


Daily Inspiration



Don't stop
ignoring the data

I'm InspiroBot.

I am an artificial intelligence dedicated to generating unlimited amounts of unique inspirational quotes for endless enrichment of pointless human existence.



Today

- Recap yesterday
- Logistic regression: using regression tools for classification
- Neural network basics

Yesterday

- Cost function: (differentiable) function that shows how wrong an estimate is for given parameters.
- Gradient descent: one common way to minimise the cost function automatically, i.e. to get optimal parameters
- Linear regression: very simple model that assumes that value to predict is linear combination of input features.
- Overfitting and underfitting, bias and variance: want our model to work well for unseen data. Need just enough model freedom given the complexity of our problem. How:
 - Cross-validation to measure ability to generalise + get best hyperparameters
 - Use learning curves to diagnose bias vs. variance

Gradient descent in linear algebra

- Goal gradient descent: take a small step in every parameter such that you get closer to the minimum of the cost. Return new theta's.

$$\theta_{0new} = \theta_0 - \frac{\alpha}{m} \sum_{i=1}^m [(h_{\theta}(x^{(i)}) - y^{(i)}) \cdot 1]$$

$$\theta_{1new} = \theta_1 - \frac{\alpha}{m} \sum_{i=1}^m [(h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)}]$$

$$\theta_{2new} = \theta_2 - \frac{\alpha}{m} \sum_{i=1}^m [(h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_2^{(i)}]$$

Gradient descent in linear algebra

We have data, known values, and initial theta's:

$$X = \begin{bmatrix} 1 & feat_1 val_1 & feat_2 val_1 \\ 1 & feat_1 val_2 & feat_2 val_2 \end{bmatrix}; y = \begin{bmatrix} 10.23 \\ -4 \end{bmatrix}; params = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

Get predicted values:

$$\begin{bmatrix} 1 & feat_1 val_1 & feat_2 val_1 \\ 1 & feat_1 val_2 & feat_2 val_2 \end{bmatrix} @ \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} 9.23 \\ -2.5 \end{bmatrix}$$

2 by 3 times 3 by 1 gives 2 by 1 (rows by columns)

Get errors:

$$errs = \begin{bmatrix} 9.23 \\ -2.5 \end{bmatrix} - y = \begin{bmatrix} 9.23 \\ -2.5 \end{bmatrix} - \begin{bmatrix} 10.23 \\ -4 \end{bmatrix} = \begin{bmatrix} -1 \\ 1.5 \end{bmatrix}$$

Gradient descent in linear algebra

$$errs = \begin{bmatrix} -1 \\ 1.5 \end{bmatrix}$$

$$\theta_{0new} = \theta_0 - \frac{a}{m} \sum_{i=1}^m ((h_{\theta}(x^{(i)}) - y^{(i)}) \cdot 1)$$

$$\theta_{1new} = \theta_1 - \frac{a}{m} \sum_{i=1}^m ((h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)})$$

$$\theta_{2new} = \theta_2 - \frac{a}{m} \sum_{i=1}^m ((h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_2^{(i)})$$

Gradient descent in linear algebra

Calculate, for each feature, sum of each error times that feature:

$$\begin{bmatrix} -1 \\ 1.5 \end{bmatrix}^T = [-1 \quad 1.5]$$

$$errs = \begin{bmatrix} -1 \\ 1.5 \end{bmatrix}$$

$$[-1 \quad 1.5] @ \begin{bmatrix} 1 & feat_1 val_1 & feat_2 val_1 \\ 1 & feat_1 val_2 & feat_2 val_2 \end{bmatrix} =$$

$$[-1 \cdot 1 + 1.5 \cdot 1 \quad -1 \cdot feat_1 val_1 + 1.5 \cdot feat_1 val_2 \quad -1 \cdot feat_2 val_1 + 1.5 \cdot feat_2 val_2]$$

Gradient descent in linear algebra

Calculate, for each feature, sum of each error times that feature:

$$\begin{bmatrix} -1 \\ 1.5 \end{bmatrix}^T = [-1 \quad 1.5]$$

$$errs = \begin{bmatrix} -1 \\ 1.5 \end{bmatrix}$$

$$[-1 \quad 1.5] @ \begin{bmatrix} 1 & feat_1 val_1 & feat_2 val_1 \\ 1 & feat_1 val_2 & feat_2 val_2 \end{bmatrix} =$$

$$[-1 \cdot 1 + 1.5 \cdot 1]$$

$$[-1 \cdot feat_1 val_1 + 1.5 \cdot feat_1 val_2]$$

$$[-1 \cdot feat_2 val_1 + 1.5 \cdot feat_2 val_2]$$

$$\theta_{0new} = \theta_0 - \frac{a}{m} \sum_{i=1}^m ((h_{\theta}(x^{(i)}) - y^{(i)}) \cdot 1)$$

$$\theta_{1new} = \theta_1 - \frac{a}{m} \sum_{i=1}^m ((h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)})$$

$$\theta_{2new} = \theta_2 - \frac{a}{m} \sum_{i=1}^m ((h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_2^{(i)})$$

Gradient descent in linear algebra

Now all that we need to do is multiply with α/m and subtract from our old theta's:

$$\begin{aligned} & \alpha/m \cdot \begin{bmatrix} -1 \cdot 1 + 1.5 \cdot 1 & -1 \cdot feat_1 val_1 + 1.5 \cdot feat_1 val_2 & -1 \cdot feat_2 val_1 + 1.5 \cdot feat_2 val_2 \end{bmatrix} \\ &= \begin{bmatrix} \frac{\alpha}{m}(-1 \cdot 1 + 1.5 \cdot 1) & \frac{\alpha}{m}(-1 \cdot feat_1 val_1 + 1.5 \cdot feat_1 val_2) & \frac{\alpha}{m}(-1 \cdot feat_2 val_1 + 1.5 \cdot feat_2 val_2) \end{bmatrix} \end{aligned}$$

Transpose it:

$$\begin{bmatrix} \frac{\alpha}{m}(-1 \cdot 1 + 1.5 \cdot 1) & \frac{\alpha}{m}(-1 \cdot feat_1 val_1 + 1.5 \cdot feat_1 val_2) & \frac{\alpha}{m}(-1 \cdot feat_2 val_1 + 1.5 \cdot feat_2 val_2) \end{bmatrix}^T = \begin{bmatrix} \frac{\alpha}{m}(-1 \cdot 1 + 1.5 \cdot 1) \\ \frac{\alpha}{m}(-1 \cdot feat_1 val_1 + 1.5 \cdot feat_1 val_2) \\ \frac{\alpha}{m}(-1 \cdot feat_2 val_1 + 1.5 \cdot feat_2 val_2) \end{bmatrix}$$

So finally:

$$\begin{bmatrix} \theta_{0old} \\ \theta_{1old} \\ \theta_{2old} \end{bmatrix} - \begin{bmatrix} \frac{\alpha}{m}(-1 \cdot 1 + 1.5 \cdot 1) \\ \frac{\alpha}{m}(-1 \cdot feat_1 val_1 + 1.5 \cdot feat_1 val_2) \\ \frac{\alpha}{m}(-1 \cdot feat_2 val_1 + 1.5 \cdot feat_2 val_2) \end{bmatrix} = \begin{bmatrix} \theta_{0new} \\ \theta_{1new} \\ \theta_{2new} \end{bmatrix}$$

$$\theta_{1new} = \theta_{1old} - \frac{\alpha}{m} \sum_{i=1}^m ((h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)})$$

Logistic regression

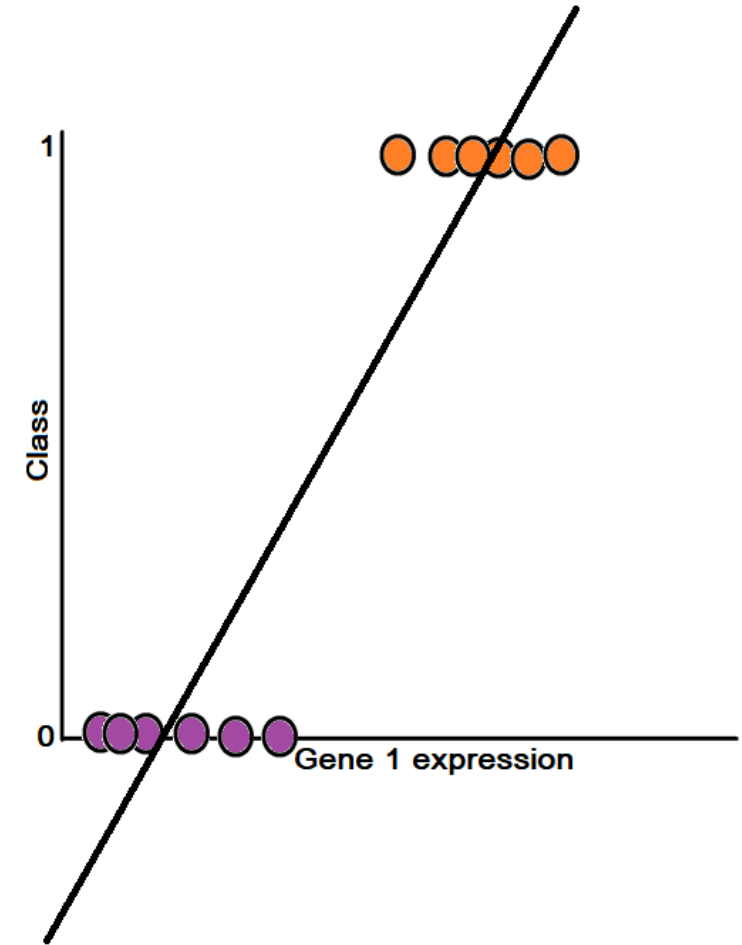
- You tell me: what is logistic regression?

Logistic regression

- Use regression-like framework for classification

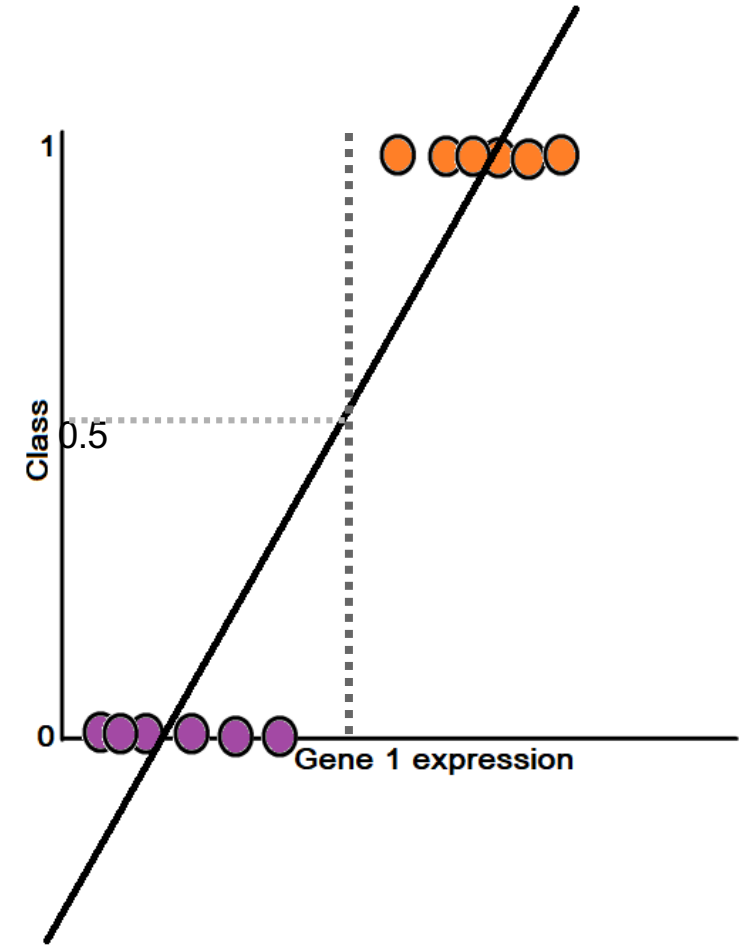
Logistic regression

- Naïve idea:
Train a linear regression. If
 $\text{Class} \geq 0.5$, predict class 1.
Otherwise, class 0.



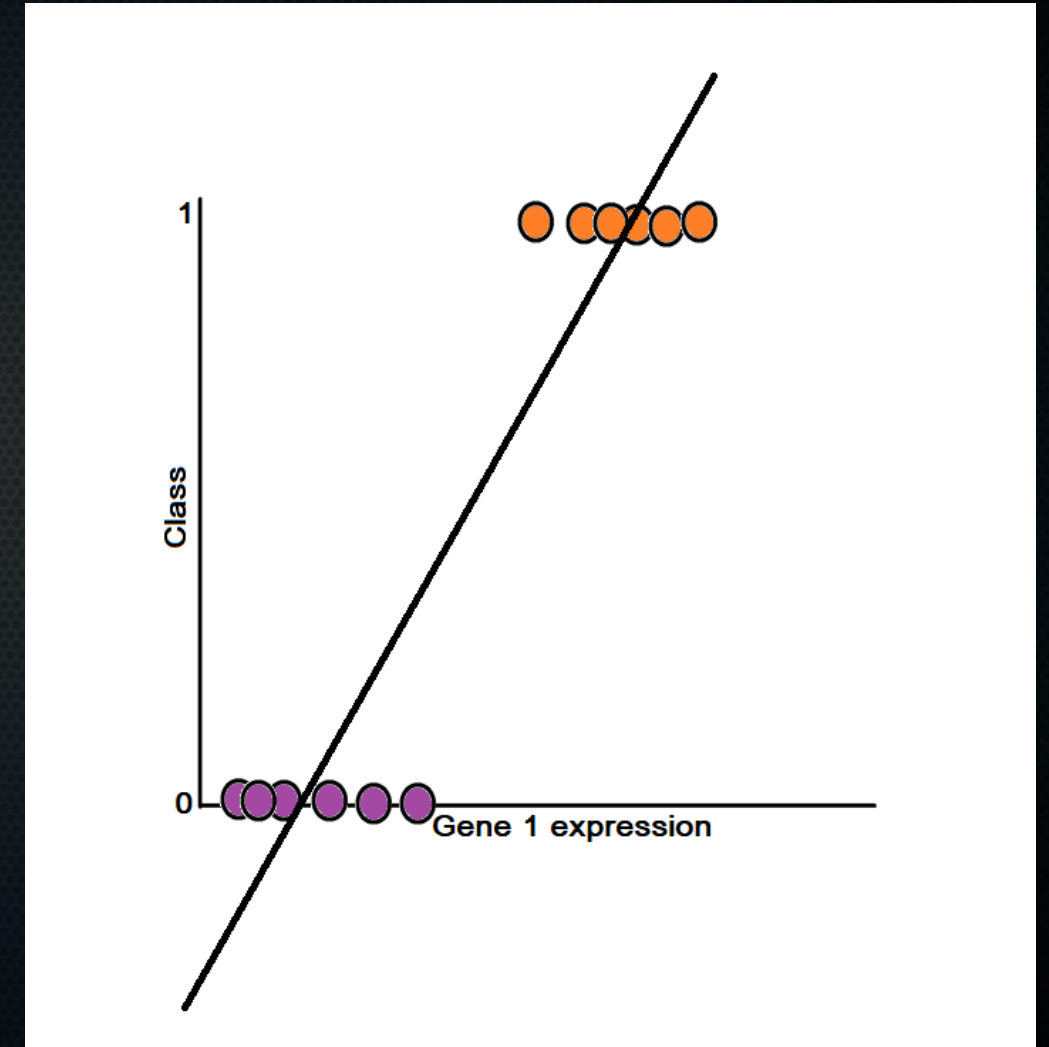
Logistic regression

- Naïve idea:
Train a linear regression. If
Class ≥ 0.5 , predict class 1.
Otherwise, class 0.



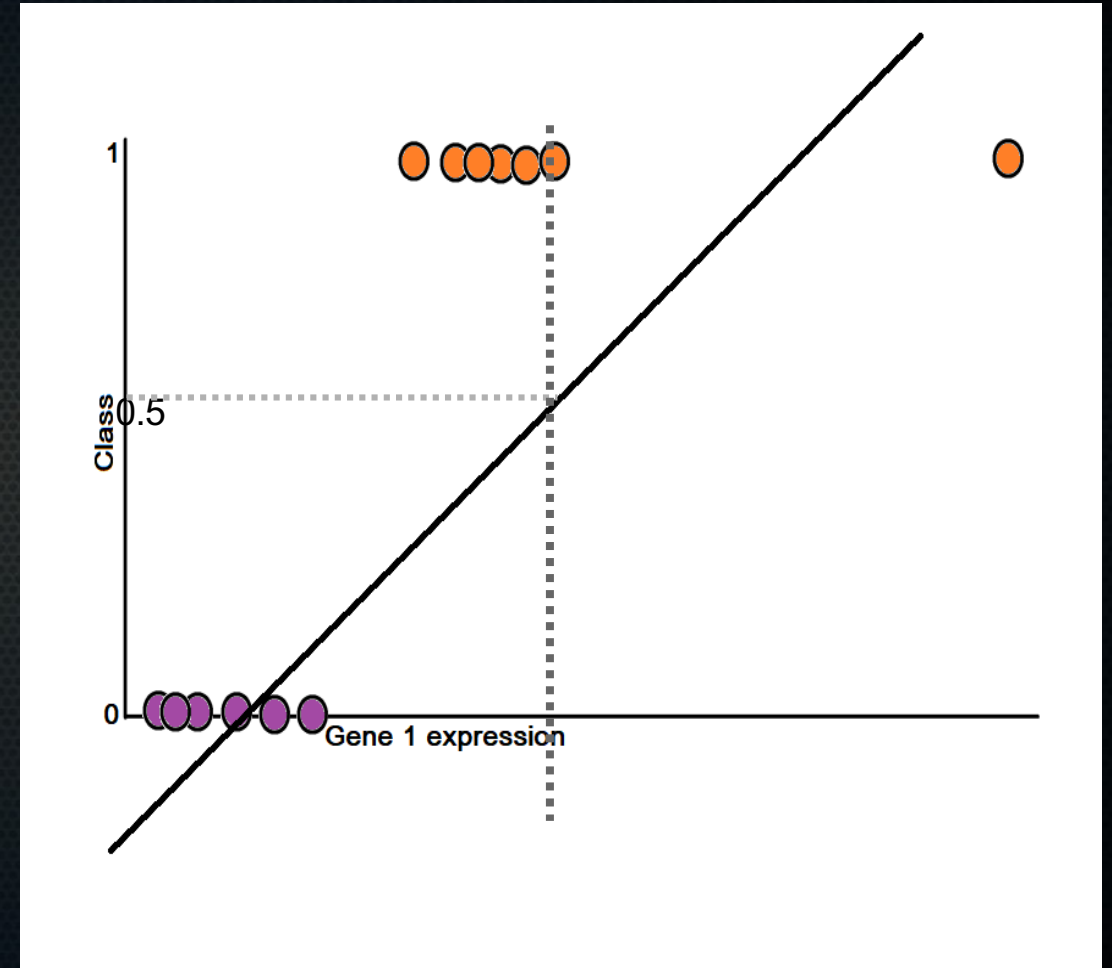
Logistic regression

- Naïve idea:
Train a linear regression. If $\text{Class} \geq 0.5$, predict class 1. Otherwise, class 0.
- Problems:
 - You can predict class > 1 and < 0 , while that is not possible in reality.



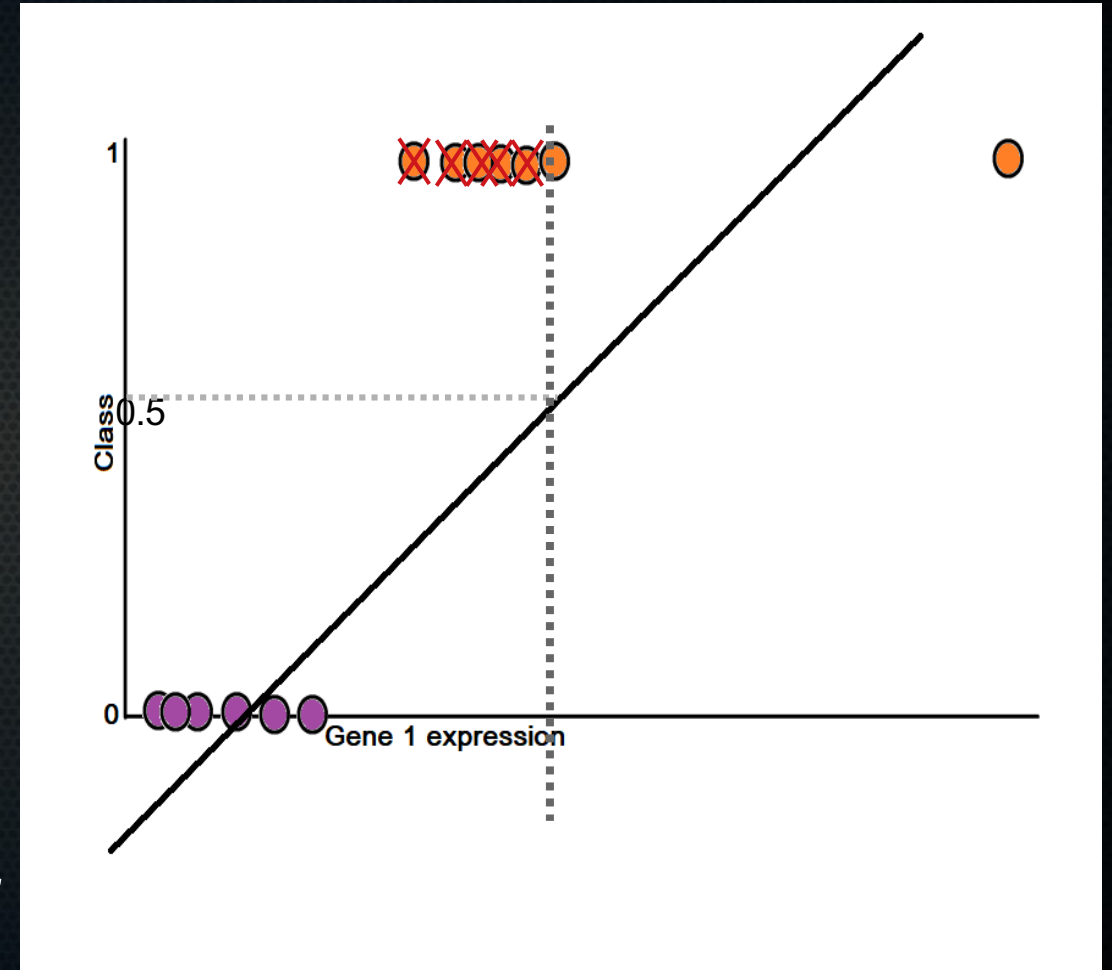
Logistic regression

- Naïve idea:
Train a linear regression. If $\text{Class} \geq 0.5$, predict class 1. Otherwise, class 0.
- Problems:
 - You can predict class > 1 and < 0 , while that is not possible in reality.
 - This example seemed to work, but quickly breaks down \rightarrow



Logistic regression

- Naïve idea:
Train a linear regression. If $\text{Class} \geq 0.5$, predict class 1. Otherwise, class 0.
- Problems:
 - You can predict class > 1 and < 0 , while that is not possible in reality.
 - This example seemed to work, but quickly breaks down \rightarrow get what is basically confirmation of hypothesis, but perform worse!



Logistic regression

- What we want:
 - Use the information that we only have two classes, 0 or 1.
 - Hypothesis function should output only numbers between 0 or 1.

Sigmoid or logistic function

- Before, our hypothesis function was of the form:

$$h_{\theta}(x) = \theta^T \cdot x$$

Sigmoid or logistic function

- Before, our hypothesis function was of the form:

$$h_{\theta}(x) = \theta^T \cdot x \quad \longrightarrow \quad [0.5 \quad 3 \quad -1.5] \cdot \begin{bmatrix} 1 \\ 3 \\ 8 \end{bmatrix}$$

Sigmoid or logistic function

- Before, our hypothesis function was of the form:

$$h_{\theta}(x) = \theta^T \cdot x \quad \longrightarrow \quad \underbrace{[0.5 \quad 3 \quad -1.5]}_{\text{Learned parameters (theta 0 – theta 2)}} \cdot \underbrace{\begin{bmatrix} 1 \\ 3 \\ 8 \end{bmatrix}}_{\text{Features for one sample (x0 = 1, intercept term, 2 data-derived features x1 and x2)}}$$

Sigmoid or logistic function

- Before, our hypothesis function was of the form:

$$h_{\theta}(x) = \theta^T \cdot x \quad \longrightarrow \quad [0.5 \quad 3 \quad -1.5] \cdot \begin{bmatrix} 1 \\ 3 \\ 8 \end{bmatrix} = 0.5 \cdot 1 + 3 \cdot 3 - 1.5 \cdot 8 = -2.5$$

Sigmoid or logistic function

- Before, our hypothesis function was of the form:

$$h_{\theta}(x) = \theta^T \cdot x$$

- Change that to the following:

$$h_{\theta}(x) = g(\theta^T \cdot x)$$

Sigmoid or logistic function

- Before, our hypothesis function was of the form:

$$h_{\theta}(x) = \theta^T \cdot x$$

- Change that to the following:

$$h_{\theta}(x) = g(\theta^T \cdot x) \longrightarrow g(z) = \frac{1}{1 + e^{-z}}$$

Sigmoid or logistic function

- Before, our hypothesis function was of the form:

$$h_{\theta}(x) = \theta^T \cdot x$$

- Change that to the following:

$$h_{\theta}(x) = g(\theta^T \cdot x) \longrightarrow g(z) = \frac{1}{1 + e^{-z}}$$

- What does that look like?

Sigmoid or logistic function

- Before, our hypothesis function was of the form:

$$h_{\theta}(x) = \theta^T \cdot x$$

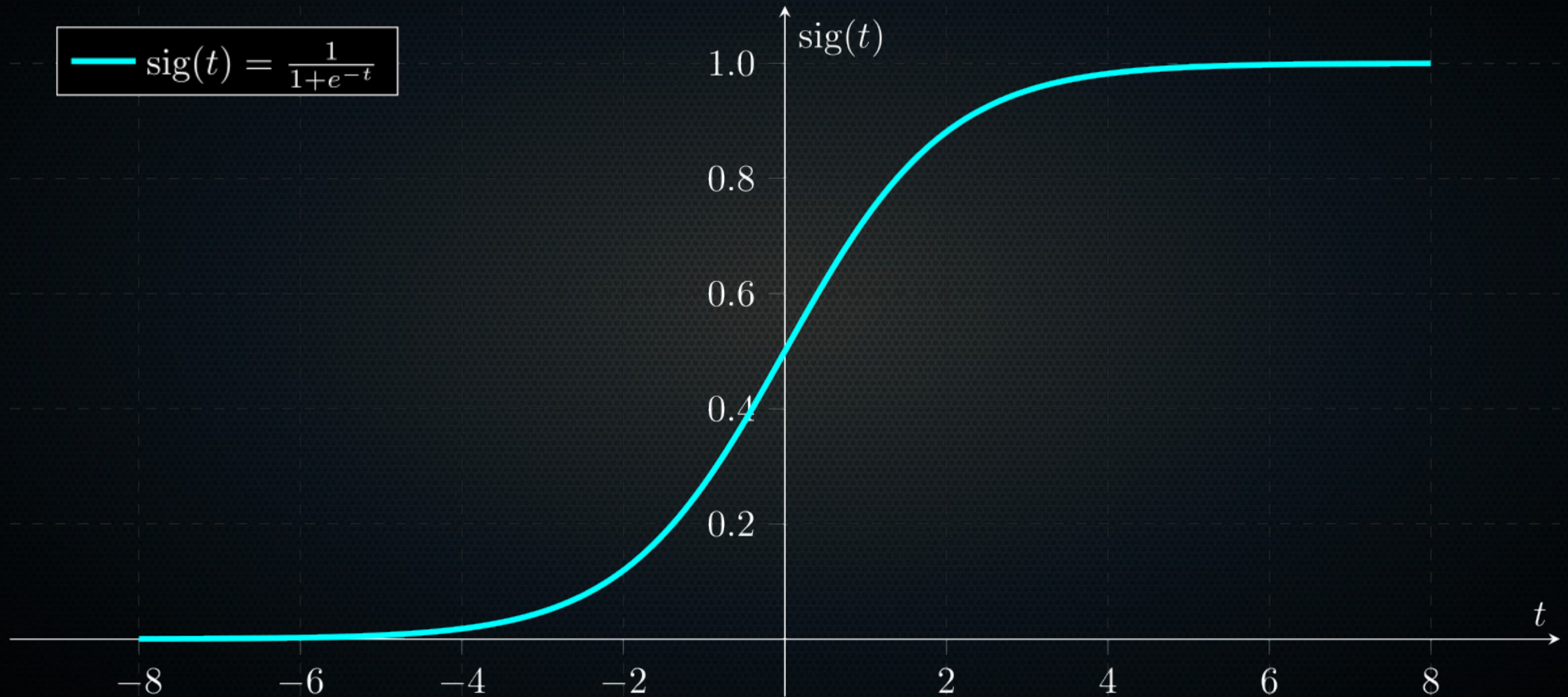
- Change that to the following:

$$h_{\theta}(x) = g(\theta^T \cdot x) \longrightarrow g(z) = \frac{1}{1 + e^{-z}}$$

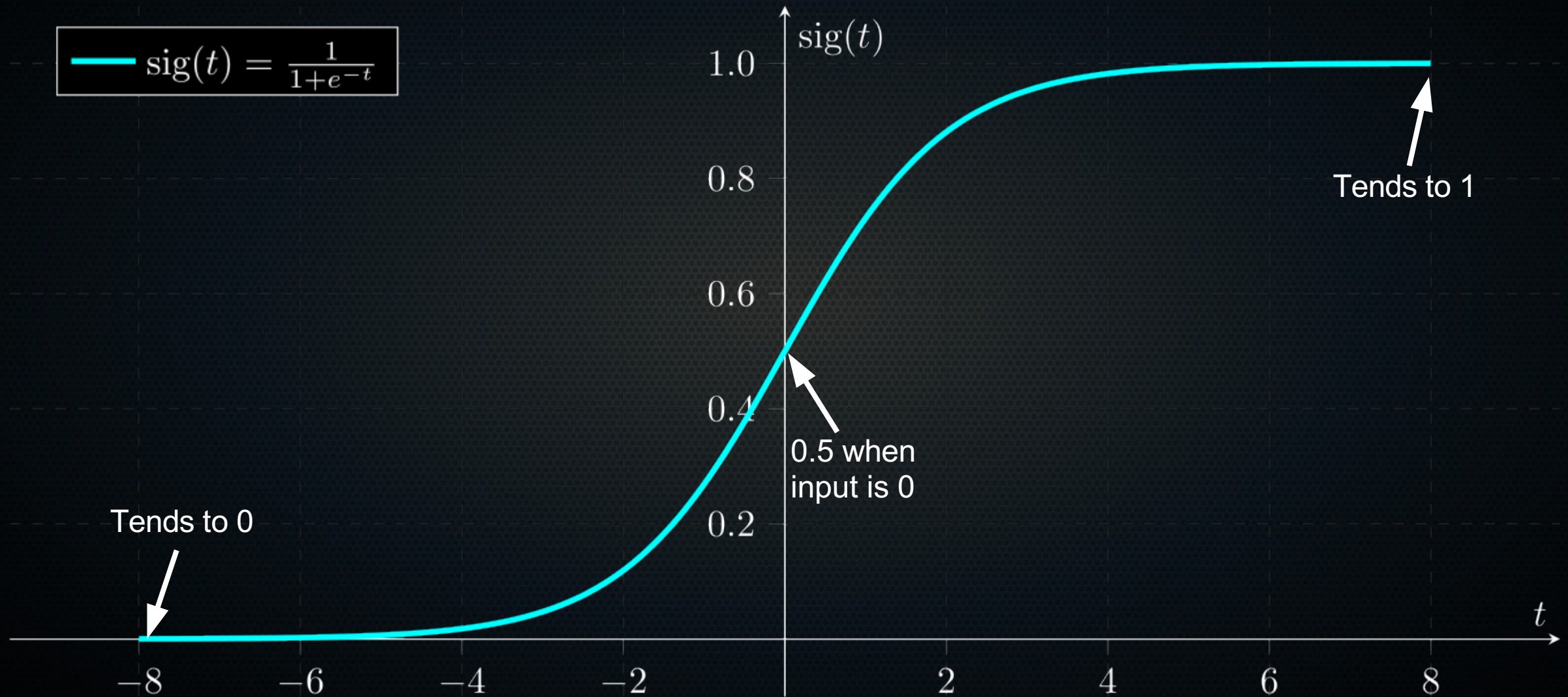
- What does that look like? $z \rightarrow \infty, e^{-z} \rightarrow 0$

$$z \rightarrow -\infty, e^{-z} \rightarrow \infty$$

What does the sigmoid function look like?

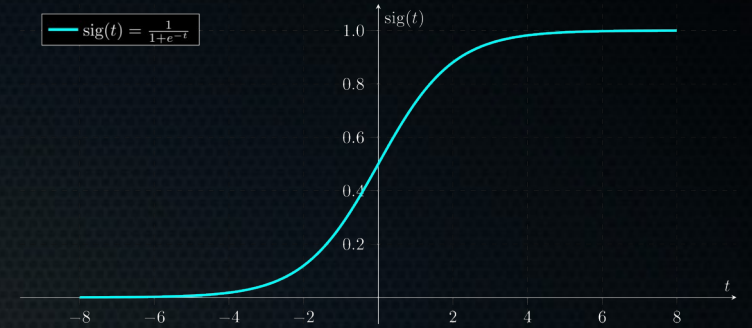


What does the sigmoid function look like?



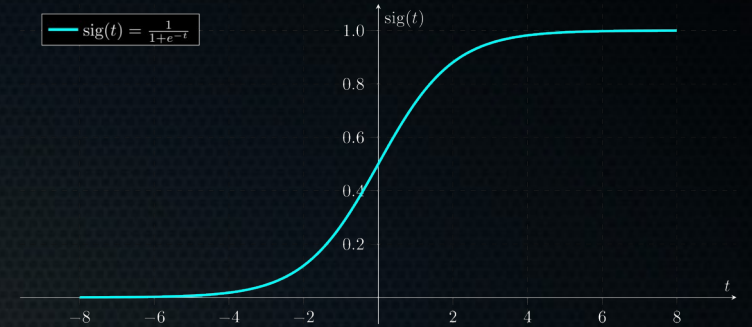
Sigmoid or logistic function

- How do we work with this? $h_{\theta}(x) = \frac{1}{1 + e^{-(\theta^T \cdot x)}}$



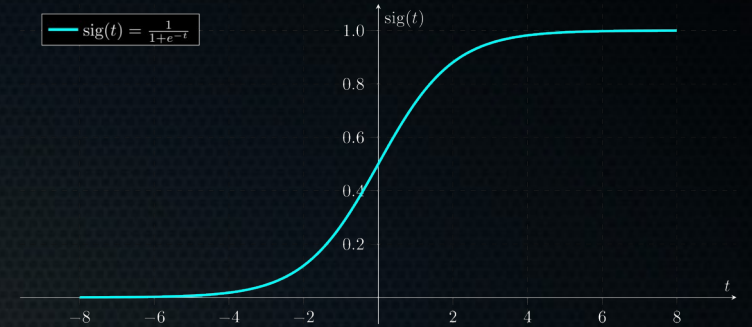
Sigmoid or logistic function

- How do we work with this? $h_{\theta}(x) = \frac{1}{1 + e^{-(\theta^T \cdot x)}}$
 - Interpret outcome of $h_{\theta}(x)$ as probability that class = 1 given the features.



Sigmoid or logistic function

- How do we work with this? $h_{\theta}(x) = \frac{1}{1 + e^{-(\theta^T \cdot x)}}$



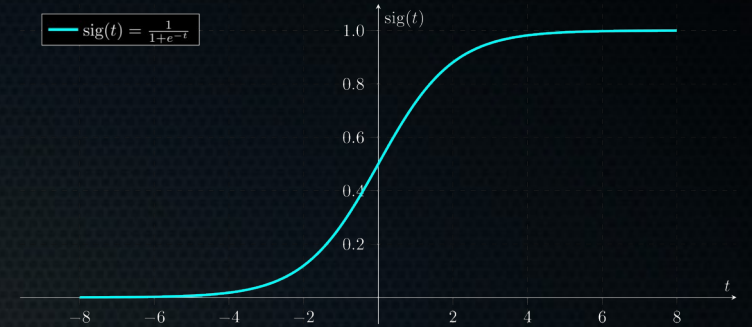
- Interpret outcome of $h_{\theta}(x)$ as probability that class = 1 given the features. Example:

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{Tumor size} \\ \text{Neovascularisation level} \end{bmatrix}$$

$h_{\theta}(x) = 0.8 \longrightarrow$ 80% chance of tumor being malignant

Sigmoid or logistic function

- How do we work with this? $h_{\theta}(x) = \frac{1}{1 + e^{-(\theta^T \cdot x)}}$



- Interpret outcome of $h_{\theta}(x)$ as probability that class = 1 given the features. Example:

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{Tumor size} \\ \text{Neovascularisation level} \end{bmatrix}$$

$h_{\theta}(x) = 0.8 \longrightarrow$ 80% chance of tumor being malignant (class 1)
100% - 80% \rightarrow 20 % chance of being benign (class 0)

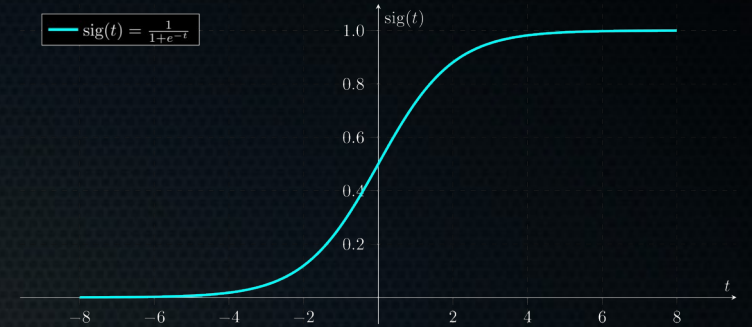
Sigmoid or logistic function

• How do we work with this? $h_{\theta}(x) = \frac{1}{1 + e^{-(\theta^T \cdot x)}}$

- Interpret outcome of $h_{\theta}(x)$ as probability that class = 1 given the features.
- Formally:

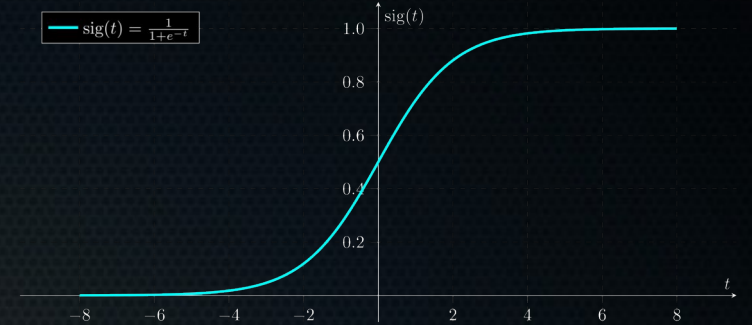
$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta^T \cdot x)}} = p(y=1 | x; \theta)$$

$$p(y=0 | x; \theta) = 1 - h_{\theta}(x)$$



Sigmoid or logistic function

- How do we work with this? $h_{\theta}(x) = \frac{1}{1 + e^{-(\theta^T \cdot x)}}$



- Interpret outcome of $h_{\theta}(x)$ as probability that class = 1 given the features.

- Formally:

$$\left. \begin{aligned} h_{\theta}(x) &= \frac{1}{1 + e^{-(\theta^T \cdot x)}} = p(y=1|x;\theta) \\ p(y=0|x;\theta) &= 1 - h_{\theta}(x) \end{aligned} \right\} \frac{p(y=1)}{1 - p(y=1)}$$

Log odds

Coefficients:

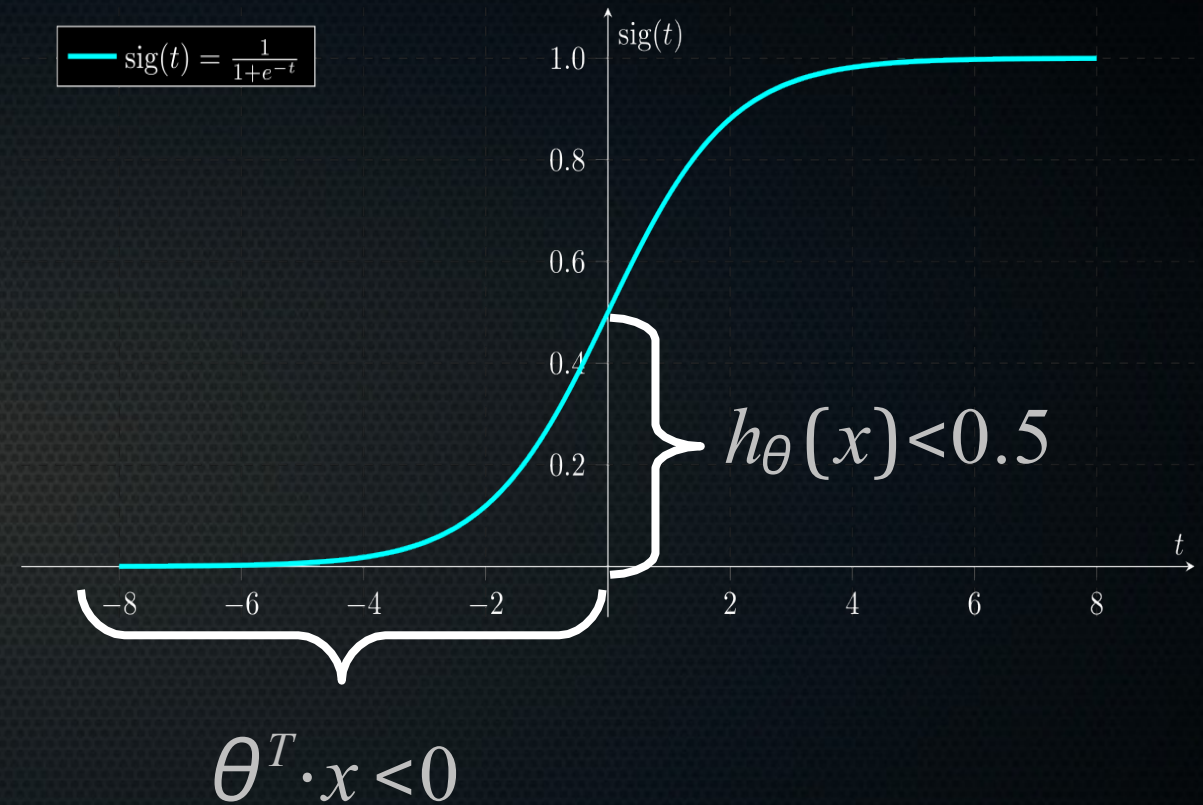
	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-8.922213	0.879485	-10.145	< 2e-16 ***
pregnant	0.106275	0.039885	2.665	0.00771 **
glucose	0.036555	0.004546	8.041	8.88e-16 ***
pressure	-0.008019	0.006230	-1.287	0.19801
triceps	0.001909	0.008258	0.231	0.81721
insulin	-0.001394	0.001098	-1.269	0.20445
mass	0.082665	0.017927	4.611	4.00e-06 ***
pedigree	0.901691	0.374350	2.409	0.01601 *
age	0.020417	0.010854	1.881	0.05995 .

$\text{logOdds}(\text{diabetes}) = -8.9 + (0.106 * \text{pregnant}) + (0.037 * \text{glucose}) + \dots + (0.02 * \text{age})$



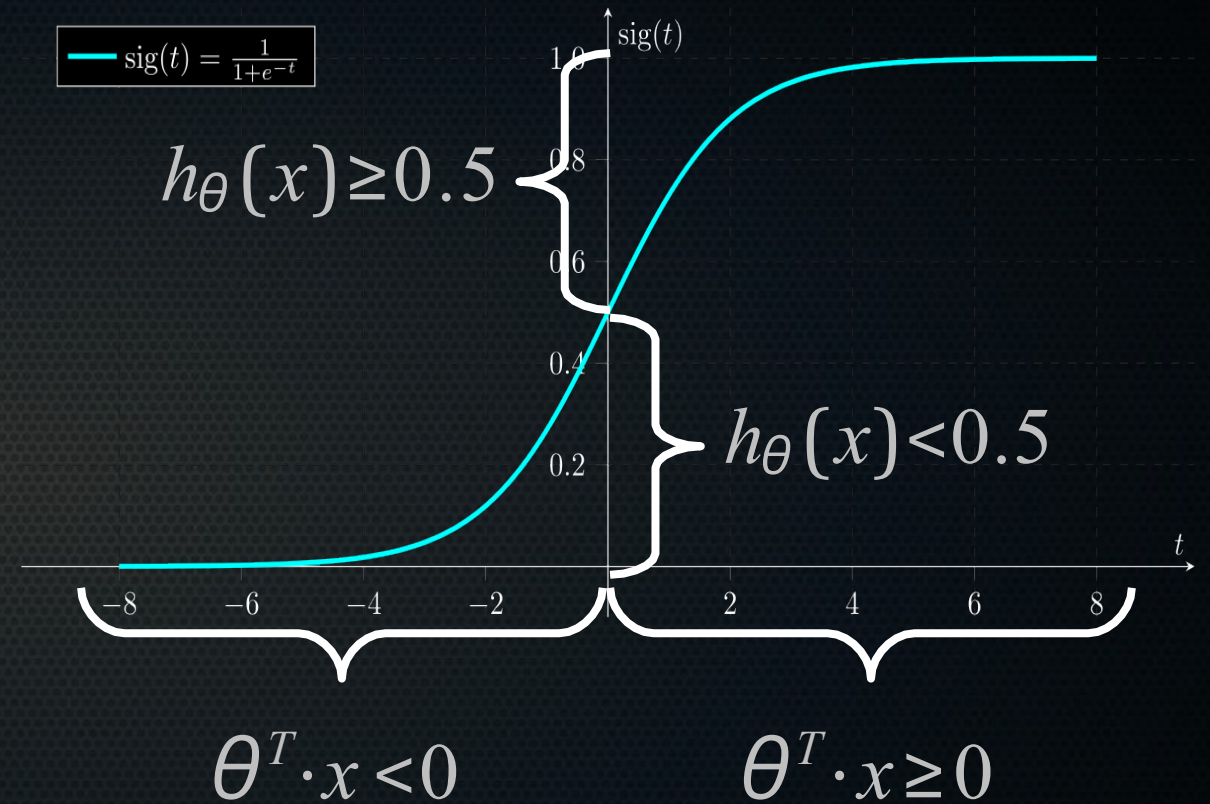
Decision boundary

- Threshold:



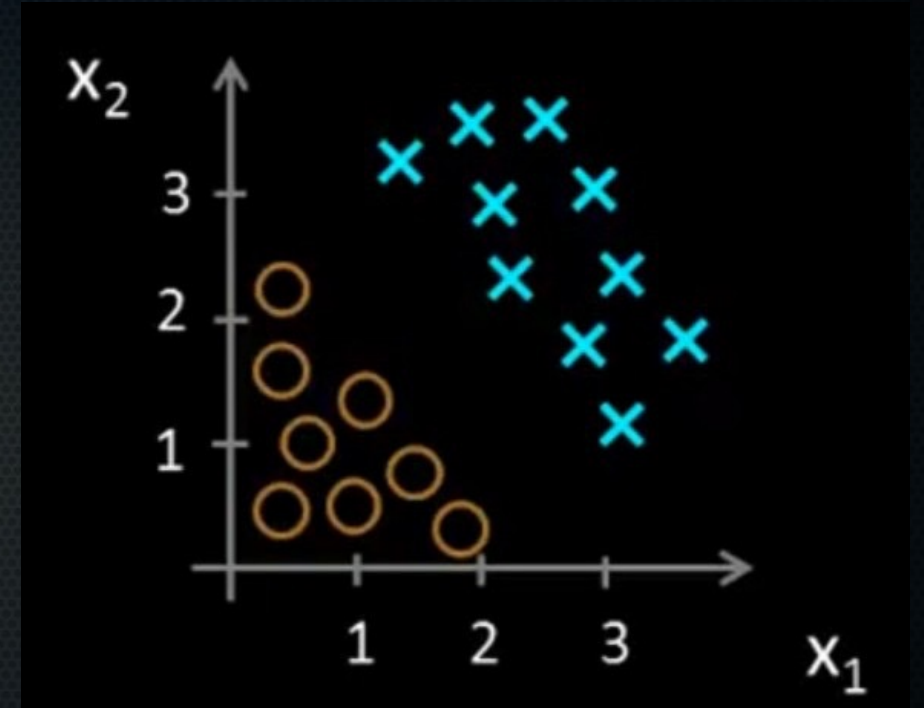
Decision boundary

- Threshold:



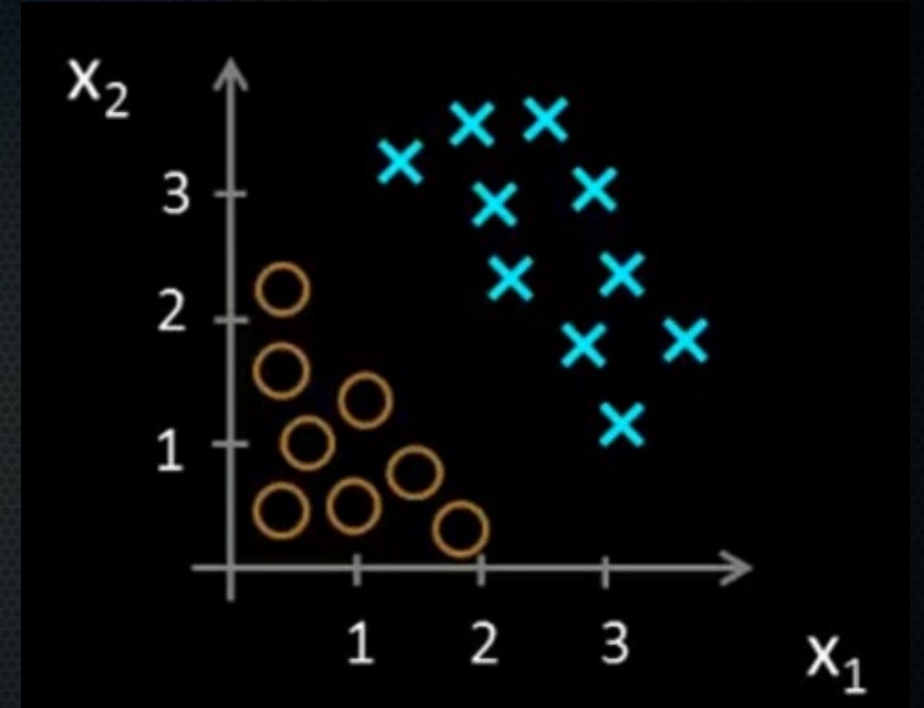
Decision boundary

- How does it look? $g(z) = \frac{1}{1 + e^{-z}}$
 $h_{\theta}(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$



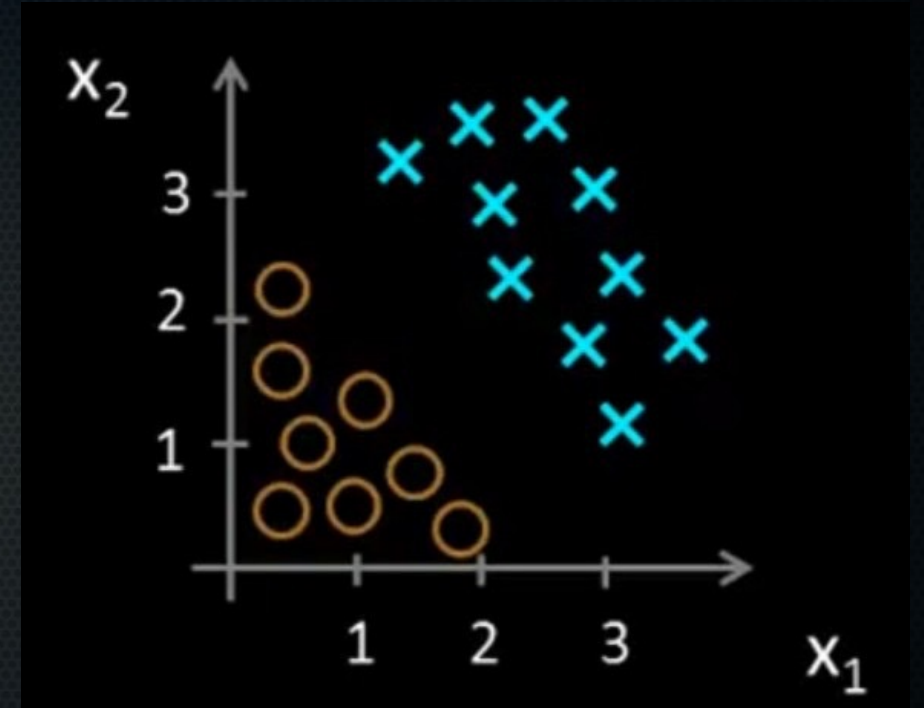
Decision boundary

- How does it look? $g(z) = \frac{1}{1 + e^{-z}}$
 $h_{\theta}(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$
 $\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$



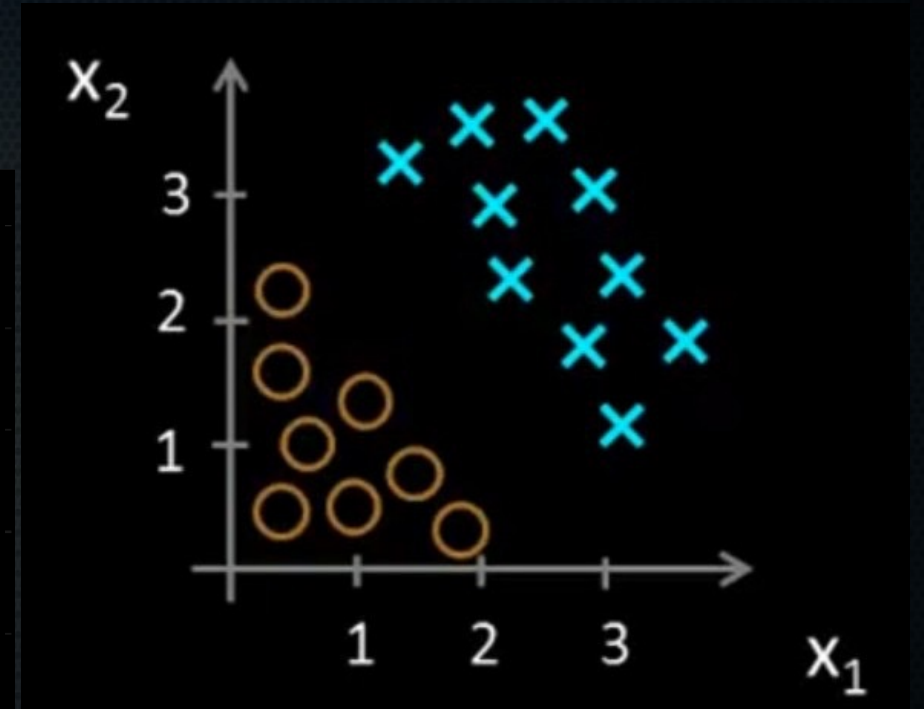
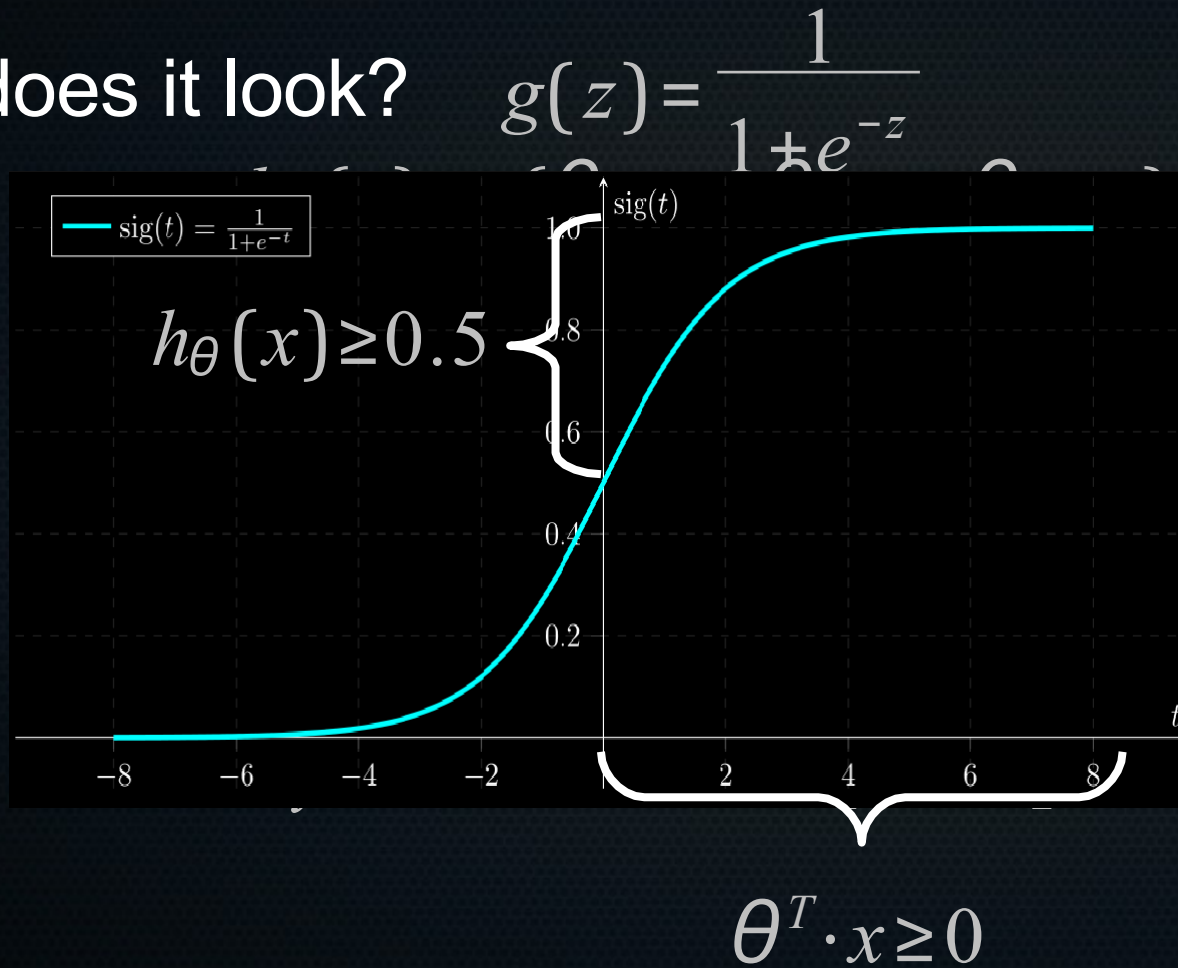
Decision boundary

- How does it look? $g(z) = \frac{1}{1 + e^{-z}}$
 $h_{\theta}(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$
 $\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$
 $y = 1$ if $-3 \cdot 1 + 1 \cdot x_1 + 1 \cdot x_2 \geq 0$



Decision boundary

- How does it look?



Decision boundary

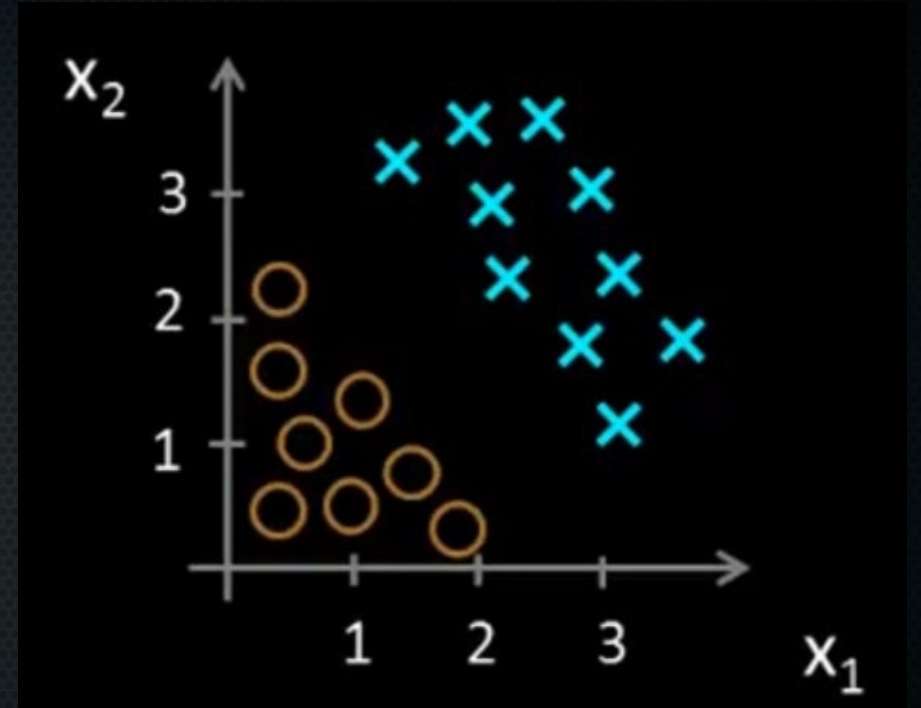
- How does it look? $g(z) = \frac{1}{1 + e^{-z}}$
 $h_{\theta}(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$

$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

$$y = 1 \text{ if } -3 \cdot 1 + 1 \cdot x_1 + 1 \cdot x_2 \geq 0$$

$$-3 \cdot 1 + \cancel{1 \cdot x_1} + \cancel{1 \cdot x_2} \geq 0$$

$$x_1 + x_2 = 3$$



Decision boundary

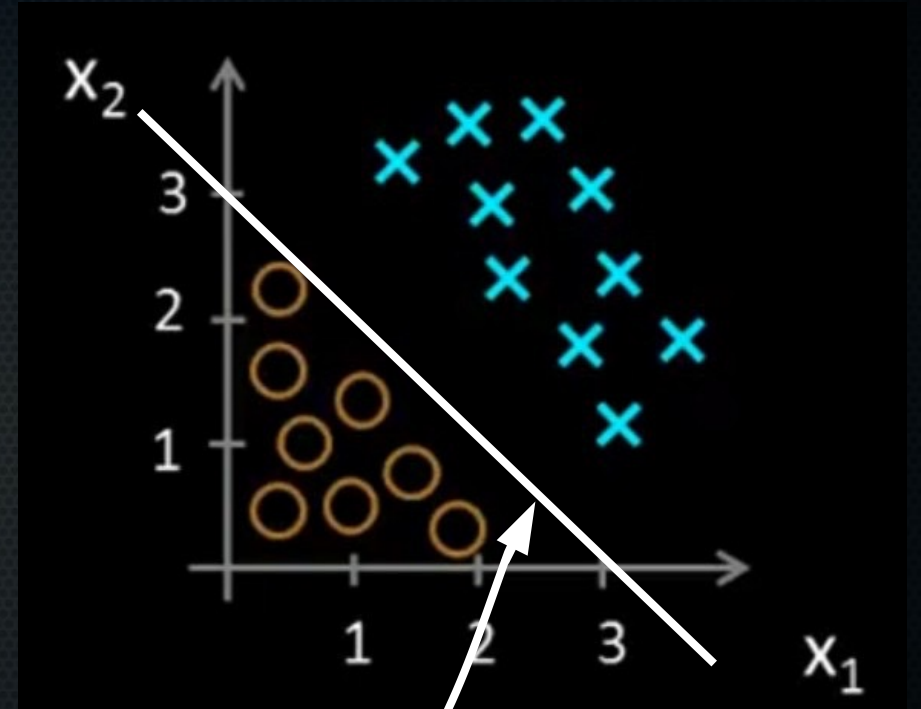
- How does it look? $g(z) = \frac{1}{1 + e^{-z}}$
 $h_{\theta}(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$

$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

$$y = 1 \text{ if } -3 \cdot 1 + 1 \cdot x_1 + 1 \cdot x_2 \geq 0$$

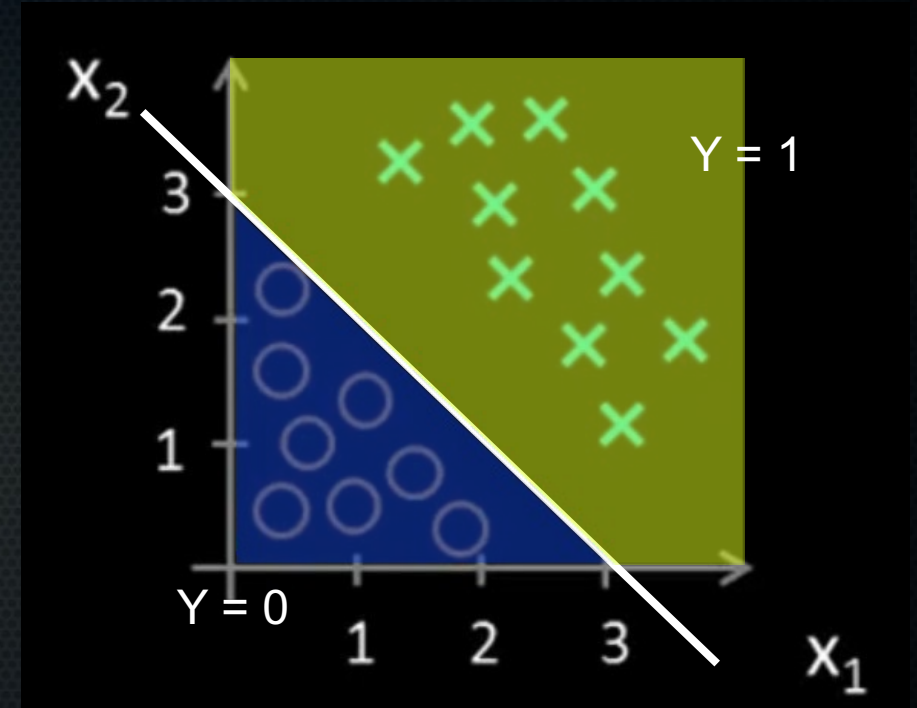
$$-3 \cdot 1 + \cancel{1 \cdot x_1} + \cancel{1 \cdot x_2} \geq 0$$

$$x_1 + x_2 = 3$$



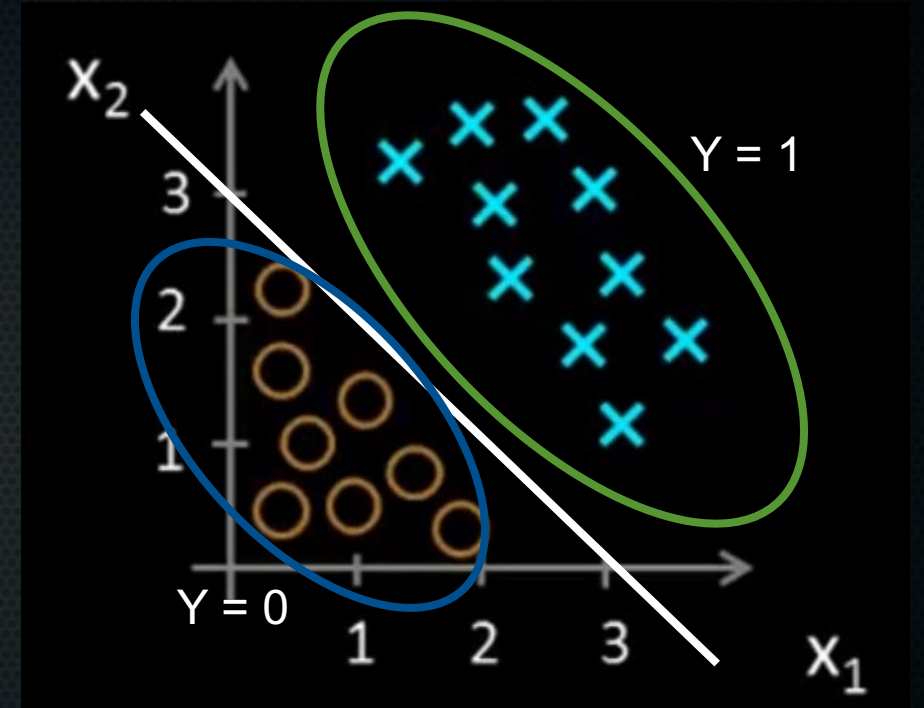
Decision boundary

- How does it look? $g(z) = \frac{1}{1 + e^{-z}}$
 $h_{\theta}(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$
 $\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$
 $y = 1$ if $-3 \cdot 1 + 1 \cdot x_1 + 1 \cdot x_2 \geq 0$



Decision boundary

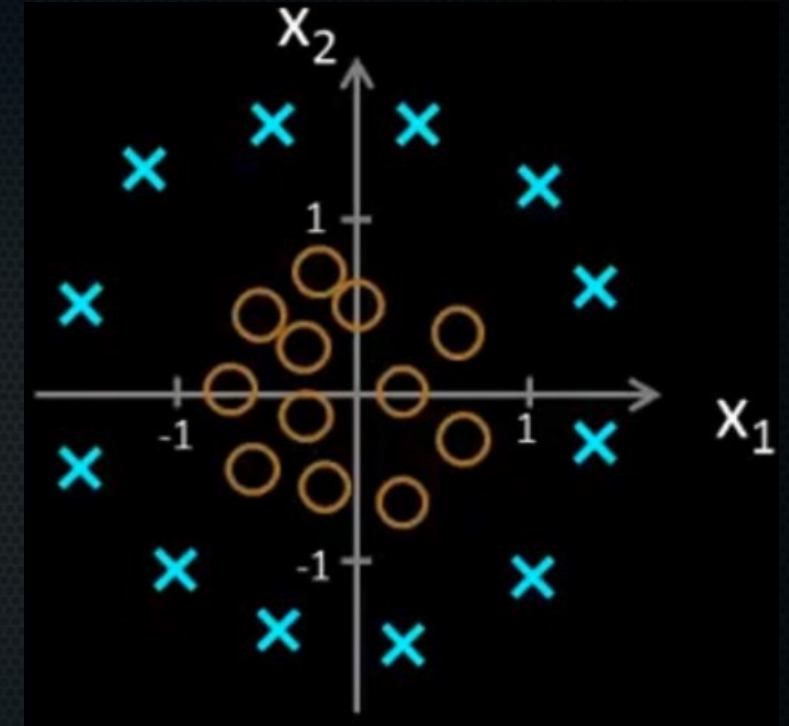
- How does it look? $g(z) = \frac{1}{1 + e^{-z}}$
 $h_{\theta}(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$
 $\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$
 $y = 1$ if $-3 \cdot 1 + 1 \cdot x_1 + 1 \cdot x_2 \geq 0$



Decision boundary

- How does it look? $g(z) = \frac{1}{1 + e^{-z}}$
 $h_{\theta}(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$

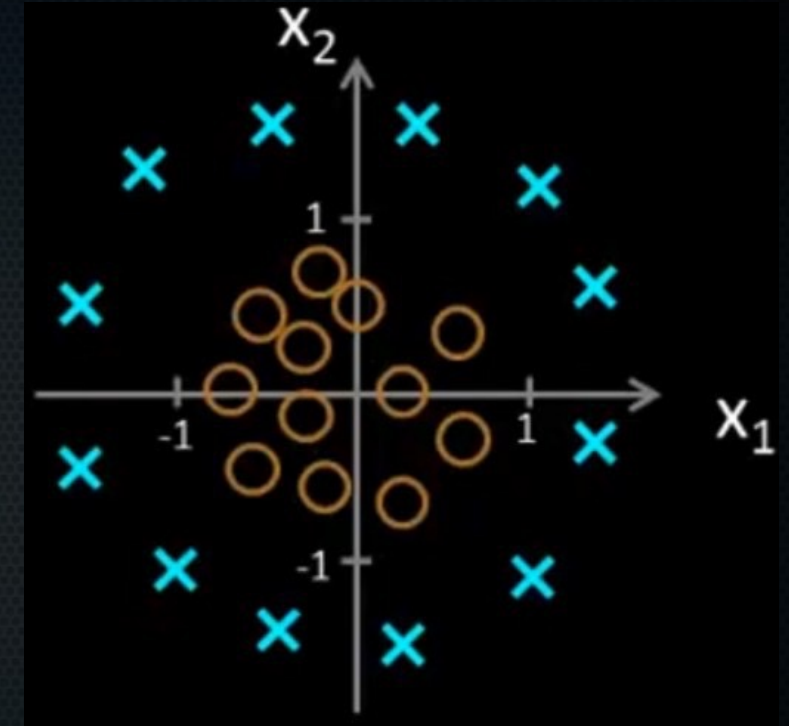
$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$



$$y = 1 \text{ if } -3 \cdot 1 + 1 \cdot x_1 + 1 \cdot x_2 \geq 0$$

Non-linear decision boundary

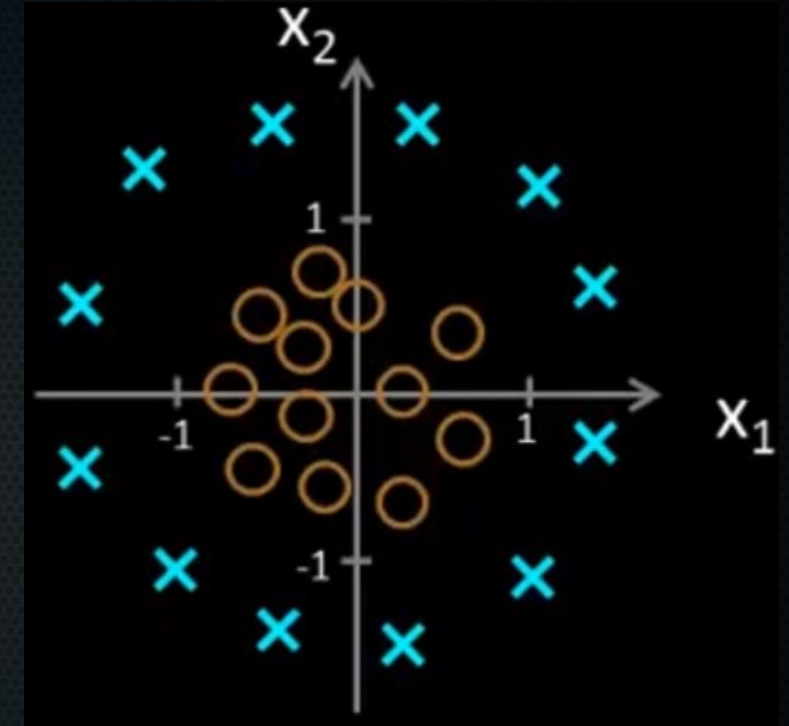
- How does it look? $g(z) = \frac{1}{1 + e^{-z}}$
 $h_{\theta}(x) = g(\theta_0 x_0, \theta_1 x_1, \theta_2 x_2, \theta_3 x_1^2, \theta_4 x_2^2)$
- Add two polynomial features



Non-linear decision boundary

- How does it look? $g(z) = \frac{1}{1 + e^{-z}}$
 $h_{\theta}(x) = g(\theta_0 x_0, \theta_1 x_1, \theta_2 x_2, \theta_3 x_1^2, \theta_4 x_2^2)$
- Add two polynomial features

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$



Non-linear decision boundary

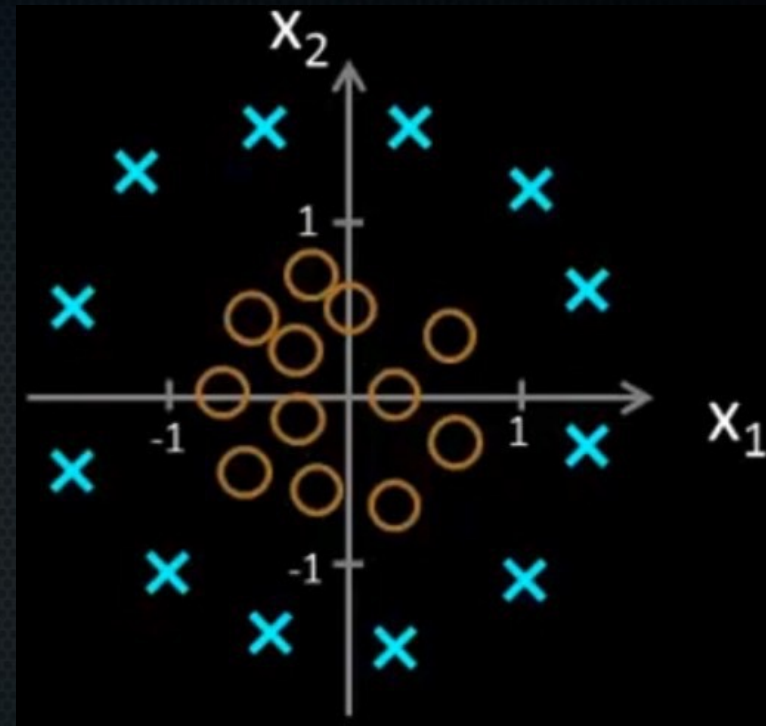
- How does it look? $g(z) = \frac{1}{1 + e^{-z}}$

$$h_{\theta}(x) = g(\theta_0 x_0, \theta_1 x_1, \theta_2 x_2, \theta_3 x_1^2, \theta_4 x_2^2)$$

- Add two polynomial features

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

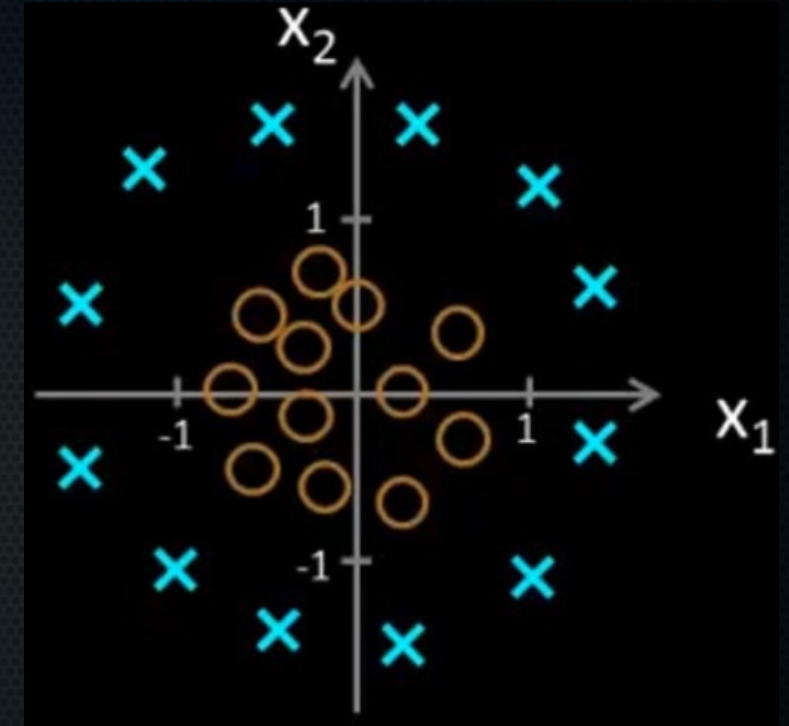
Can you work out what the decision boundary will be?



Non-linear decision boundary

- How does it look? $g(z) = \frac{1}{1 + e^{-z}}$
 $h_{\theta}(x) = g(\theta_0 x_0, \theta_1 x_1, \theta_2 x_2, \theta_3 x_1^2, \theta_4 x_2^2)$
- Add two polynomial features

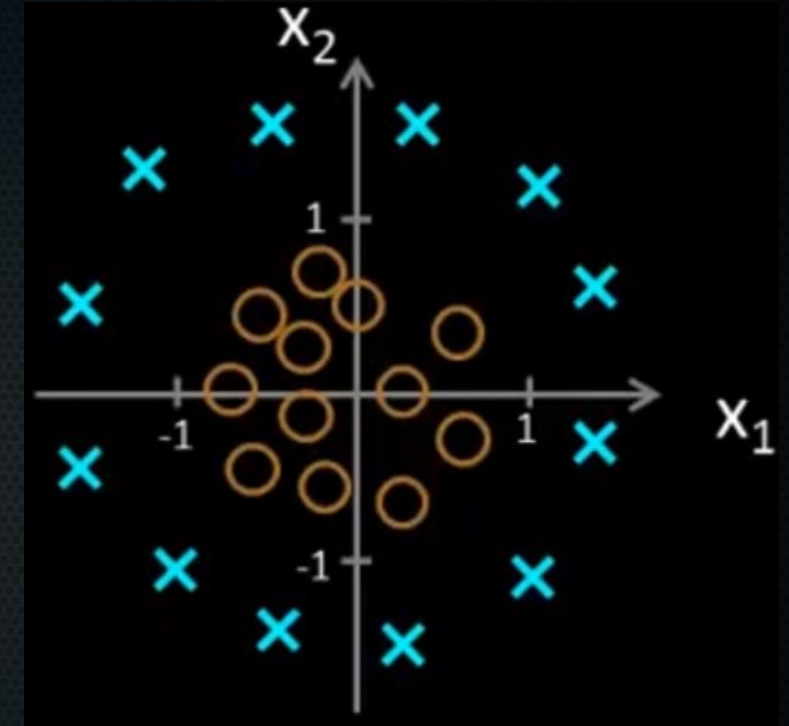
$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \longrightarrow \begin{matrix} y=1 \text{ if} \\ -1 + x_1^2 + x_2^2 \geq 0 \end{matrix}$$



Non-linear decision boundary

- How does it look? $g(z) = \frac{1}{1 + e^{-z}}$
 $h_{\theta}(x) = g(\theta_0 x_0, \theta_1 x_1, \theta_2 x_2, \theta_3 x_1^2, \theta_4 x_2^2)$
- Add two polynomial features

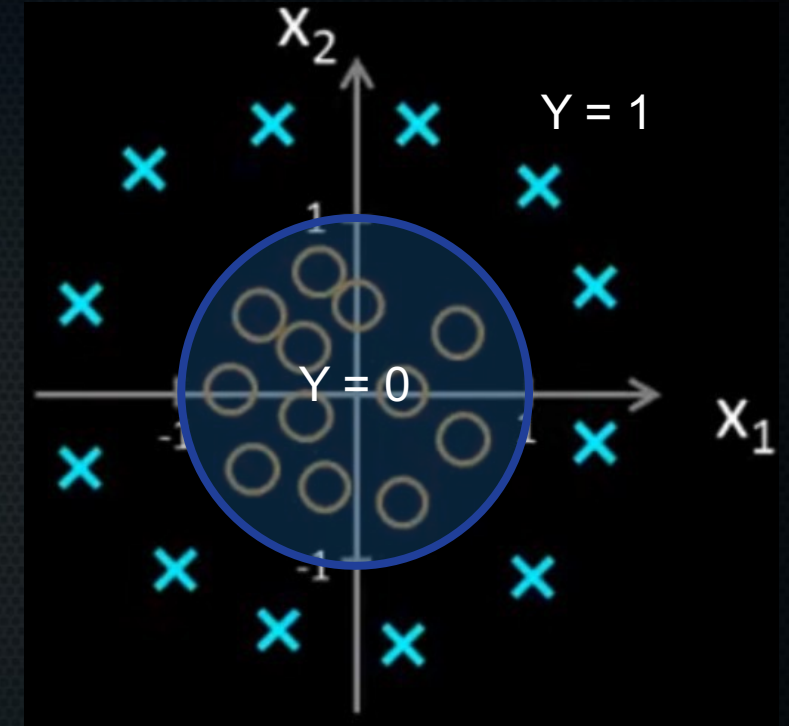
$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \longrightarrow -1 + x_1^2 + x_2^2 \geq 0$$
$$\downarrow$$
$$x_1^2 + x_2^2 \geq 1$$



Non-linear decision boundary

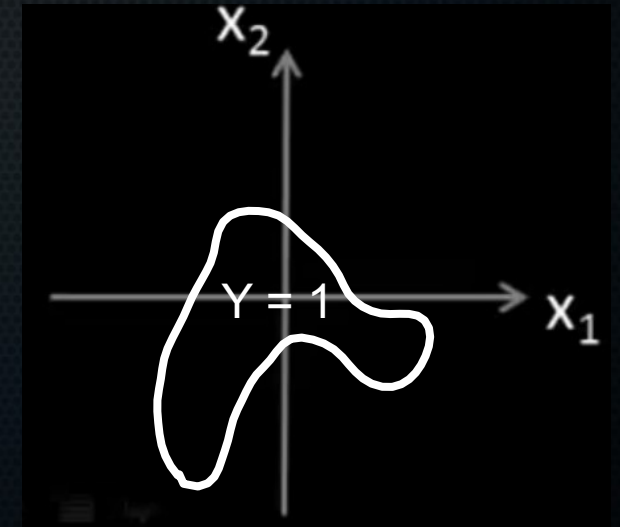
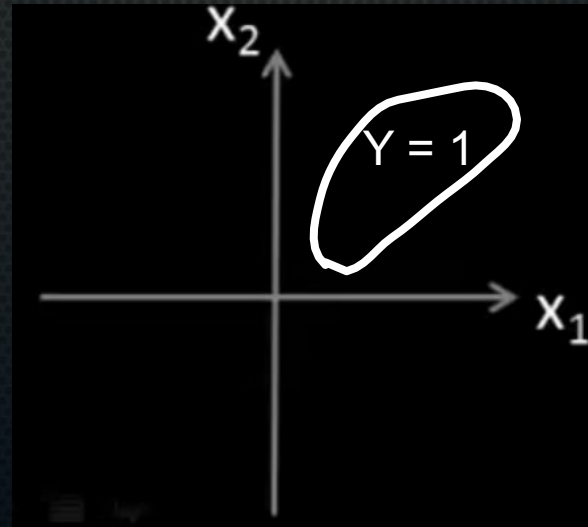
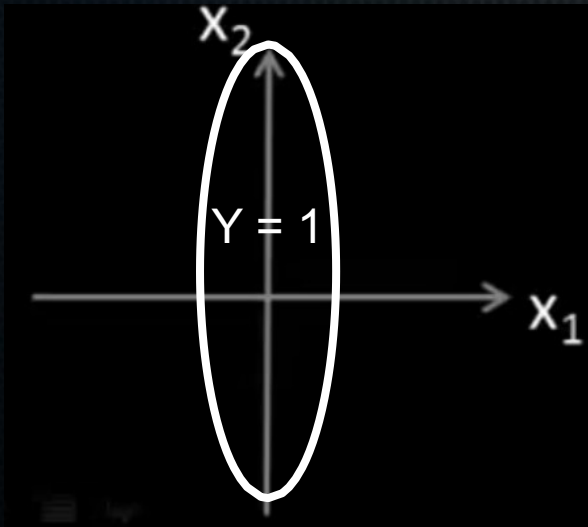
- How does it look? $g(z) = \frac{1}{1 + e^{-z}}$
 $h_{\theta}(x) = g(\theta_0 x_0, \theta_1 x_1, \theta_2 x_2, \theta_3 x_1^2, \theta_4 x_2^2)$
- Add two polynomial features

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \longrightarrow -1 + x_1^2 + x_2^2 \geq 0$$
$$\downarrow$$
$$x_1^2 + x_2^2 \geq 1$$



Non-linear decision boundary

- How does it look?
- If you add more and higher-order polynomial features, you can get complex boundaries:



So how do we get theta's?

- Need a cost function

So how do we get theta's?

- Need a cost function

- Before: $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

So how do we get theta's?

- Need a cost function

- Before: $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$




$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

So how do we get theta's?

- Need a cost function

- Before: $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left(\frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right)$$


$$\text{Cost}(x) = \frac{1}{2} (h_{\theta}(x) - y)^2$$

So how do we get theta's?

- Need a cost function

- Before: $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left\{ \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right\}$

$\text{Cost}(x) = \frac{1}{2} (h_{\theta}(x) - y)^2$

$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$

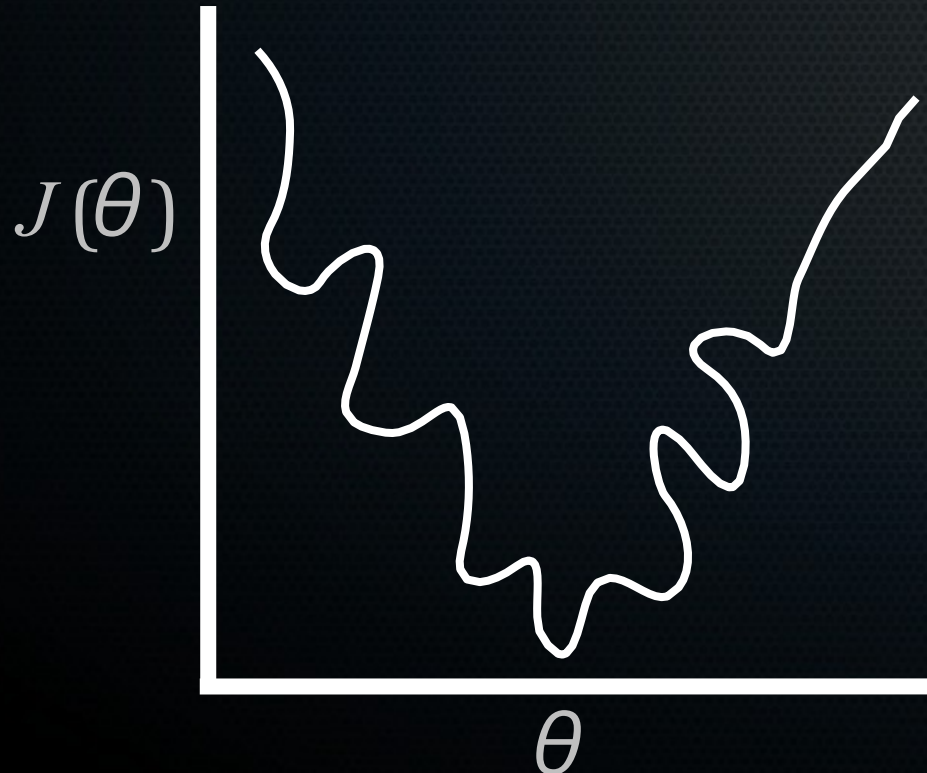
So how do we get theta's?

- Need a cost function $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$ $\text{Cost}(x) = \frac{1}{2} (h_{\theta}(x) - y)^2$
- Why not MSE? \rightarrow not convex

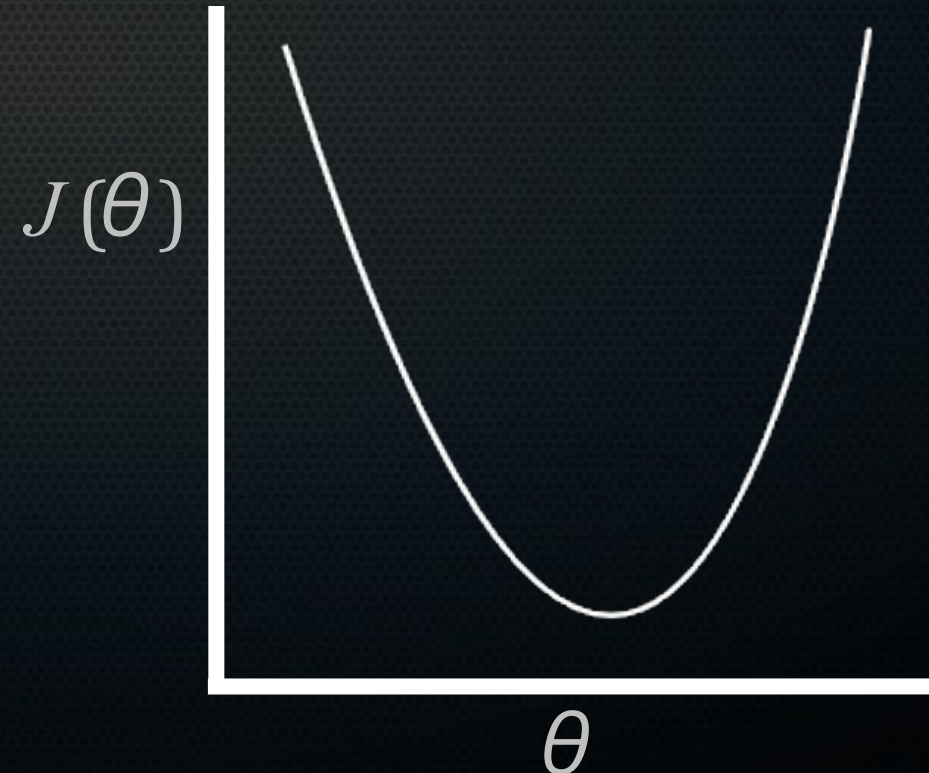
So how do we get theta's?

- Need a cost function $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$ $\text{Cost}(x) = \frac{1}{2} (h_{\theta}(x) - y)^2$
- Why not MSE? → not convex

non-convex



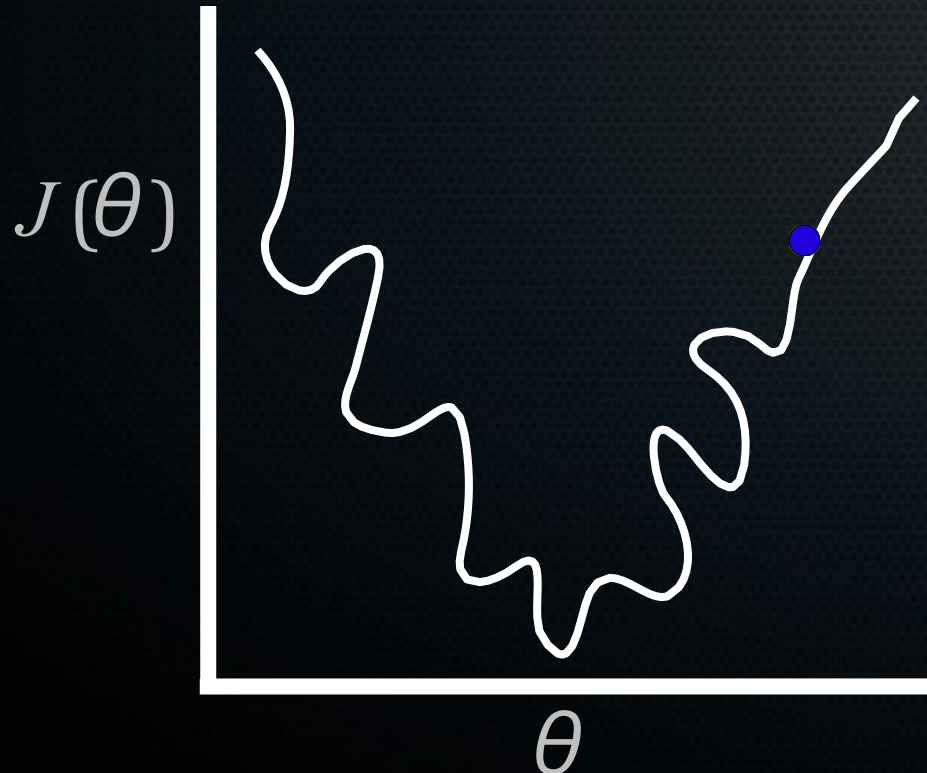
convex



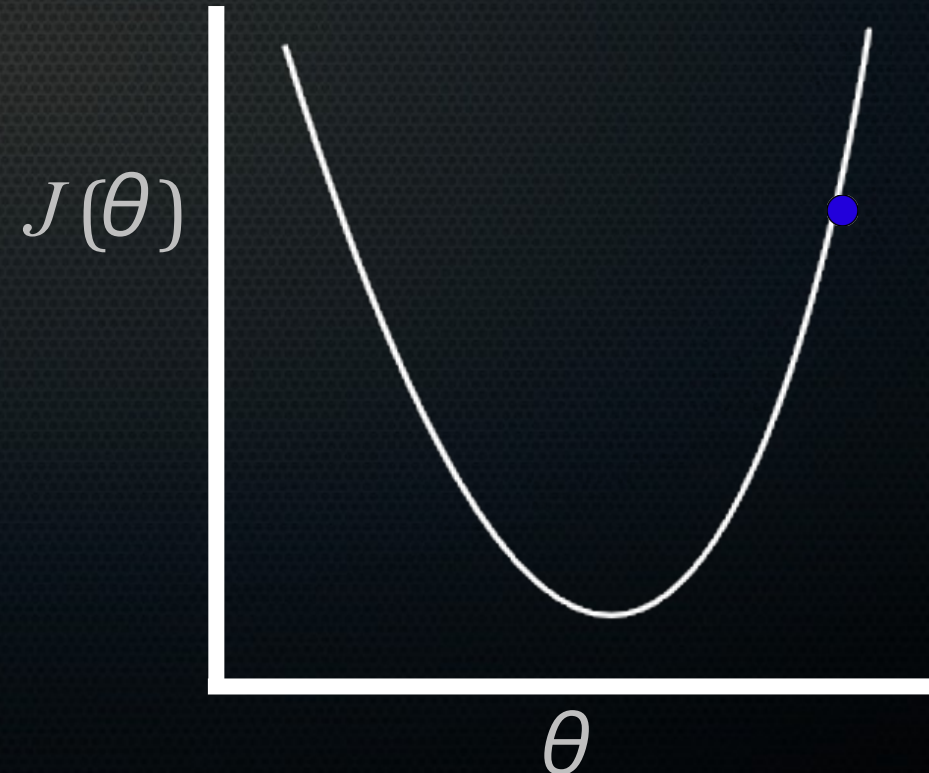
So how do we get theta's?

- Need a cost function $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$ $\text{Cost}(x) = \frac{1}{2} (h_{\theta}(x) - y)^2$
- Why not MSE? → not convex

non-convex



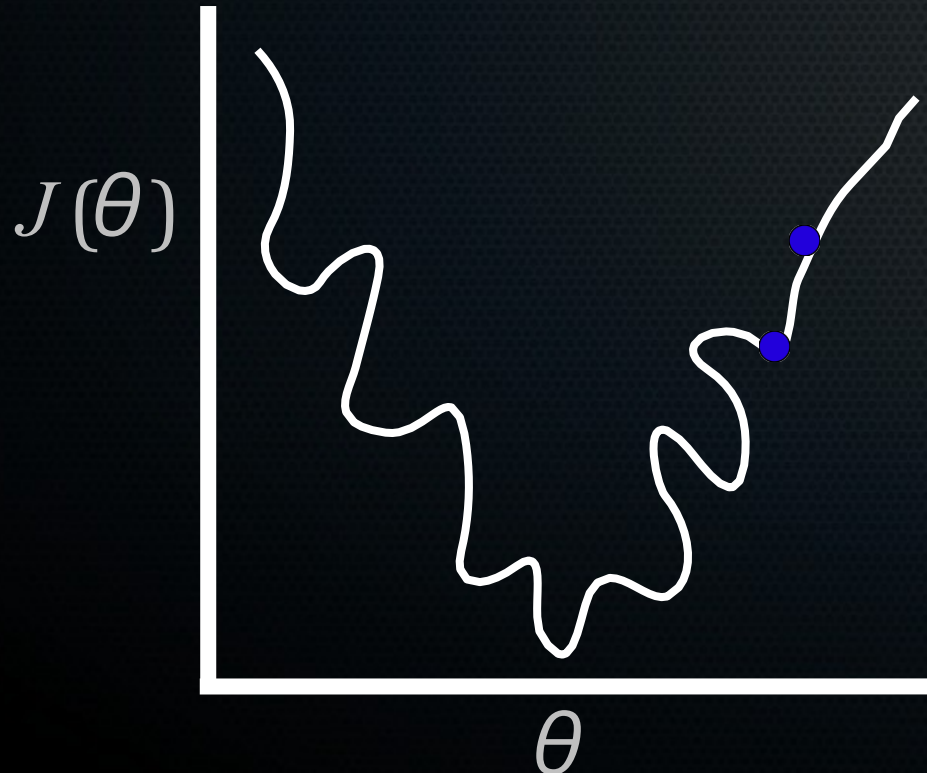
convex



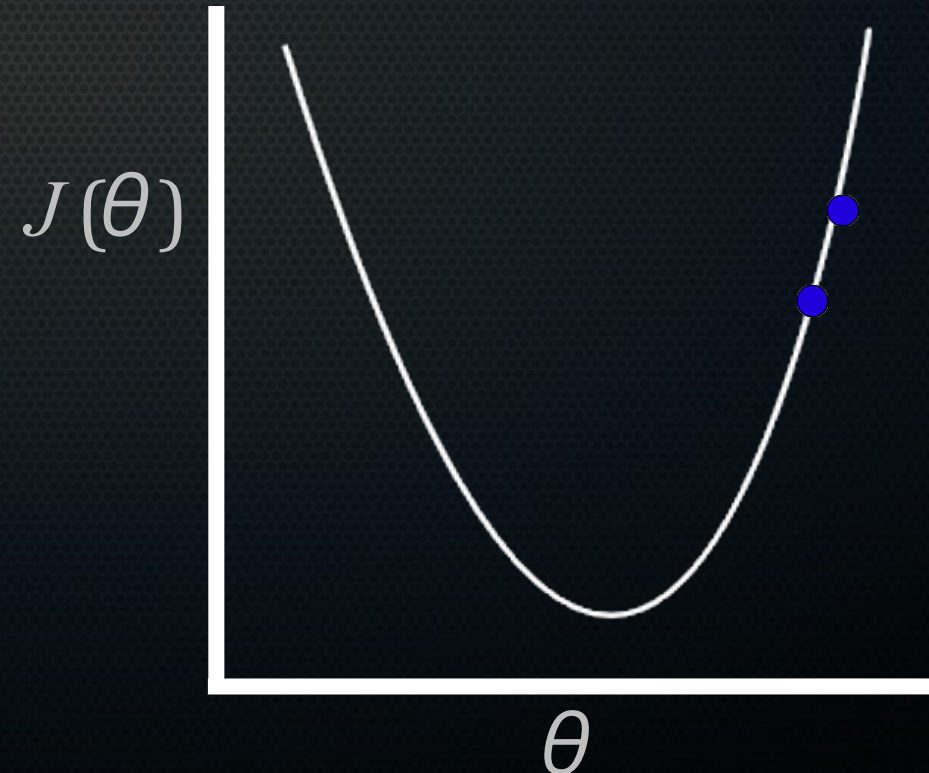
So how do we get theta's?

- Need a cost function $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$ $\text{Cost}(x) = \frac{1}{2} (h_{\theta}(x) - y)^2$
- Why not MSE? \rightarrow not convex

non-convex



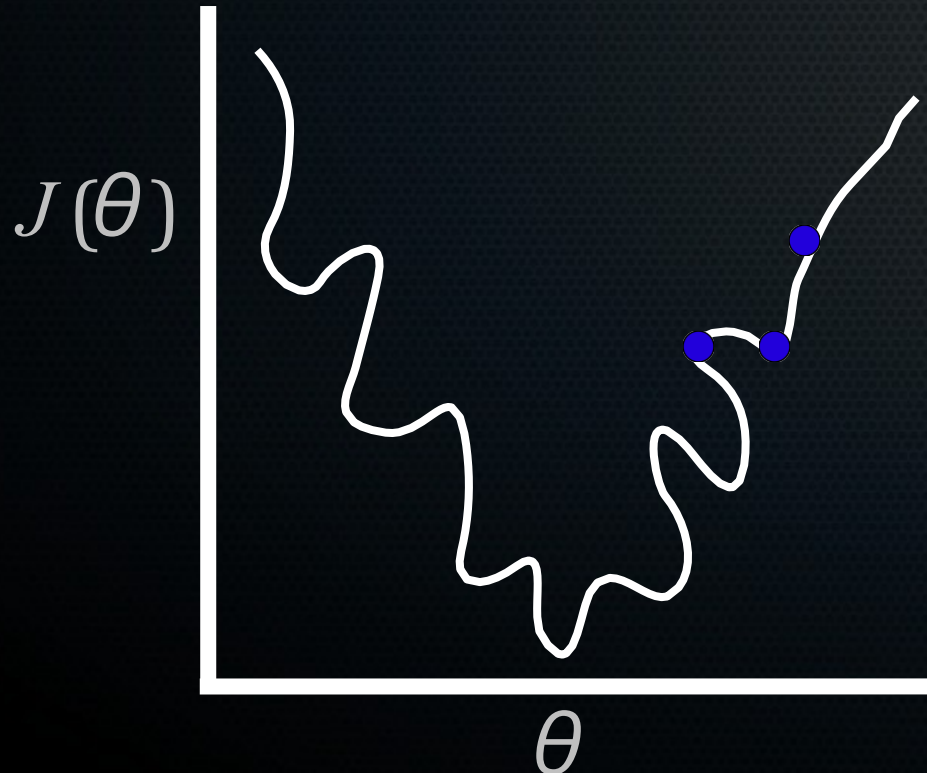
convex



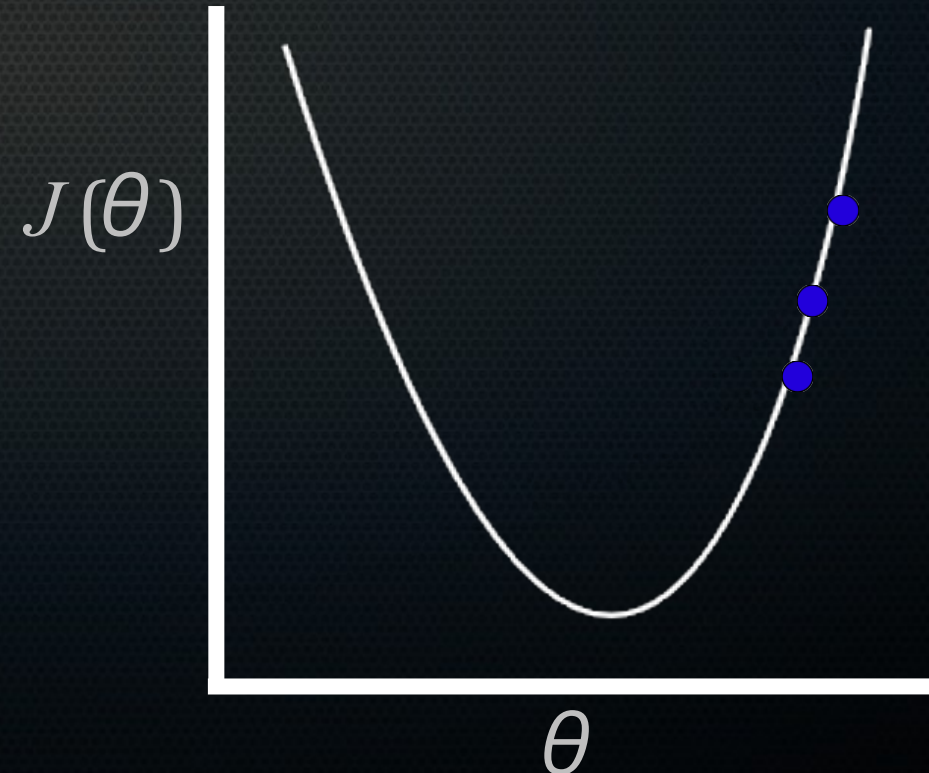
So how do we get theta's?

- Need a cost function $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$ $\text{Cost}(x) = \frac{1}{2} (h_{\theta}(x) - y)^2$
- Why not MSE? \rightarrow not convex

non-convex



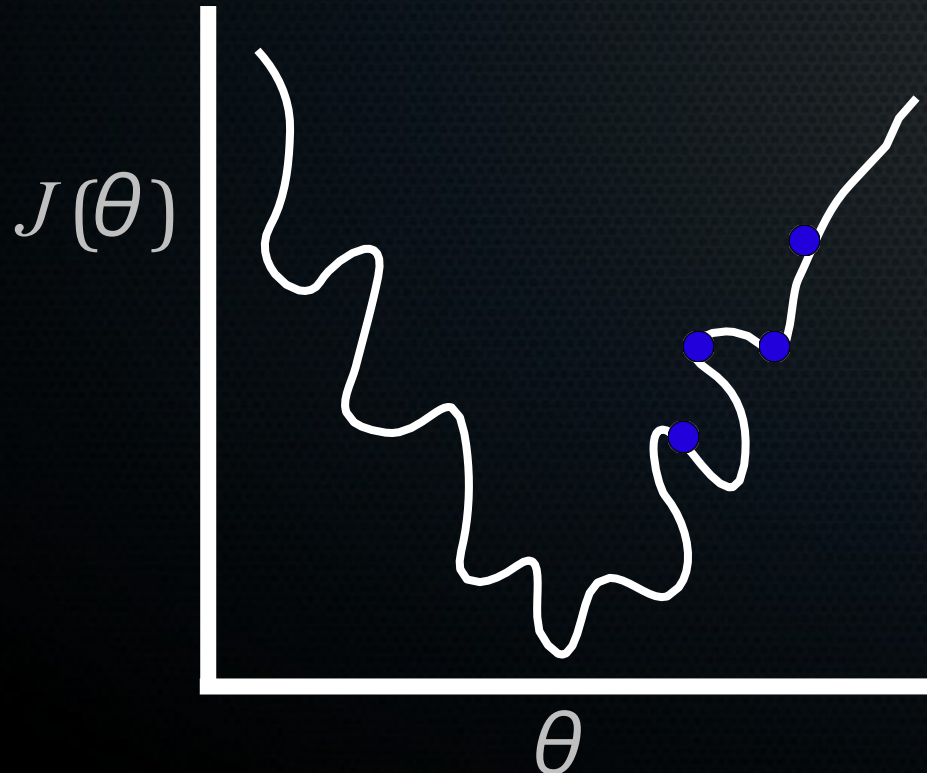
convex



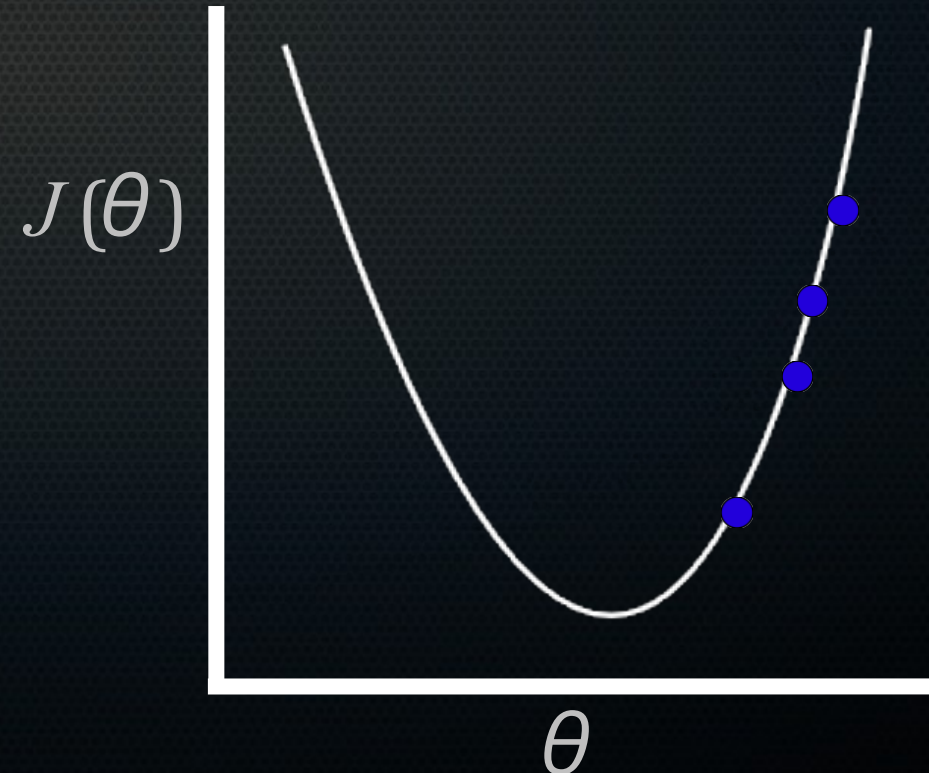
So how do we get theta's?

- Need a cost function $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$ $\text{Cost}(x) = \frac{1}{2} (h_{\theta}(x) - y)^2$
- Why not MSE? \rightarrow not convex

non-convex



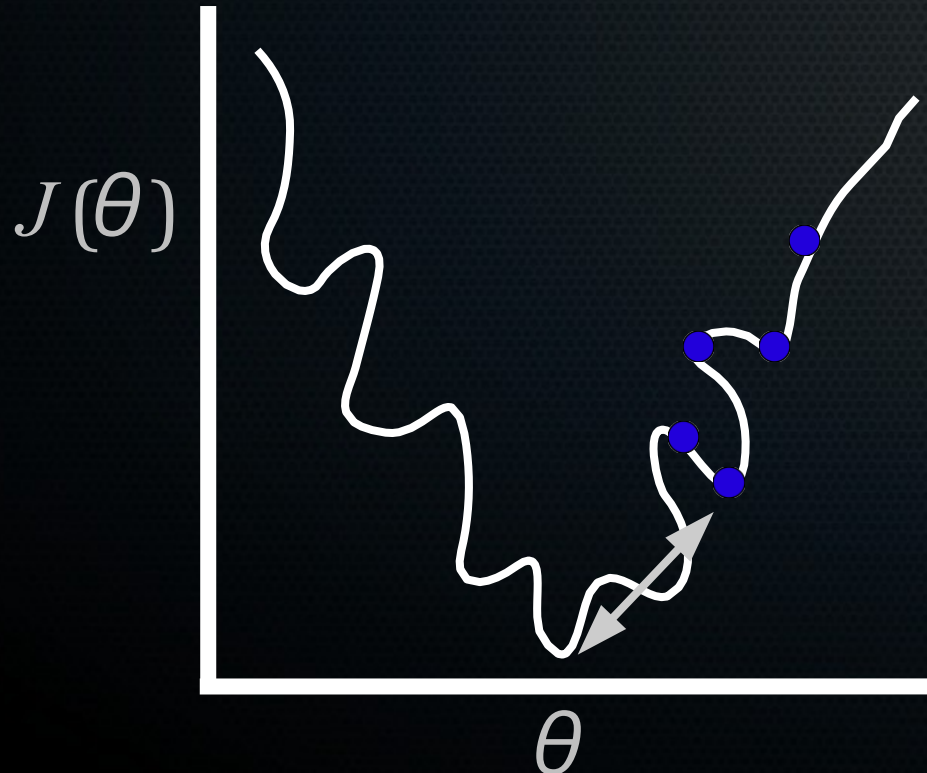
convex



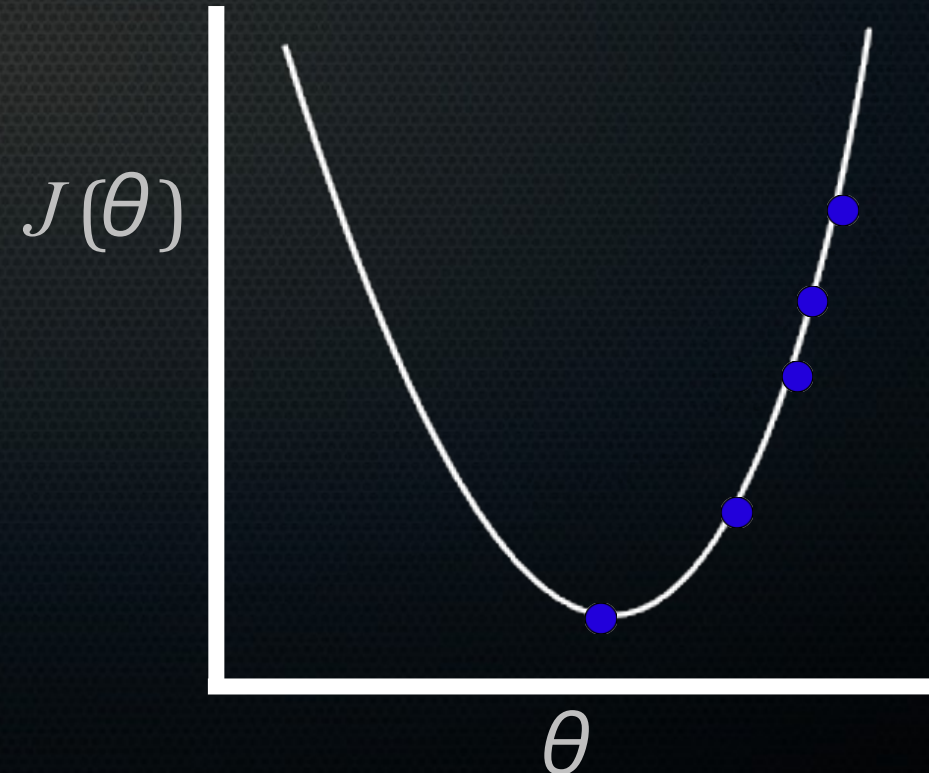
So how do we get theta's?

- Need a cost function $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$ $\text{Cost}(x) = \frac{1}{2} (h_{\theta}(x) - y)^2$
- Why not MSE? \rightarrow not convex

non-convex



convex



So how do we get theta's?

- Need a cost function $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$
 - What then?
- ~~$\text{Cost}(x) = \frac{1}{2} (h_{\theta}(x) - y)^2$~~

So how do we get theta's?

- Need a cost function $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$
- What then?

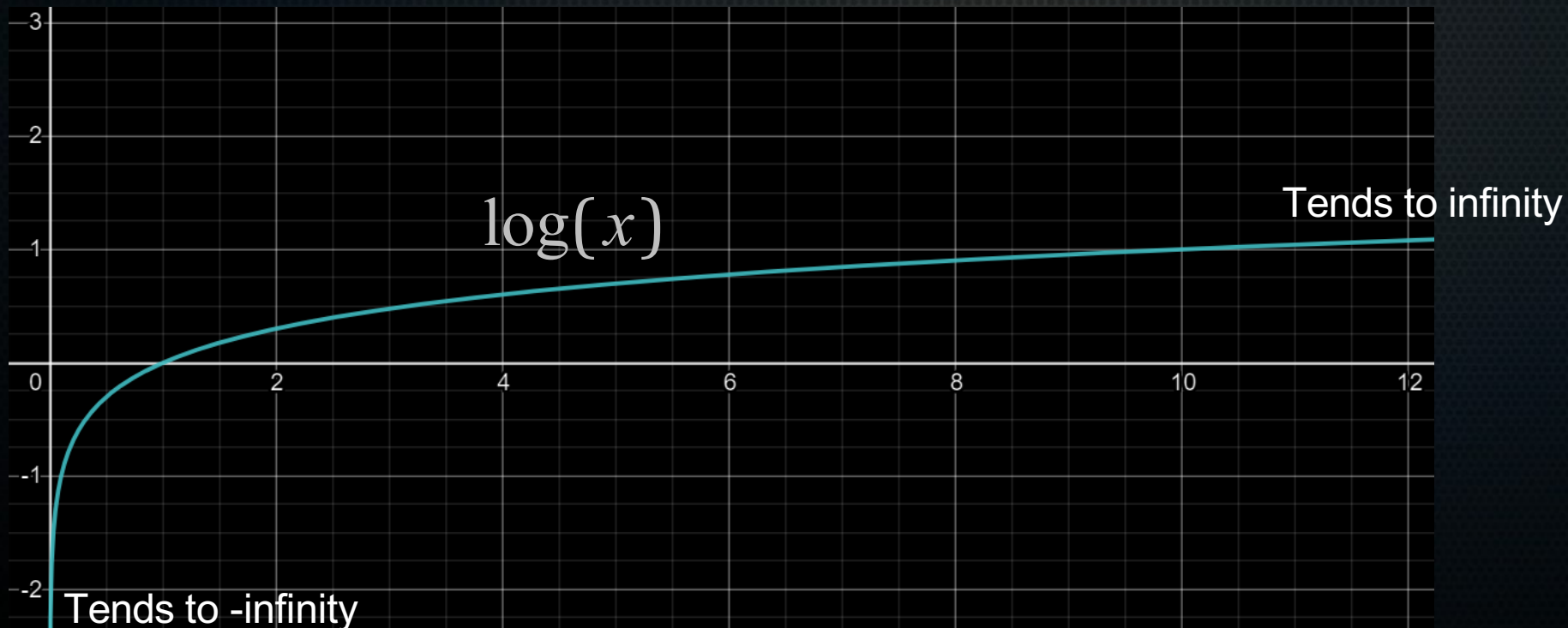
$$\text{Cost}(x) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

~~$$\text{Cost}(x) = \frac{1}{2} (h_{\theta}(x) - y)^2$$~~

What does this function look like?

$$\text{Cost}(x) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

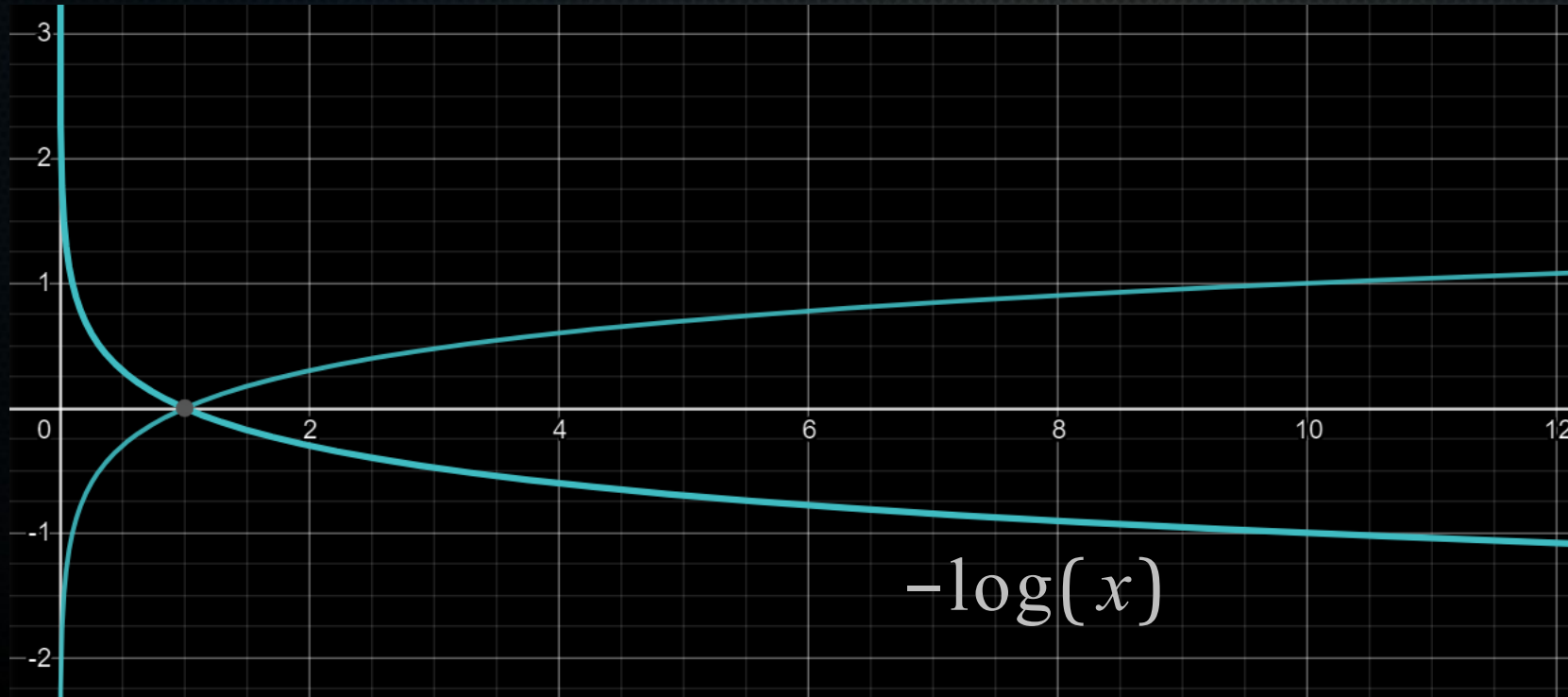
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$$



What does this function look like?

$$\text{Cost}(x) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

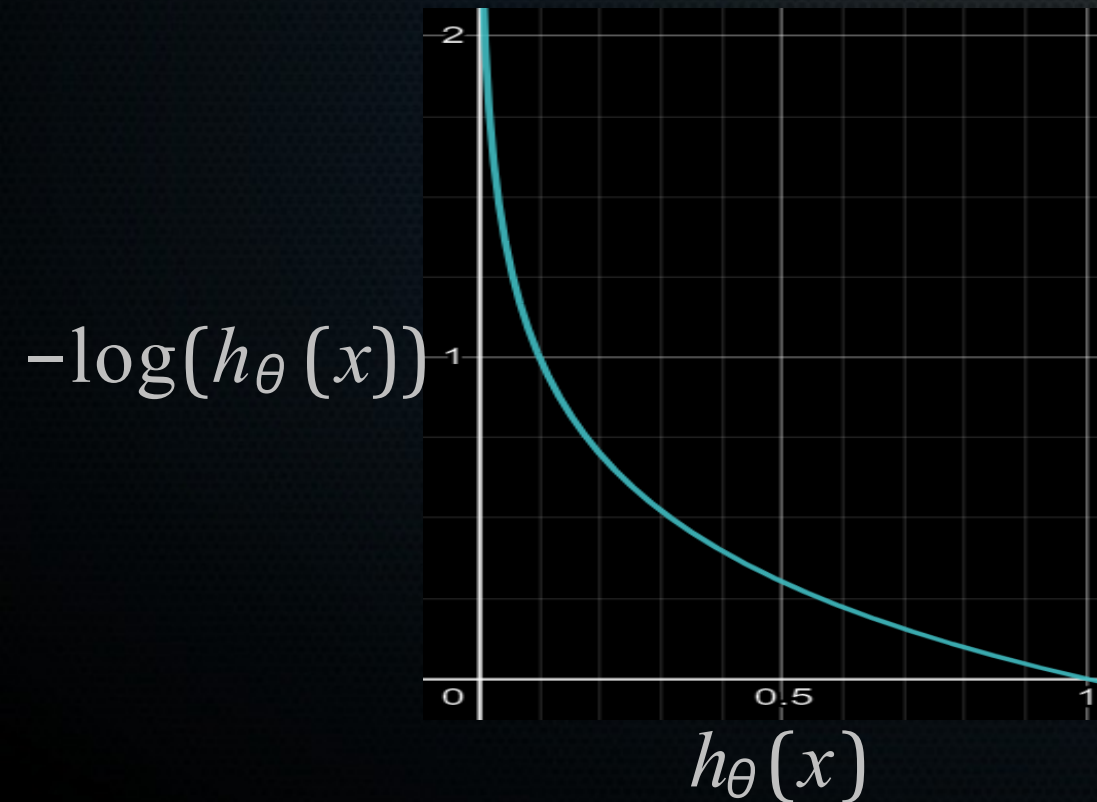
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$$



What does this function look like?

$$\text{Cost}(x) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

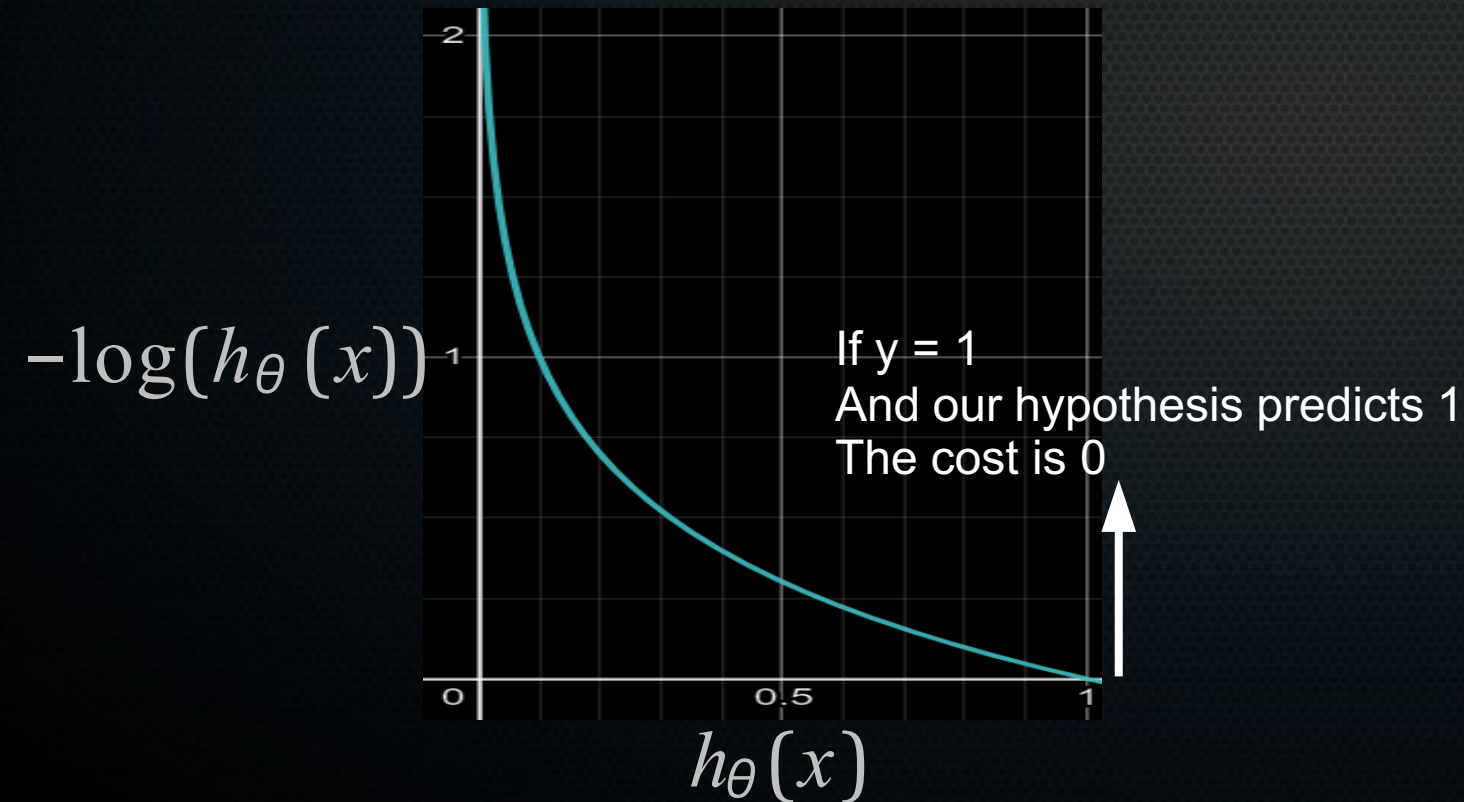
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$$



What does this function look like?

$$\text{Cost}(x) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

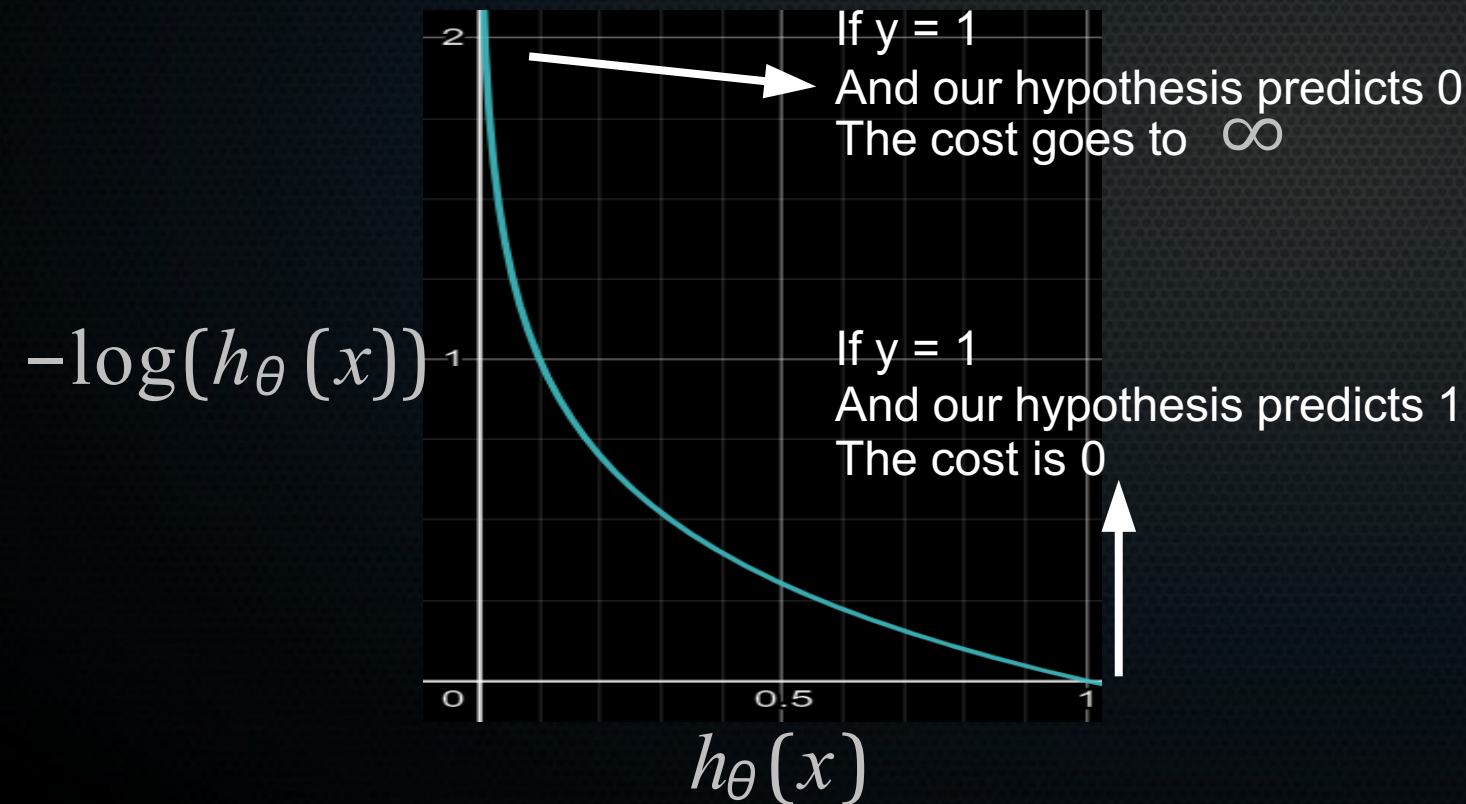
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$$



What does this function look like?

$$\text{Cost}(x) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

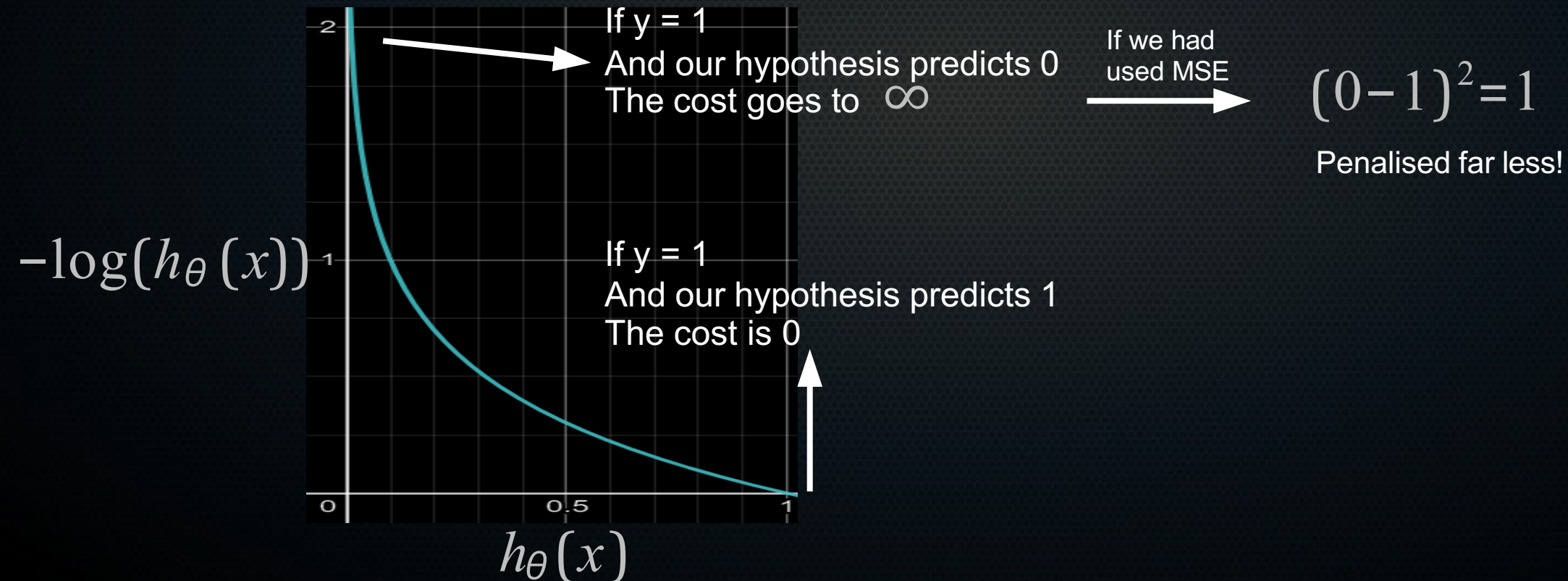
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$$



What does this function look like?

$$\text{Cost}(x) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

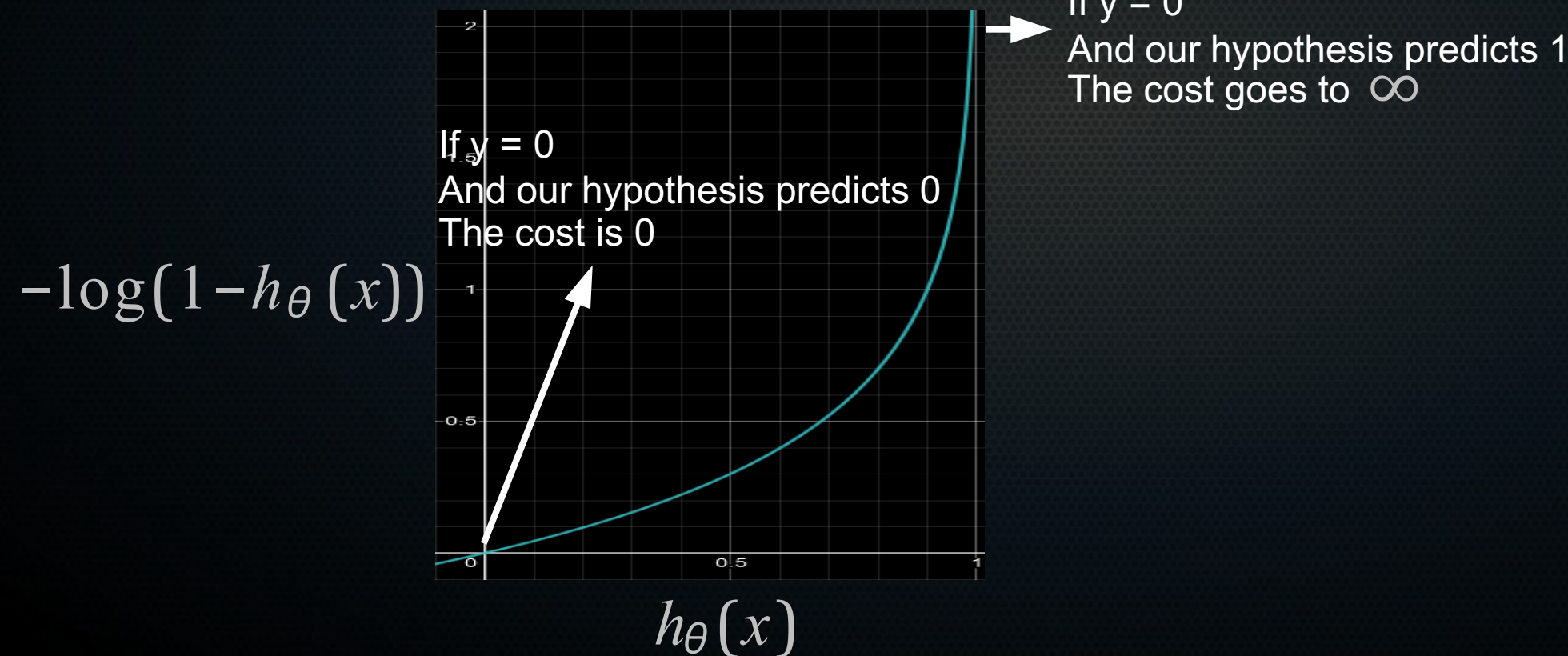
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$$



What does this function look like?

$$\text{Cost}(x) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$$



Simplified notation

$$\text{Cost}(x) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$



$$\text{Cost}(x) = -y \cdot \log(h_{\theta}(x)) - (1-y) \cdot \log(1-h_{\theta}(x))$$

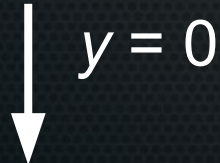
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$$

Simplified notation

$$\text{Cost}(x) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$



$$\text{Cost}(x) = -y \cdot \log(h_{\theta}(x)) - (1-y) \cdot \log(1-h_{\theta}(x))$$



$$-0 \cdot \log(h_{\theta}(x)) - (1-0) \cdot \log(1-h_{\theta}(x))$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$$

Simplified notation

$$\text{Cost}(x) = \begin{cases} -\log(h_\theta(x)) & \text{if } y=1 \\ -\log(1-h_\theta(x)) & \text{if } y=0 \end{cases}$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$$

$$\text{Cost}(x) = -y \cdot \log(h_\theta(x)) - (1-y) \cdot \log(1-h_\theta(x))$$

$$y=0$$

$$-0 \cdot \log(h_\theta(x)) - (1-0) \cdot \log(1-h_\theta(x))$$

$$-\log(1-h_\theta(x))$$

Simplified notation

$$\text{Cost}(x) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$$

$$\text{Cost}(x) = -y \cdot \log(h_{\theta}(x)) - (1-y) \cdot \log(1-h_{\theta}(x))$$

$$y = 1$$

$$-1 \cdot \log(h_{\theta}(x)) - (1-1) \cdot \log(1-h_{\theta}(x))$$

$$-\log(h_{\theta}(x))$$

Putting it all together

$$\text{Cost}(x) = -y \cdot \log(h_\theta(x)) - (1-y) \cdot \log(1-h_\theta(x)) \quad J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^{(i)})$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m -y^{(i)} \cdot \log(h_\theta(x^{(i)})) - (1-y^{(i)}) \cdot \log(1-h_\theta(x^{(i)}))$$

Putting it all together

$$\text{Cost}(x) = -y \cdot \log(h_\theta(x)) - (1-y) \cdot \log(1-h_\theta(x)) \quad J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(x^i)$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[-y^{(i)} \cdot \log(h_\theta(x^{(i)})) - (1-y^{(i)}) \cdot \log(1-h_\theta(x^{(i)})) \right]$$



$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \cdot \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \cdot \log(1-h_\theta(x^{(i)})) \right]$$

Optimising the cost function

- Same form as for linear regression (only hypothesis function differs!)

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m ((h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)})$$

$$\theta_j := \theta_j - \frac{\alpha}{m} \sum_{i=1}^m ((h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)})$$

$$h_{\theta}(x^{(i)}) = \frac{1}{1 + e^{-\theta^T * x}}$$

[Derivation
link](#)

Summary

- By using the sigmoid function as a transformation of normal regression and interpreting the output as a chance of being 0 or 1 we can do classification.
- Only the form of our hypothesis function is different
- Need a different cost function: should be smooth, and give logical values for large errors.

Break for practical

