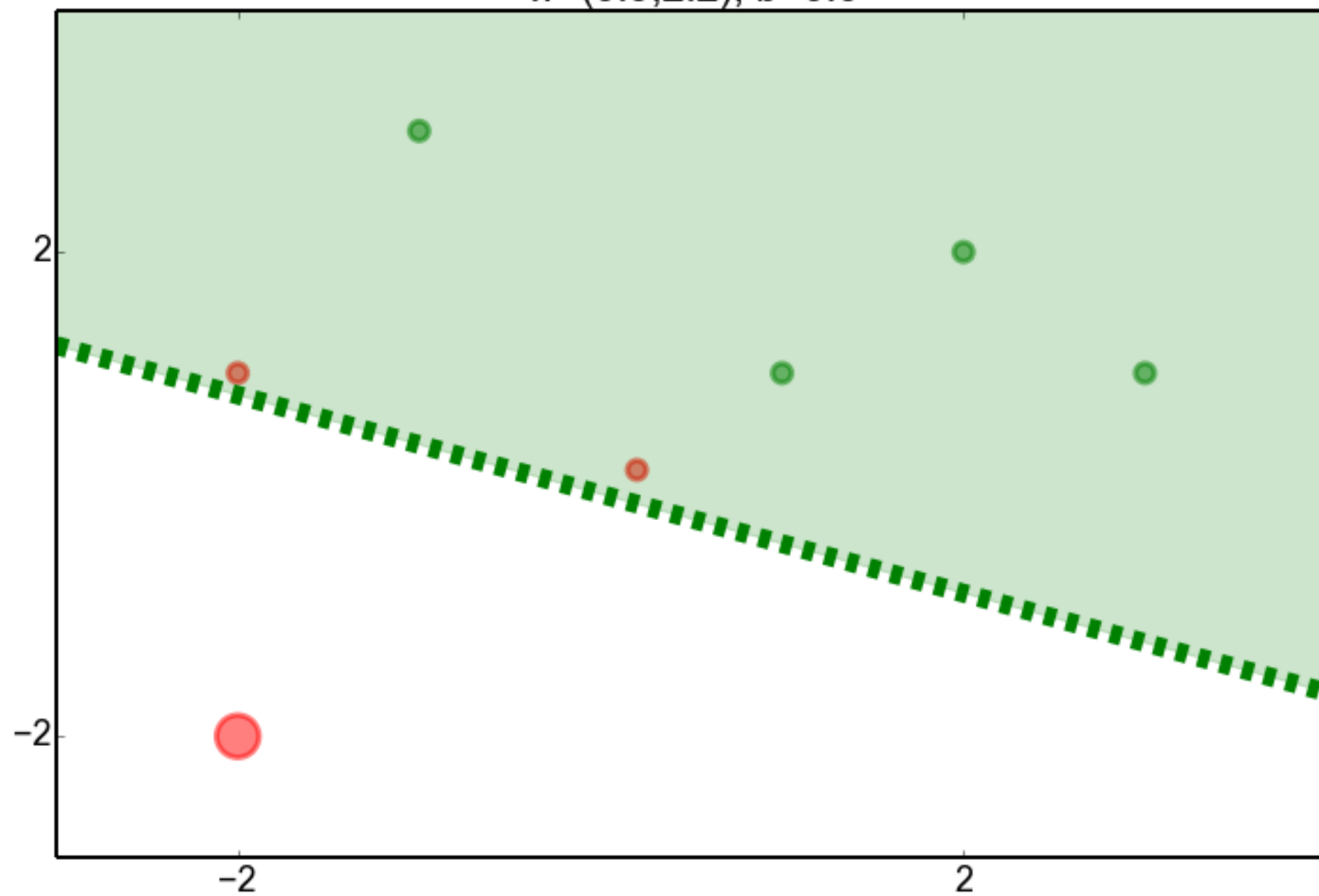




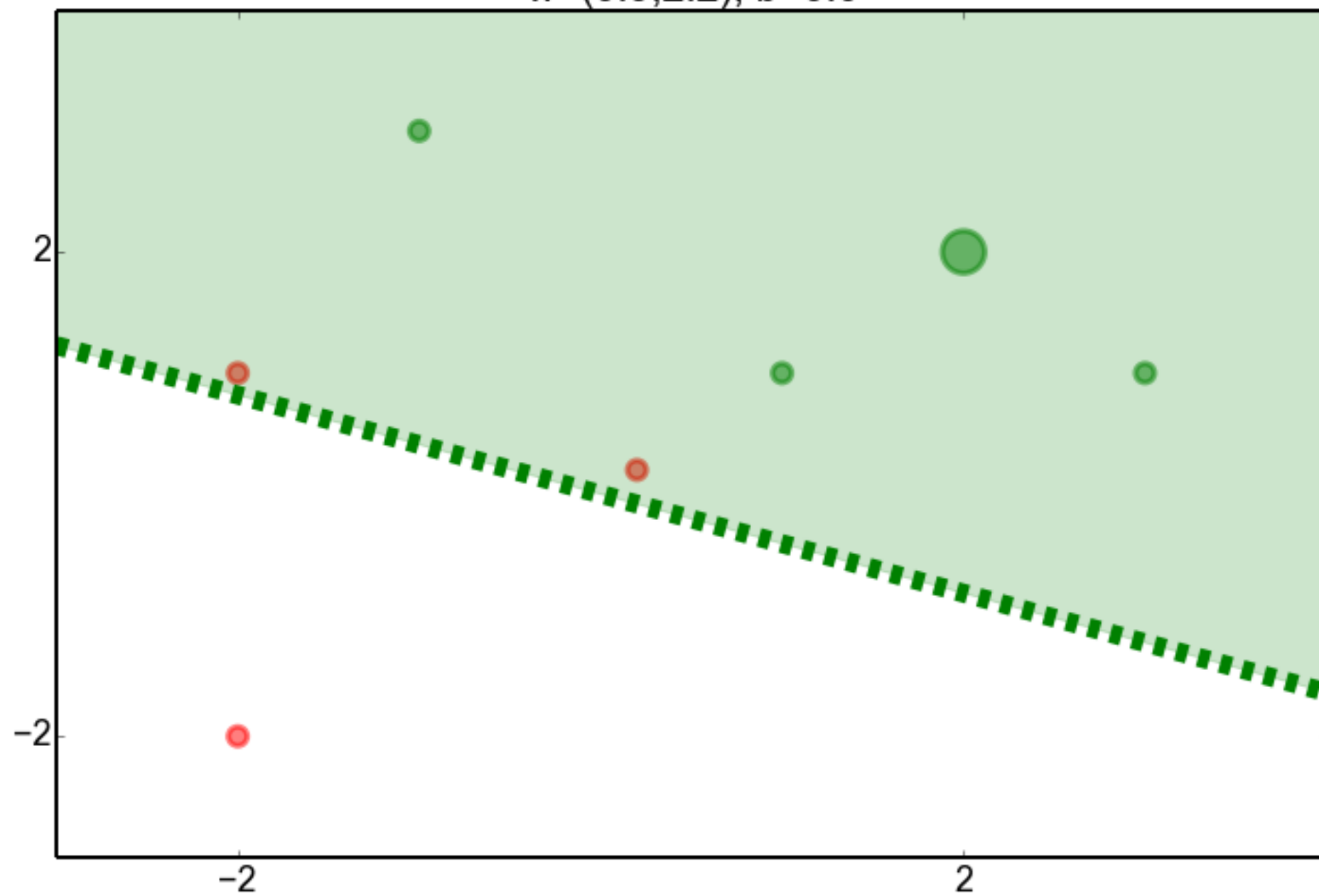
$w=(1.9,-0.8), b=-1.0$   
 $w=(0.9,2.2), b=0.0$



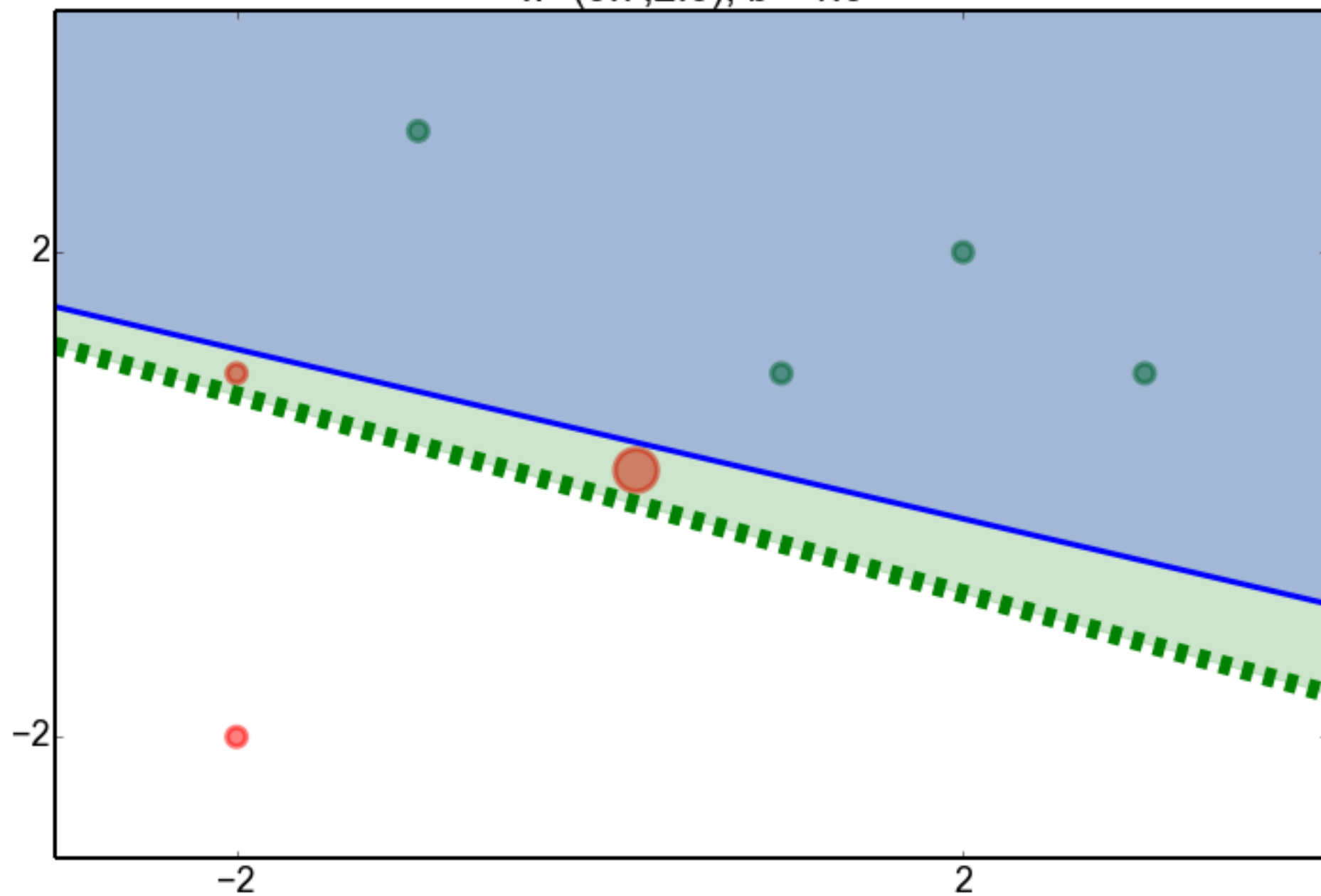
$w=(0.9,2.2)$ ,  $b=0.0$   
 $w=(0.9,2.2)$ ,  $b=0.0$



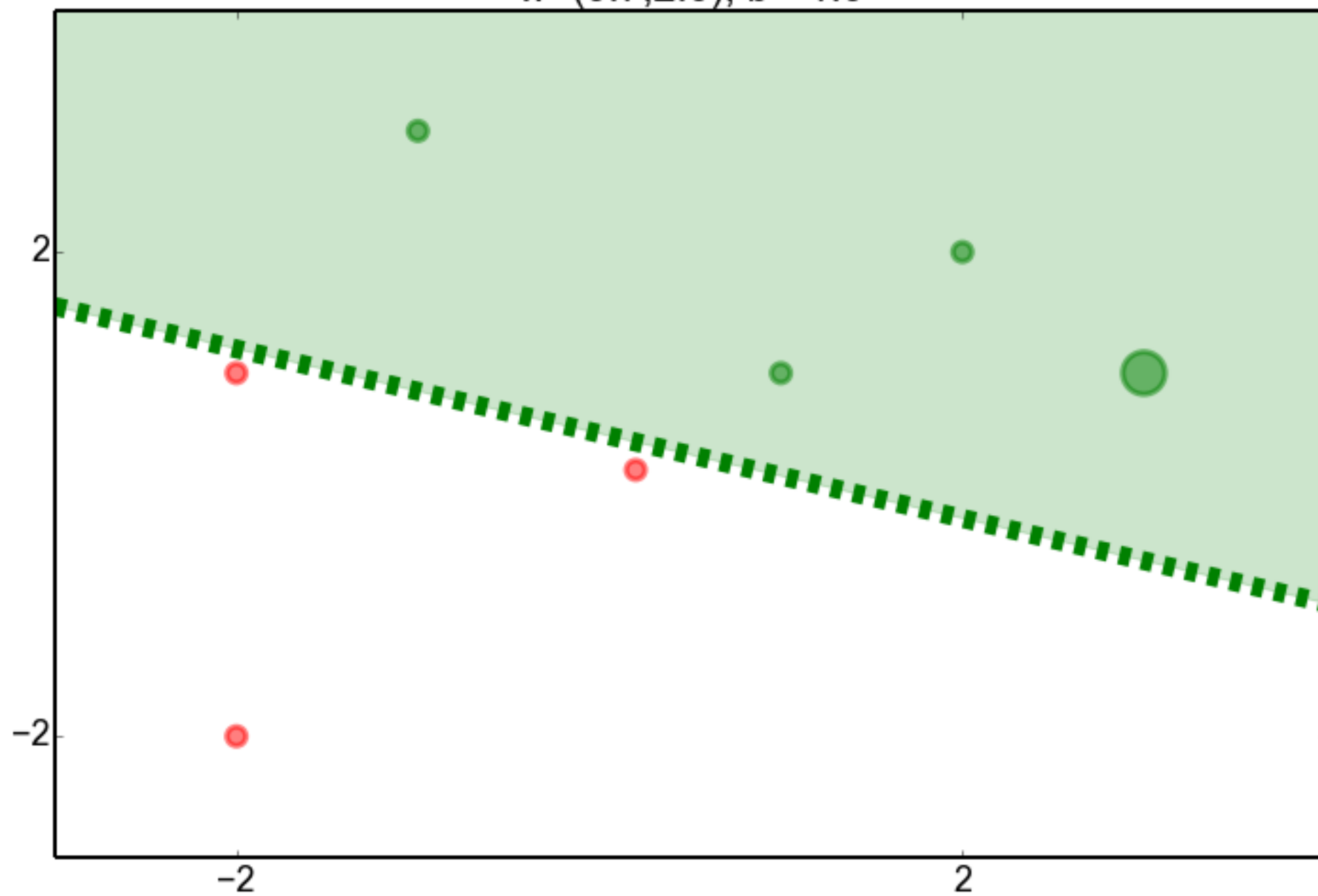
$w=(0.9,2.2)$ ,  $b=0.0$   
 $w=(0.9,2.2)$ ,  $b=0.0$



$w=(0.9,2.2), b=0.0$   
 $w=(0.7,2.0), b=-1.0$



$w=(0.7,2.0), b=-1.0$   
 $w=(0.7,2.0), b=-1.0$



# Termination

It can be proven that if there is a linear boundary separating  $+1$  from  $-1$ , the perceptron algorithm will find it.

# Agenda

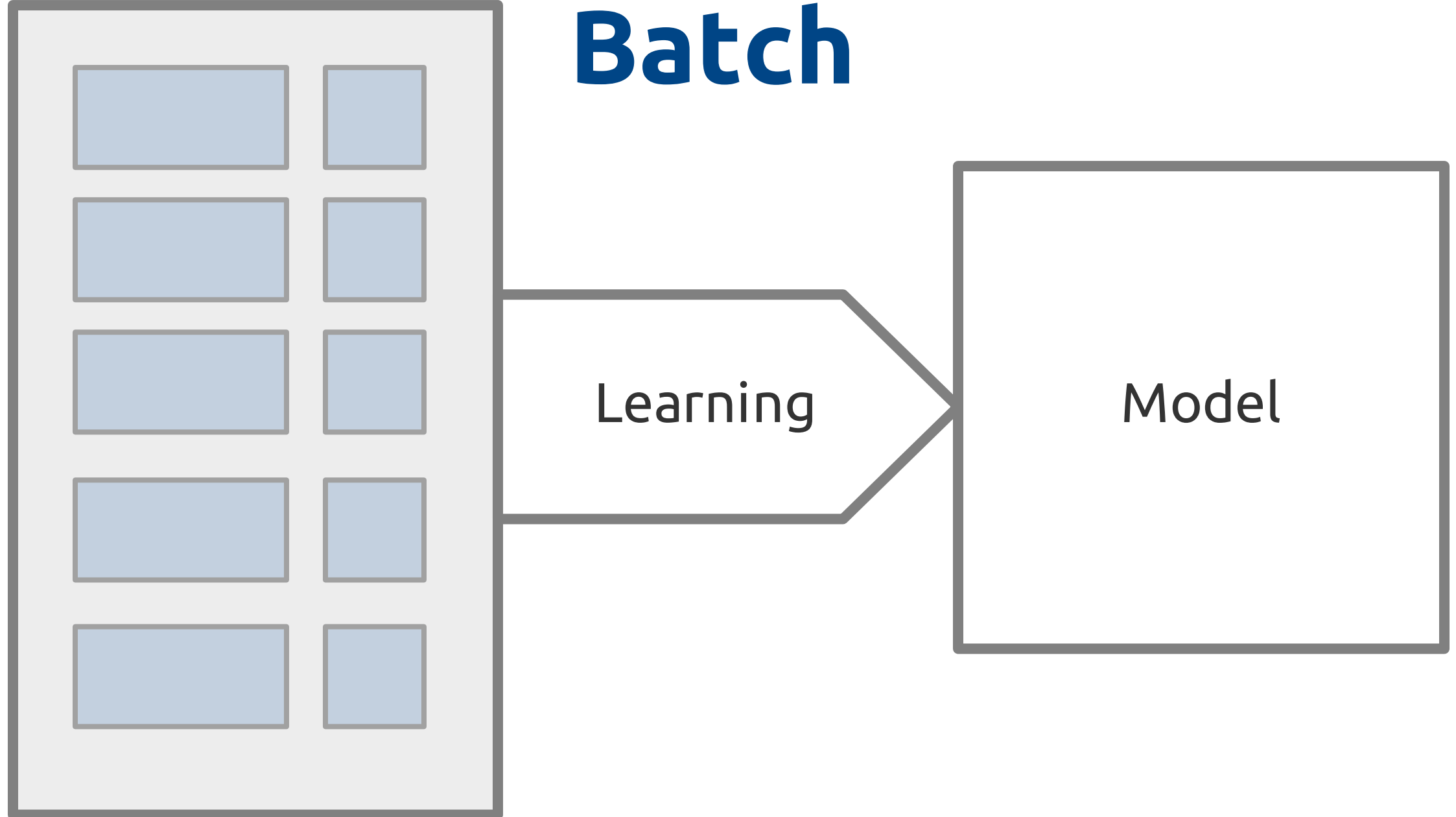
- Review of models
- Parts of a perceptron
- Making decisions with a learned model
- Learning a perceptron representation
- Batch and Online learning
- Sparse and Dense representations

# Streaming data

- Data which doesn't stop coming
  - Recordings from sensors
  - Posts to social media
  - News articles



# Batch

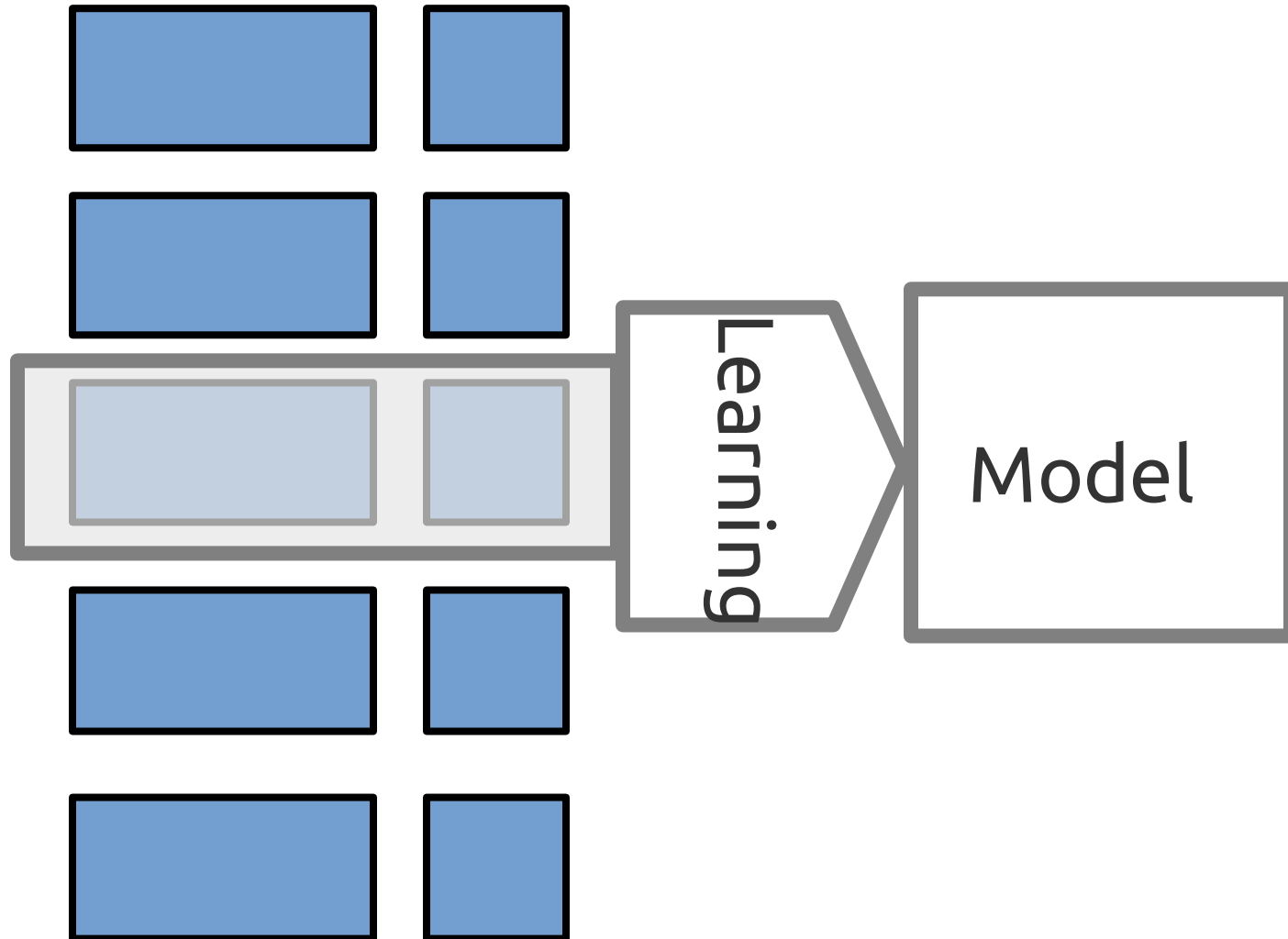


# Online learning

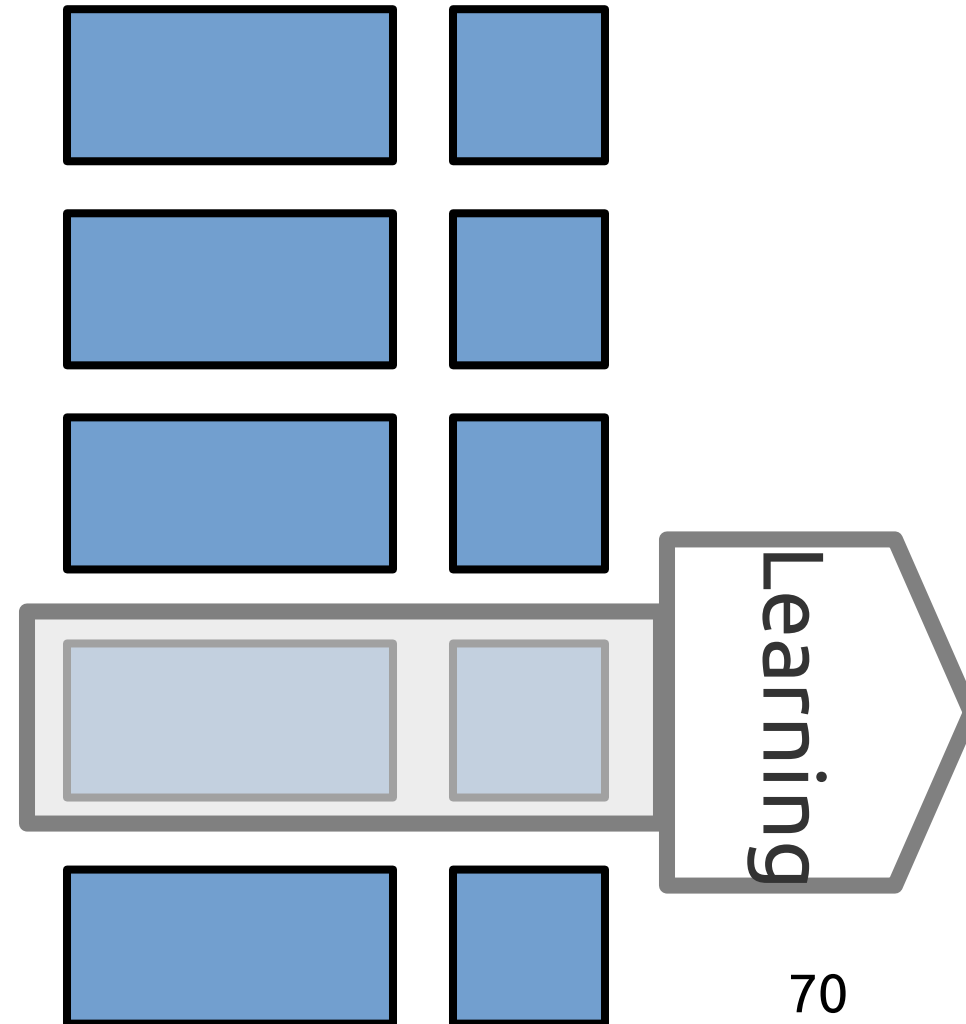
- Perceptron looks at one example at a time
- Online learners are good for streams of data
  - Social media posts
  - Photo uploads
  - User queries on a search engine

# Online

Step 3



Step 4



# Online vs Batch

- **Batch** algorithm has to remember whole dataset
- **Online** algorithm only remembers current example
- Perceptron can imitate batch learning by iterating over data several times

# Evaluation in pure online learning

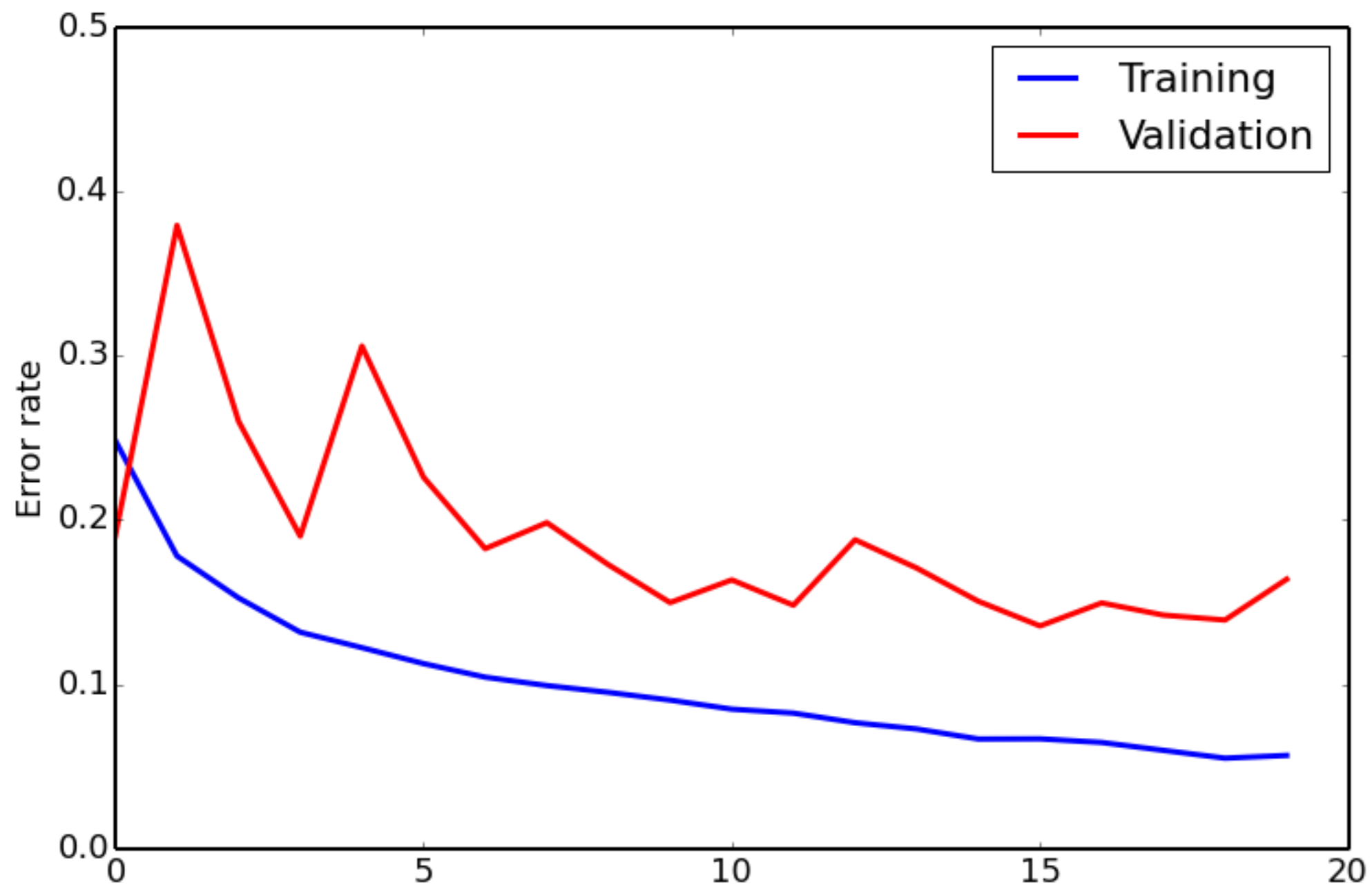
- Make prediction for current example
- Record if correct or not
- (Update model), go to next example
- At each point in time: error rate = proportion of mistakes made so far relative to total examples seen so far

# Evaluation with multiple iterations

- When simulating batch learning by using **multiple iterations**, we would be evaluating on previously seen examples
- Use separate **development** set!

# Learning ratings of movies in the sentiment dataset

- 25,000 movie reviews, positive and negative
- Use 5,000 for validation, 20,000 for training
- 20 iterations





# Early stopping

- Number of iterations is a kind of hyperparameter of the “batch” Perceptron
- Stop training when error on validation data stops dropping
- When training error goes down, but validation goes up, we're **overfitting**

# Weight averaging

- We can remember the weights from each iteration, and average them
  - Creates an ensemble of perceptrons!
- In practice, such weights generalize better than the “best” single batch

# The meaning of Percetron's most important features

- Bottom 10
  - waste worst poorly mess awful disappointment fails  
lacks annoying worse
- Top 10
  - subtitles captures enjoyable subtle noir surprisingly  
today excellent wonderfully perfect
- Around 0:
  - very character since during you're second stories  
particularly yourself hit

# Agenda

- Review of models
- Parts of a perceptron
- Making decisions with a learned model
- Learning a perceptron representation
- Batch and Online learning
- Sparse and Dense representations

# Sparsity

$$\begin{pmatrix} 0 & 3 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 7 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 & 0 & 0 & 0 & 7 \\ 0 & 0 & 0 & 5 & 0 & 0 & 0 & 8 & 0 & 0 \\ 2 & 0 & 0 & 0 & 0 & 0 & 6 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 9 & 0 & 0 & 0 & 4 & 0 \\ 0 & 1 & 0 & 7 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 4 \\ 0 & 0 & 7 & 0 & 0 & 0 & 0 & 5 & 0 & 0 \end{pmatrix}$$

# Representing text

- Text document often represented with word counts
- How many elements in the vector?
- Many tens or hundreds of thousands
- Use a **sparse representation** which omits zero values

# Dense representation

	#good	#dark	#mediocre	#the
$\mathbf{x}^1 =$	( 2,	0,	0,	5 )
$\mathbf{x}^2 =$	( 0,	1,	2,	7 )

# Sparse representation

$V = ( \text{\#good} \text{\#dark} \text{\#mediocre} \text{\#the} )$

$$\mathbf{x}^1 = \{ 0:2, 3:5 \}$$

$$\mathbf{x}^2 = \{ 1:1, 2:2, 3:7 \}$$

**All absent values are implicitly zero.**



# Sparse vectors in Python

- Python dictionaries
- Sparse matrices in **scipy**
- Some ML algorithms, in scikit-learn and elsewhere, can be optimized to work even more efficiently with sparse data

# Summary

- Perceptron as a linear model for classification
- Perceptron algorithm learns from feedback on mistakes
- Online vs batch learning

# Image credits

Frank Rosenblatt

[http://www.rutherfordjournal.org/images/TAHC\\_rosenblatt-sepia.jpg](http://www.rutherfordjournal.org/images/TAHC_rosenblatt-sepia.jpg)