

# Linear models II

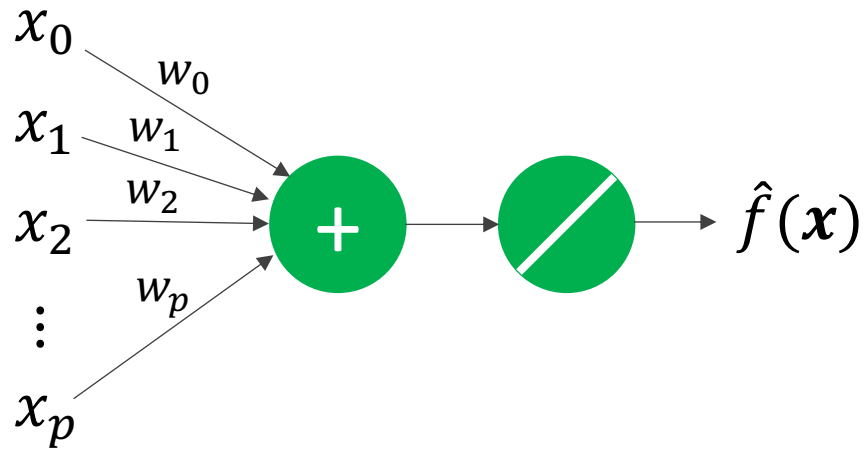
## Lecture 07

# Quiz

# Moving from regression to classification

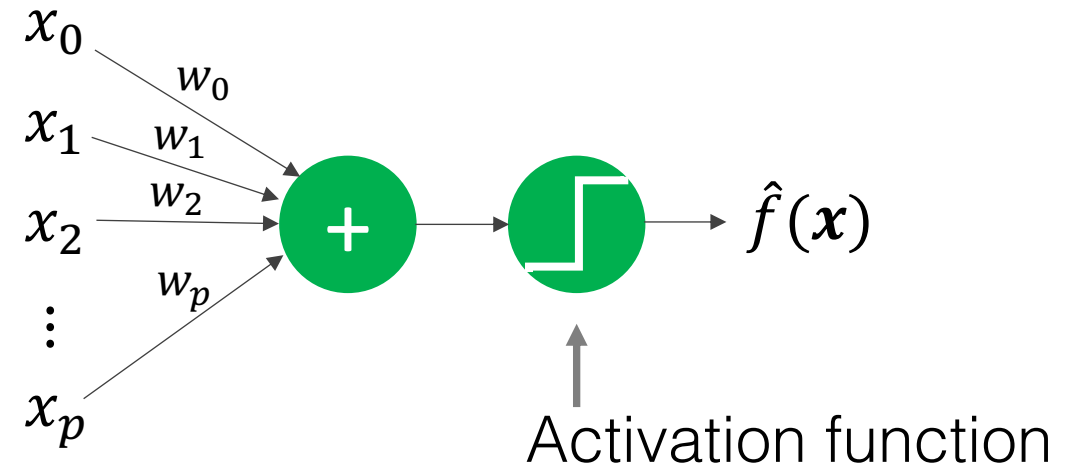
## Linear Regression

$$\hat{f}(\mathbf{x}) = \sum_{i=0}^p w_i x_i$$



## Linear Classification (perceptron)

$$\hat{f}(\mathbf{x}) = \text{sign} \left( \sum_{i=0}^p w_i x_i \right)$$



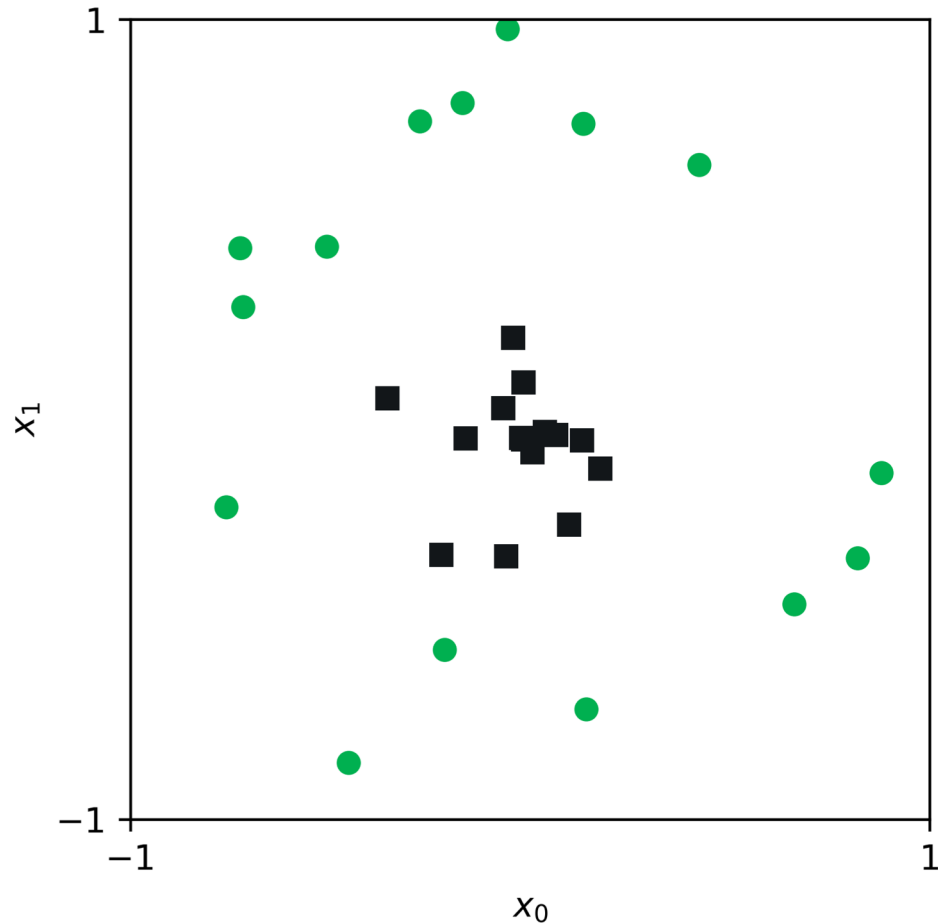
Source: Abu-Mostafa, Learning from Data, Caltech

# Can I model nonlinear relationships?

# Limitations of linear decision boundaries

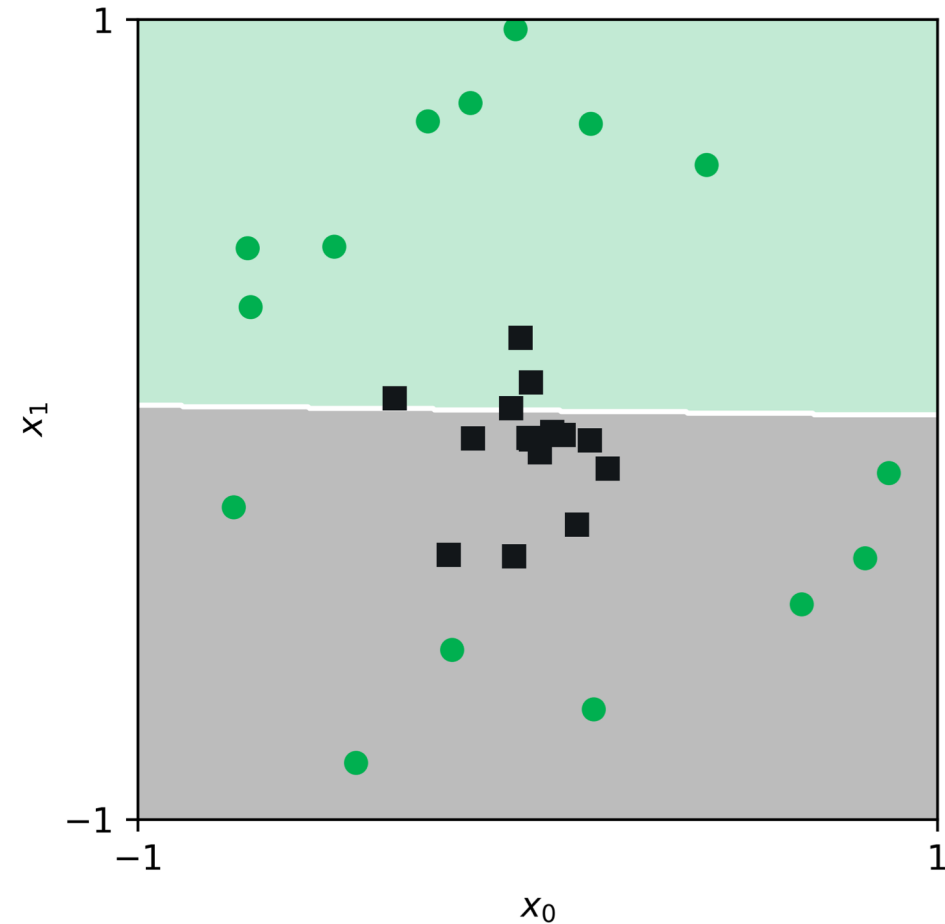
Original data

$\mathbf{x}$



Classify the features in this  $X$ -space

$$\hat{f}_{\mathbf{x}}(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$



# Transformations of features

Recall our digits example...

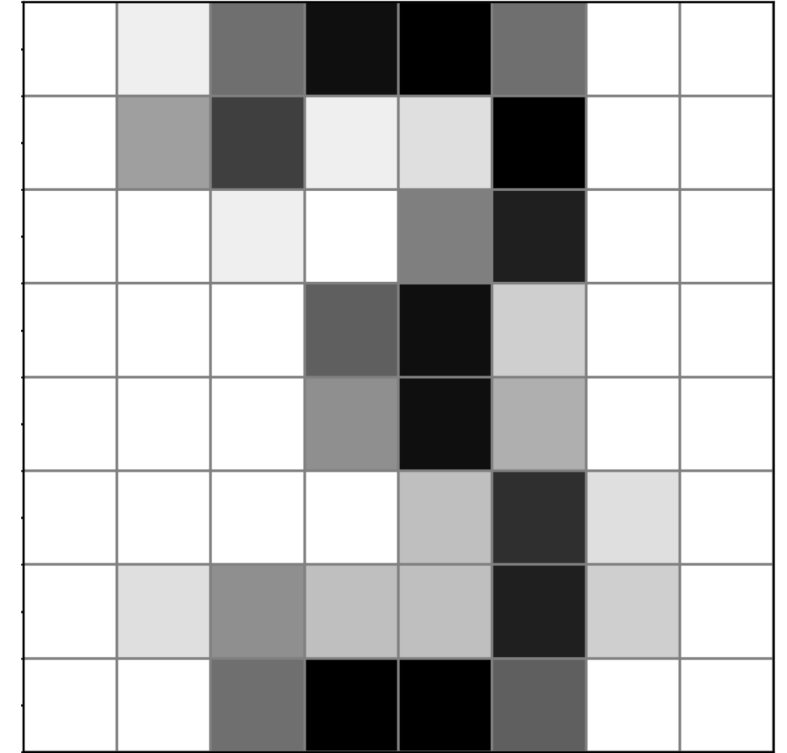
$$\mathbf{x} = [x_1, x_2, x_3, \dots, x_{64}]$$

We could create features based on the raw features. For example:

$$\mathbf{z} = [x_1 x_2, x_3^2, \frac{x_{64}}{x_{42}}]$$

Which can be written simply as variables in a new feature space:

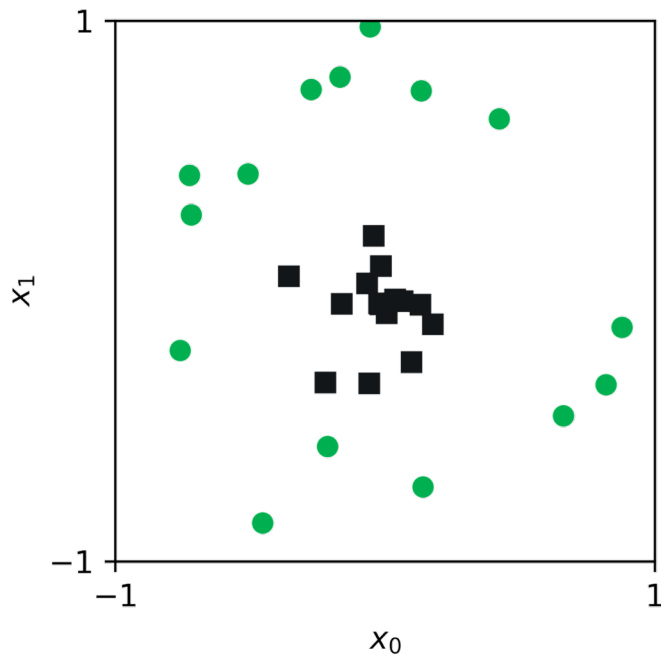
$$\mathbf{z} = [z_1, z_2, z_3]$$



Source: Abu-Mostafa, Learning from Data, Caltech

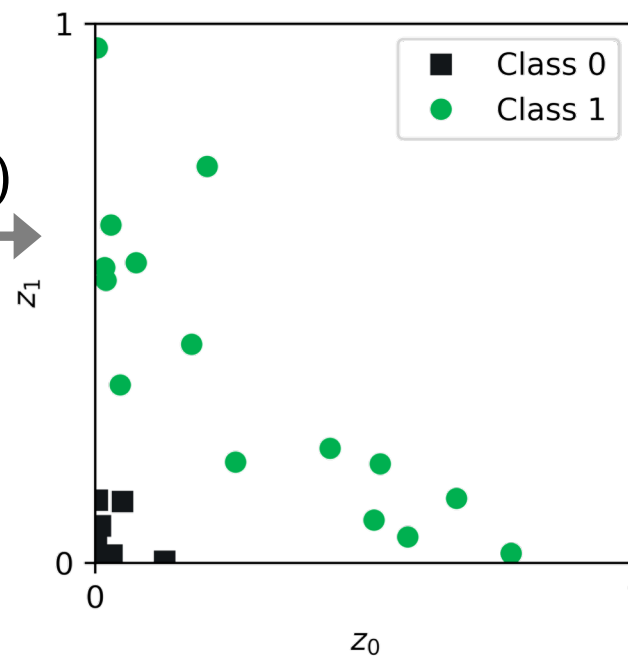
1

Original data  
 $\mathbf{x}$



transform  
the data

$$\mathbf{z} = \Phi(\mathbf{x})$$



2

This example transform  
is quadratic

$$\begin{aligned} z_i &= \Phi(x_i) = x_i^2 \\ z_0 &= x_0^2 \\ z_1 &= x_1^2 \end{aligned}$$

Classify the features  
in this Z-space

$$\hat{f}_z(\mathbf{z}) = \text{sign}(\mathbf{w}^T \mathbf{z})$$

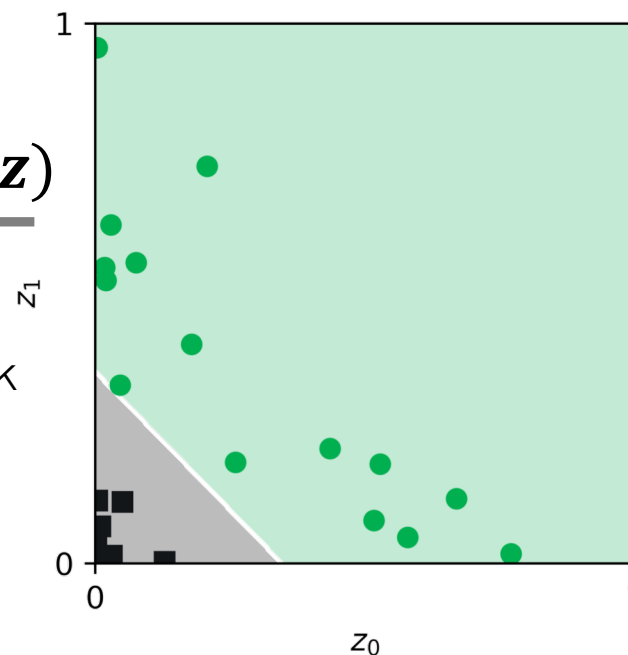
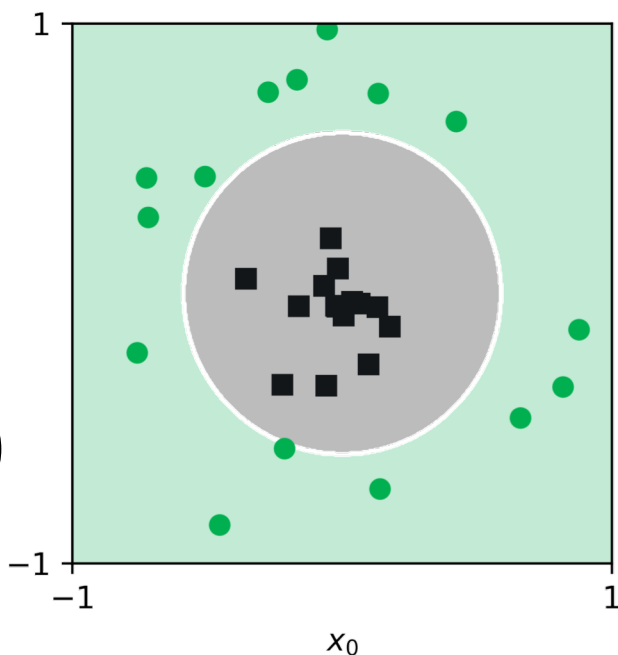
Predictions in the  $x_1$   
original X-space

$$\hat{f}(\mathbf{x}) = \hat{f}_z(\Phi(\mathbf{x}))$$

$$\mathbf{x} = \Phi^{-1}(\mathbf{z})$$

transform  
the data back

$$\begin{aligned} x_0 &= z_0^{1/2} \\ x_1 &= z_1^{1/2} \end{aligned}$$



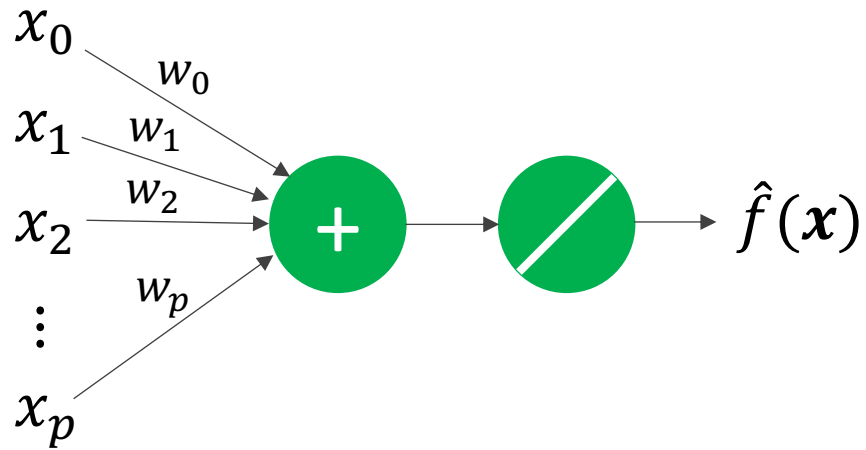
4

3

# Moving from regression to classification

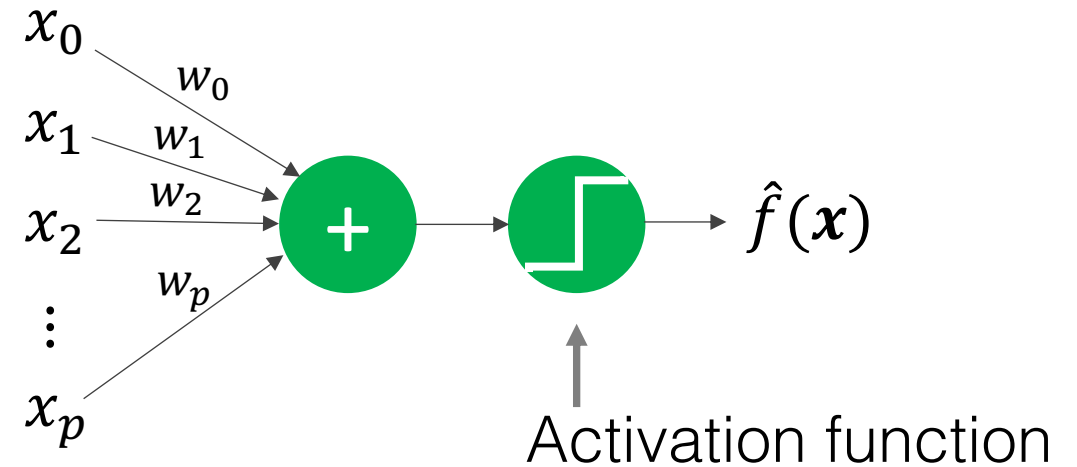
## Linear Regression

$$\hat{f}(\mathbf{x}) = \sum_{i=0}^p w_i x_i$$



## Linear Classification (perceptron)

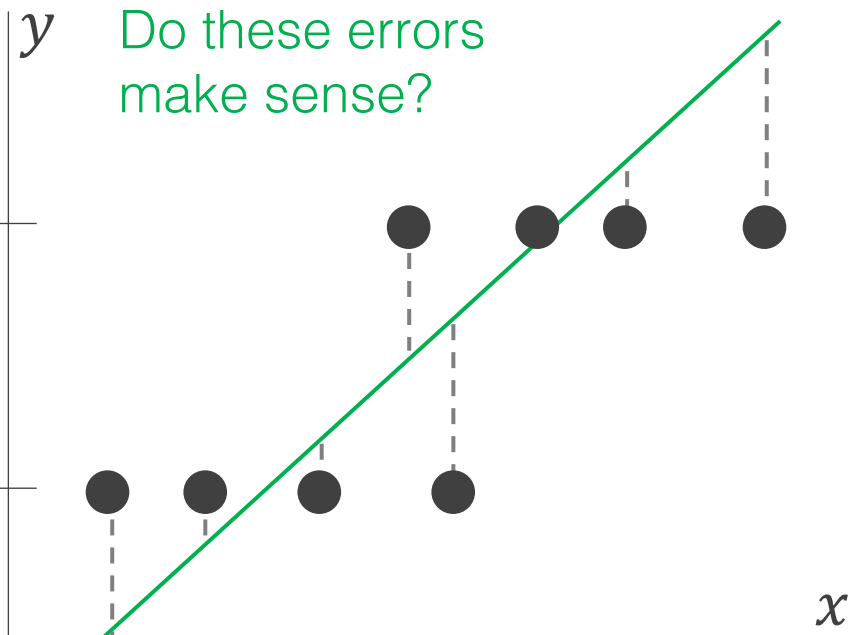
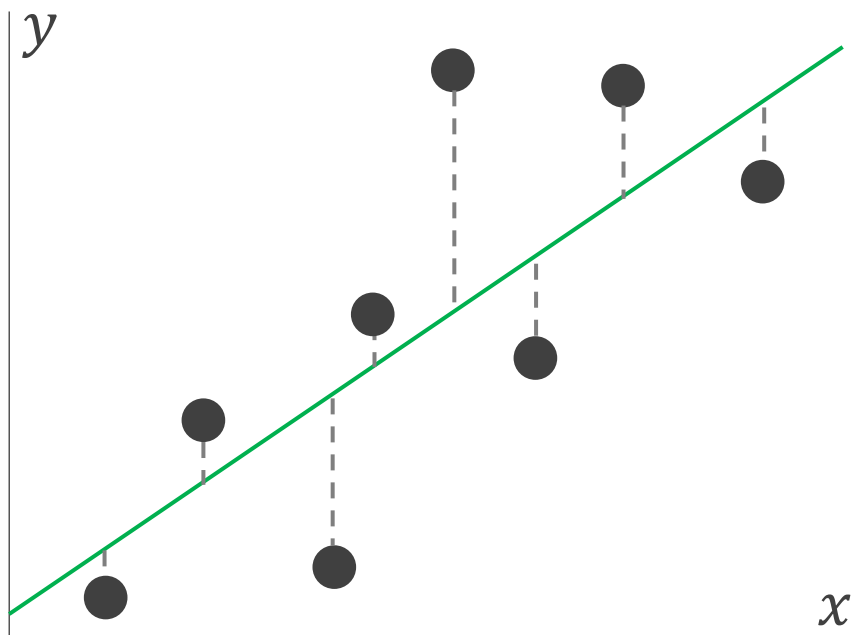
$$\hat{f}(\mathbf{x}) = \text{sign} \left( \sum_{i=0}^p w_i x_i \right)$$



Source: Abu-Mostafa, Learning from Data, Caltech

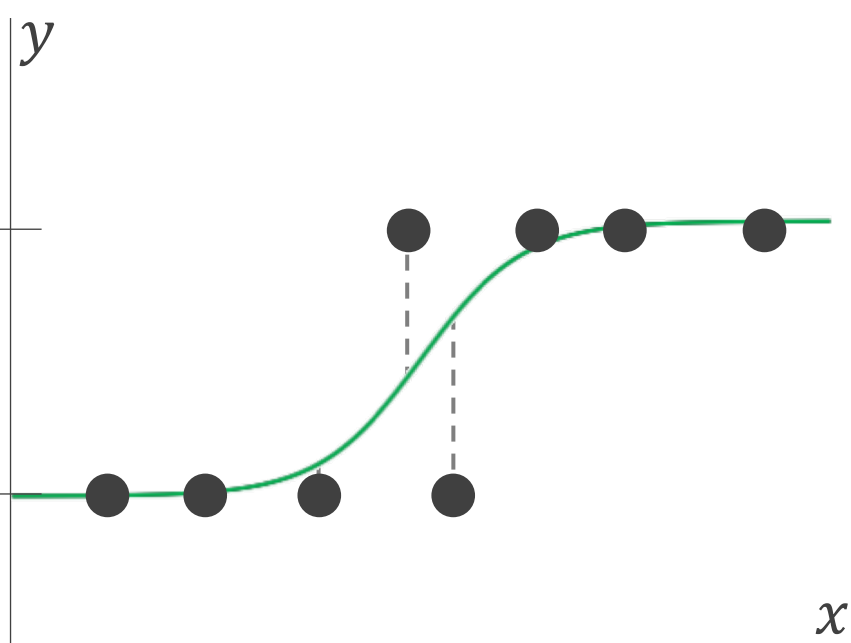
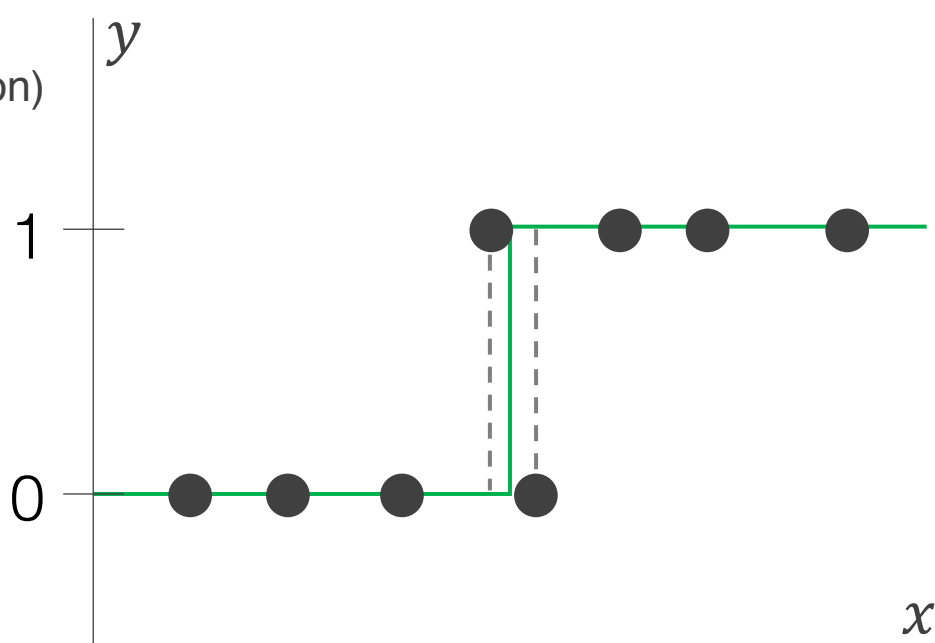


Linear regression



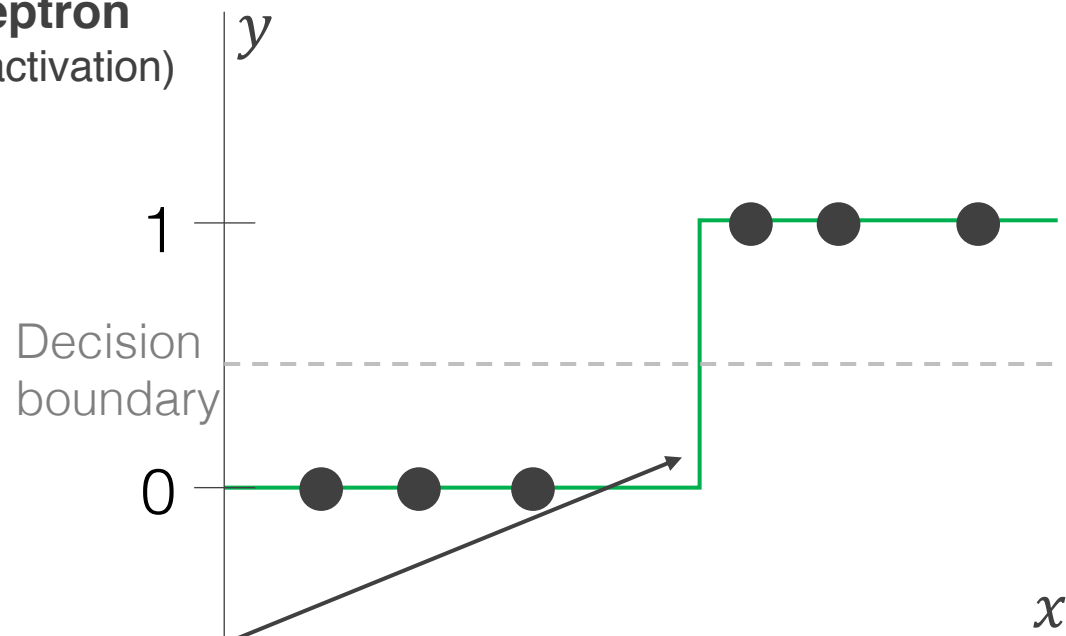
Linear regression applied to a classification problem

Perceptron (sign activation)

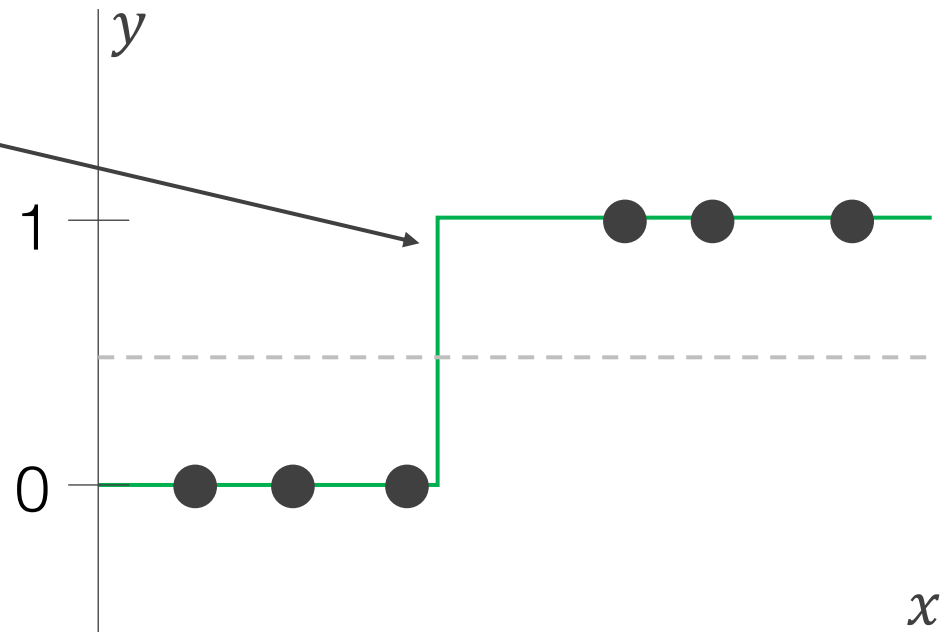


Logistic regression (sigmoid activation)

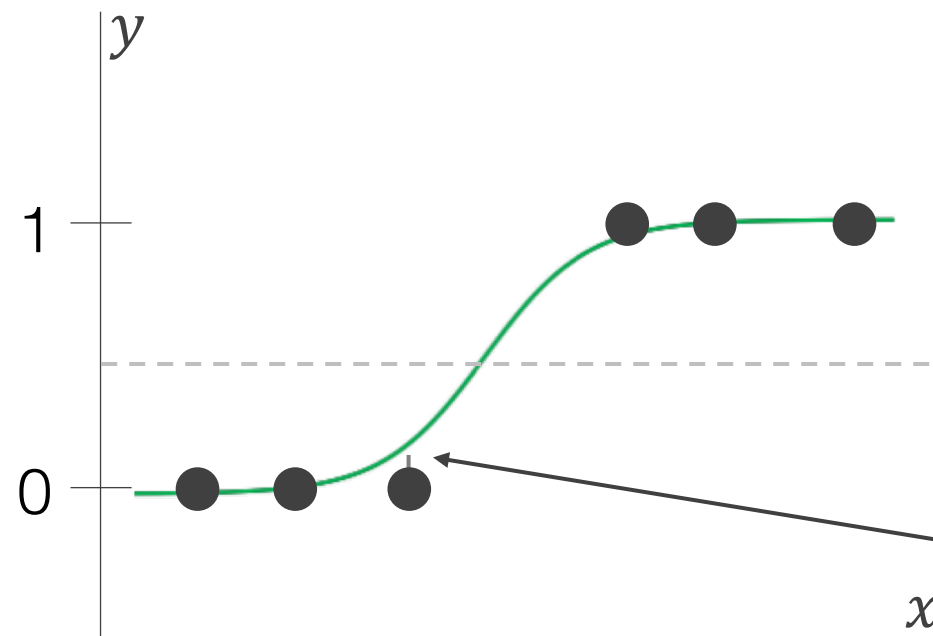
**Perceptron**  
(sign activation)



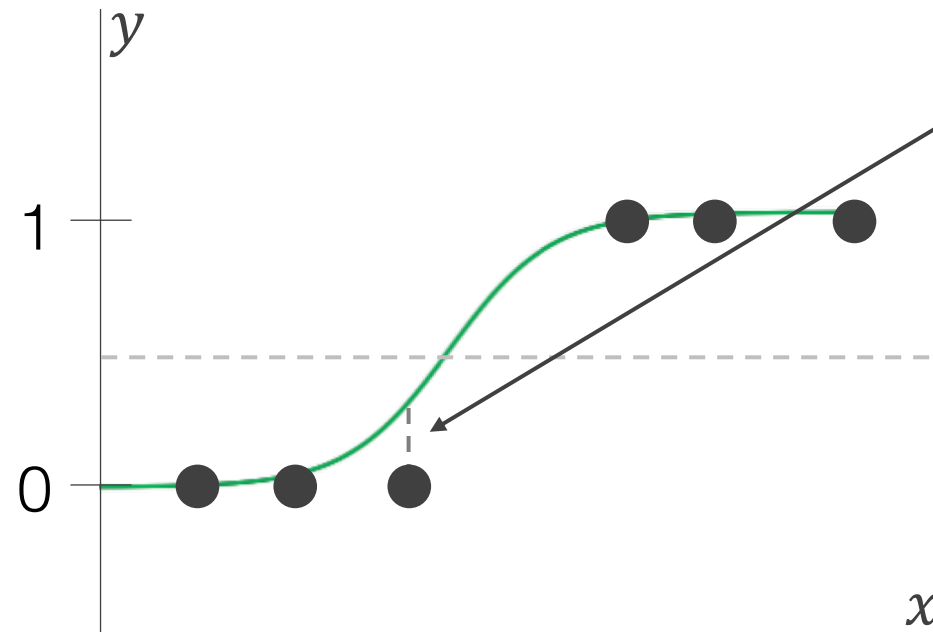
Both  
decision  
boundaries  
incur the  
same loss



**Logistic regression**  
(sigmoid activation)



The sigmoid  
assigns error  
to samples  
close to the  
margin



Favors a  
larger margin

# Sigmoid function

Definition

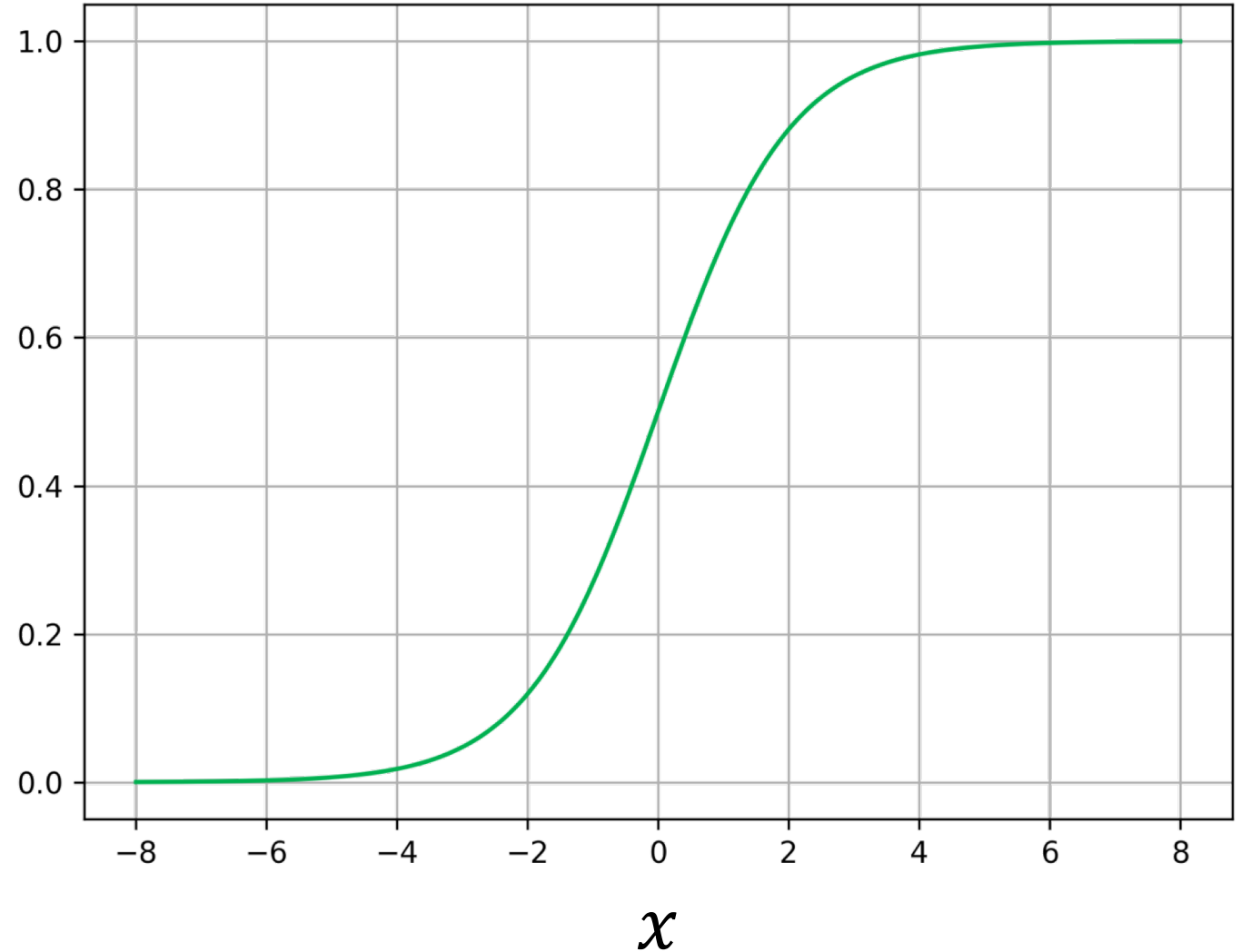
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$\sigma$

Useful properties

$$\sigma(-x) = 1 - \sigma(x)$$

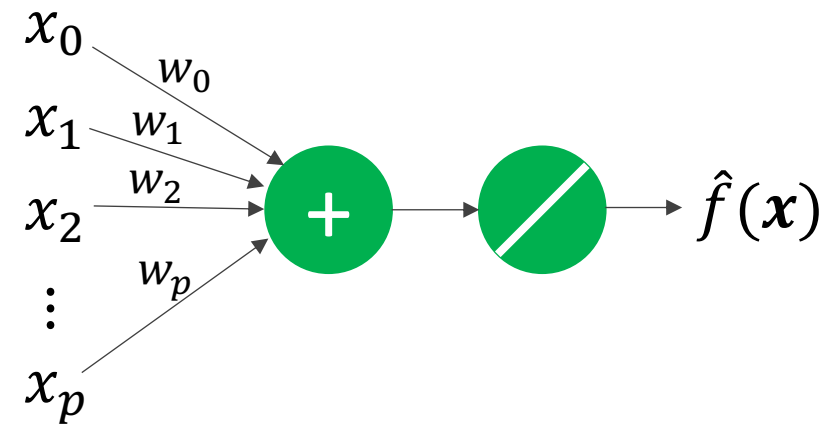
$$\frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x))$$



# Moving from regression to classification

## Linear Regression

$$\hat{f}(\mathbf{x}) = \sum_{i=0}^p w_i x_i$$

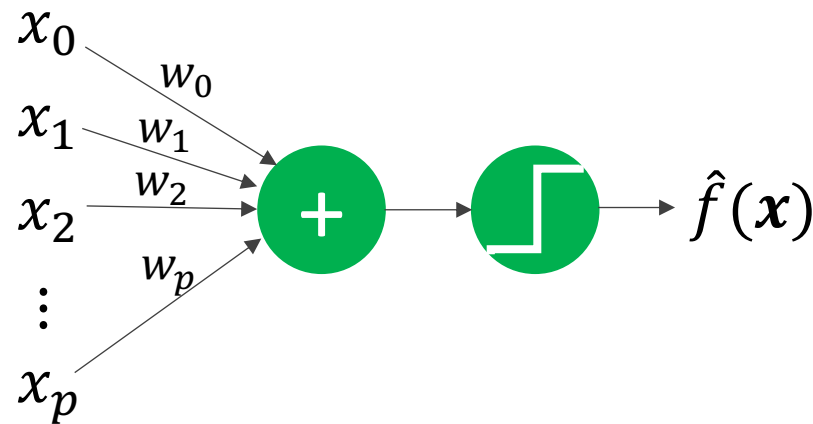


## Linear Classification

Perceptron

$$\hat{f}(\mathbf{x}) = \text{sign} \left( \sum_{i=0}^p w_i x_i \right)$$

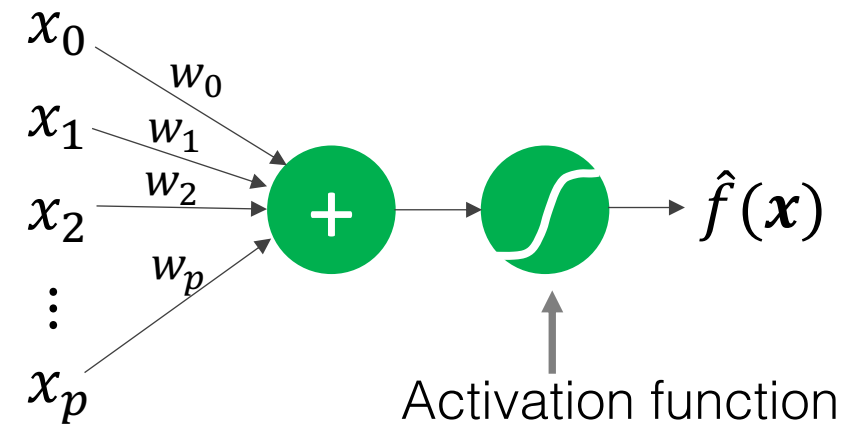
$$\text{sign}(x) = \begin{cases} 1 & x > 0 \\ -1 & \text{else} \end{cases}$$



Logistic Regression

$$\hat{f}(\mathbf{x}) = \sigma \left( \sum_{i=0}^p w_i x_i \right)$$

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



Source: Abu-Mostafa, Learning from Data, Caltech

# We take our steps to fitting our model

1. Define a cost function for measuring the fit
2. Optimize the cost function by adjusting model parameters
  - a. Calculate the gradient
  - b. Set the gradient to zero
  - c. Solve for the model parameters

# We COULD use the same cost function

Define the previous cost function

$$C(\mathbf{w}) \triangleq E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\hat{f}(\mathbf{x}_n, \mathbf{w}) - y_n)^2$$

Plug in our model

$$C(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n)^2$$

$$\hat{f}(\mathbf{x}_n, \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}_n)$$

Calculate the gradient

$$\nabla_{\mathbf{w}} C(\mathbf{w}) = \frac{2}{N} \sum_{n=1}^N [\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n] \sigma(\mathbf{w}^T \mathbf{x}_n) [\mathbf{1} - \sigma(\mathbf{w}^T \mathbf{x}_n)] \mathbf{x}_n$$

Set the gradient to zero and solve for  $\mathbf{w}$

$$\nabla_{\mathbf{w}} C(\mathbf{w}) = \mathbf{0}$$

**But we don't for logistic regression...**

There's a another cost function to use...

# Refresher: Maximum Likelihood Estimation

We purchase a bunch of scratch tickets (1,000 of them) and want to determine the probability of them being a winner

Assume we have  $N = 1,000$  **independent** Bernoulli random variables

$$\begin{aligned}P(X = 1) &= p \\P(X = 0) &= 1 - p\end{aligned}$$

**Goal:** find the value of  $p$  that maximizes the likelihood of our data



**Goal:** find the value of  $p$  that maximizes the likelihood of our data

$$P(X = 1) = p$$

$$P(X = 0) = 1 - p$$

For a **single observation**, the likelihood is:

$$L(x_i) = P(x_i|p) = p^{x_i}(1 - p)^{1-x_i}$$

For a **multiple independent observations**, the likelihood is:

$$\begin{aligned} L(\mathbf{x}) = P(\mathbf{x}|p) &= \prod_{i=1}^N P(x_i|p) \\ &= p^{\sum x_i} (1 - p)^{N - \sum x_i} \end{aligned}$$

**Goal:** find the value of  $p$  that maximizes the likelihood of our data

$$P(\mathbf{x}|p) = p^{\sum x_i} (1 - p)^{N - \sum x_i}$$

Maximizing the likelihood is equivalent to maximizing the log-likelihood

$$\ln[P(\mathbf{x}|p)] = \ln[p^{\sum x_i} (1 - p)^{N - \sum x_i}]$$

$$\ln[P(\mathbf{x}|p)] = \ln(p) \sum_{i=1}^N x_i + \ln(1 - p) \left[ N - \sum_{i=1}^N x_i \right]$$

We take the **derivative of this log likelihood and set it to zero**, then solve for  $p$

**Goal:** find the value of  $p$  that maximizes the likelihood of our data

We take the derivative of this log likelihood and set it to zero, then solve for  $p$

$$\ln[P(\mathbf{x}|p)] = \ln(p) \sum_{i=1}^N x_i + \ln(1-p) \left[ N - \sum_{i=1}^N x_i \right]$$

$$\frac{\partial \ln[P(\mathbf{x}|p)]}{\partial p} = \frac{\sum_{i=1}^N x_i}{p} + \frac{N - \sum_{i=1}^N x_i}{1-p} = 0$$

This results in our estimate being the mean of our observations:

$$\hat{p} = \frac{1}{N} \sum_{i=1}^N x_i$$

# Another interpretation of logistic regression

Our model:  $\hat{y} = \hat{f}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$

$$\sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

Logistic regression models the **probability that features belong to a class**

$$P(y_i = 1 | \mathbf{x}_i) = \sigma(\mathbf{w}^T \mathbf{x}_i)$$

$$P(y_i = 0 | \mathbf{x}_i) = 1 - \sigma(\mathbf{w}^T \mathbf{x}_i)$$

# The interpretation of the **Likelihood**

The probability of observing the class labels  $y_1, y_2, \dots, y_N$  corresponding to  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$

The likelihood for **one observation**:

$$P(y_i|\mathbf{x}_i) = P(y_i = 1|\mathbf{x}_i)^{y_i}P(y_i = 0|\mathbf{x}_i)^{1-y_i}$$

The likelihood for **all observations**:

$$P(\mathbf{y}|\mathbf{X}) = P(y_1, y_2, \dots, y_N|\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) = \prod_{i=1}^N P(y_i|\mathbf{x}_i)$$

Source: Malik Magdon-Ismail, Learning from Data

The likelihood for all observations:

$$\begin{aligned} P(\mathbf{y}|\mathbf{X}) &= \prod_{i=1}^N P(y_i|\mathbf{x}_i) = \prod_{i=1}^N P(y_i = 1|\mathbf{x}_i)^{y_i} P(y_i = 0|\mathbf{x}_i)^{1-y_i} \\ &= \prod_{i=1}^N \sigma(\mathbf{w}^T \mathbf{x}_i)^{y_i} [1 - \sigma(\mathbf{w}^T \mathbf{x}_i)]^{1-y_i} \end{aligned}$$

**This is our cost function**

(to be precise, the negative of the cost function)

We can take the logarithm, then the gradient, then set equal to zero...

**This is not solvable in closed form: need a new approach**

# Gradient descent

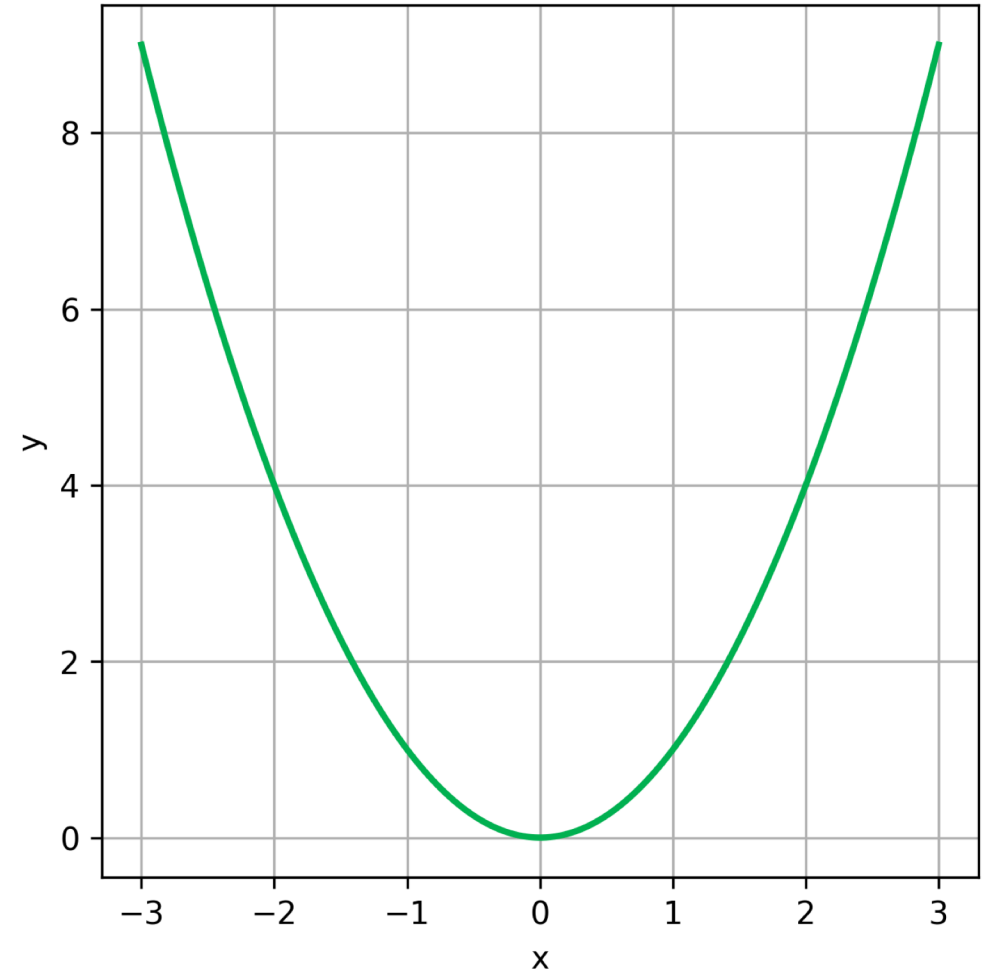
Minimize  $y = x^2$

We start at a point and want to “roll”  
down to the minimum

$$\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} + \eta \mathbf{v}$$

Learning  
rate

Direction  
to move in



# Gradient descent

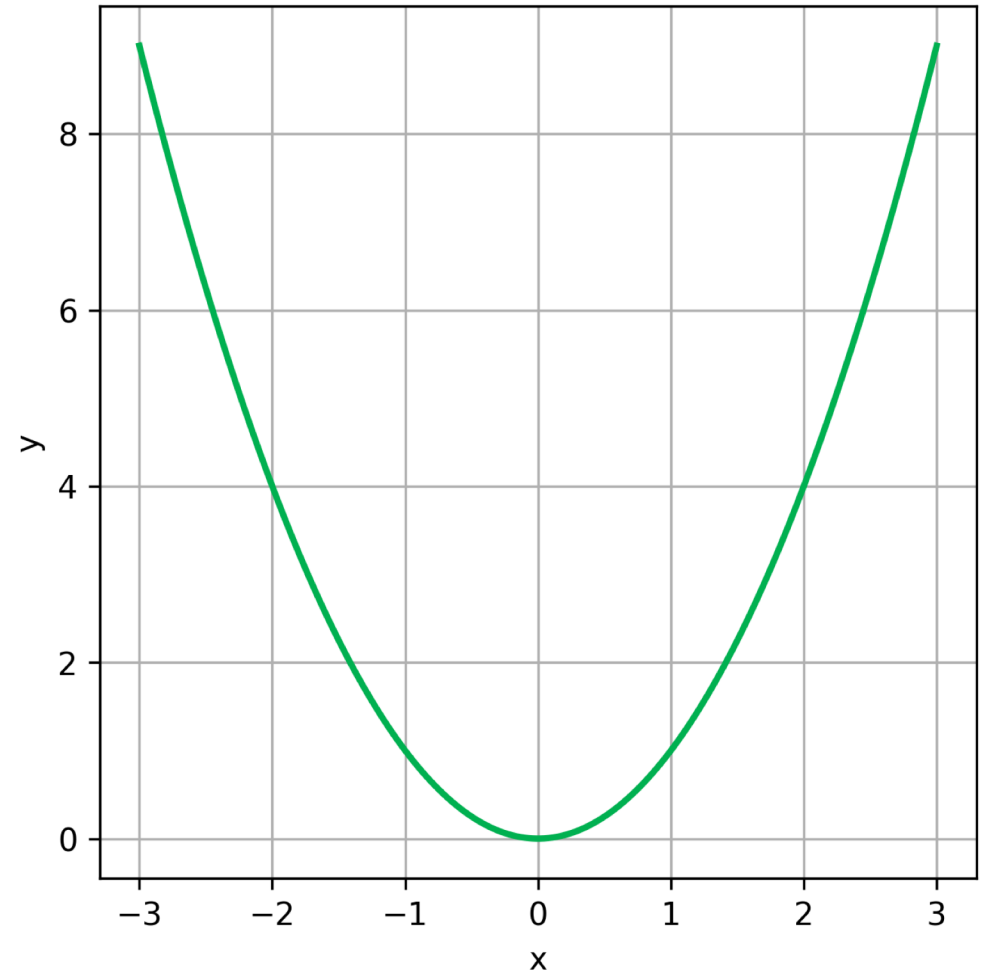
Minimize  $f(x) = x^2$

The gradient points in the direction of steepest **positive** change

$$\frac{df(x)}{dx} = 2x$$

We want to move in the **opposite** direction of the gradient

$$\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$$





# Gradient descent

Minimize  $f(x) = x^2$

Assume  $x^{(0)} = 2$  and  $\eta = 0.25$

$$x^{(i+1)} = x^{(i)} - (0.25)(2x^{(i)})$$

$$x^{(i+1)} = x^{(i)} - (0.5)x^{(i)}$$

$i$	$x^{(i)}$	$y^{(i)}$
-----	-----------	-----------

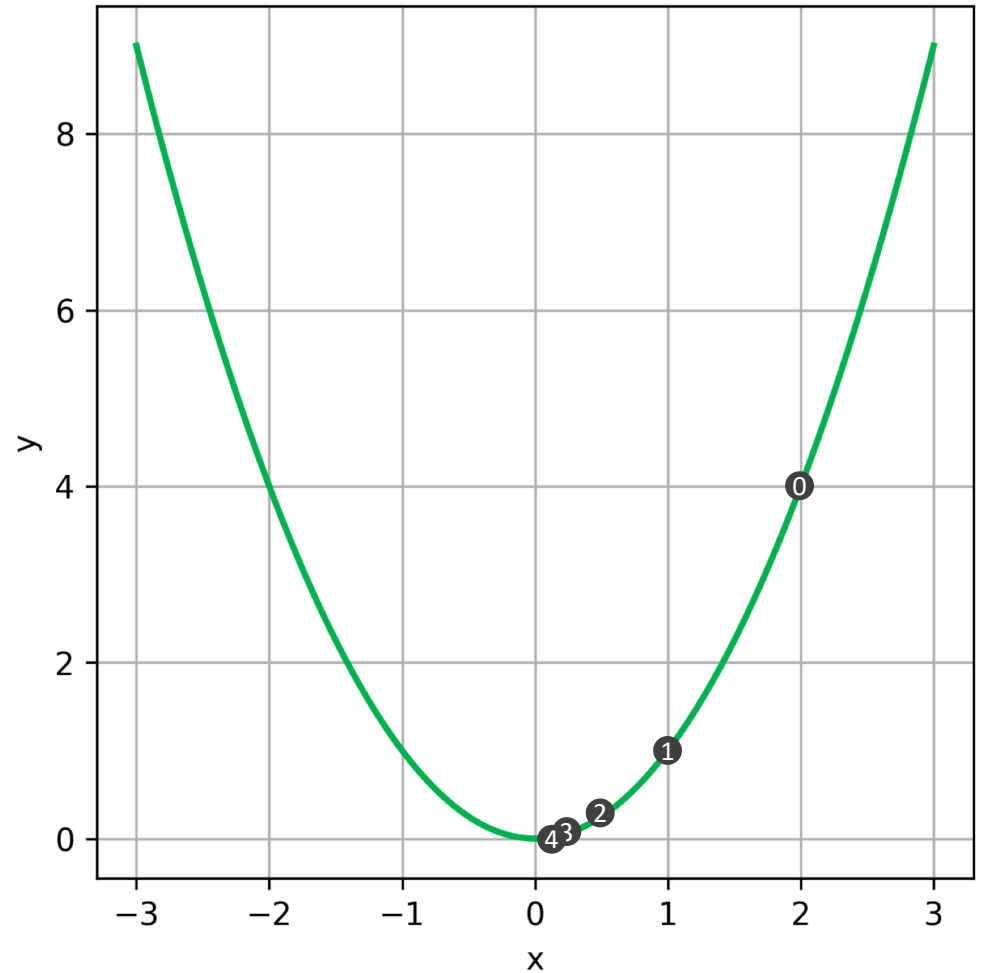
0	2	4
---	---	---

1	1	1
---	---	---

2	0.5	0.25
---	-----	------

3	0.25	0.0625
---	------	--------

4	0.125	0.0156
---	-------	--------



# Takeaways

- Transformations of features may help to overcome nonlinearities
- Logistic regression is much better suited for classification than linear regression
- Logistic regression parameters must be estimated iteratively, and a method for that optimization is gradient descent
- Gradient descent can be used for cost function optimization and there are a number of variants