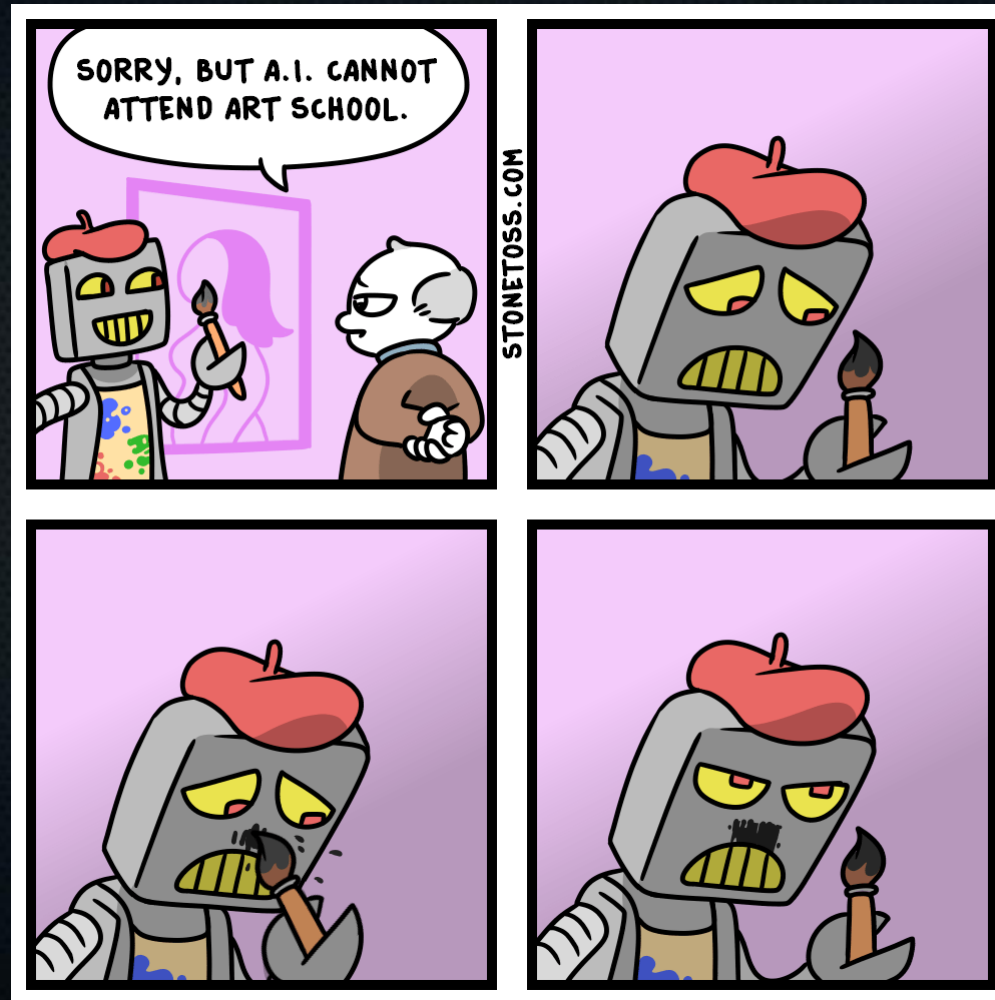# Daily Inspiration

# Daily Inspiration

# Today

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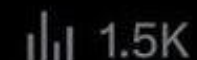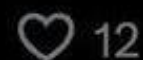- Neural network basics

# 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 wel 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

# 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
  Class >= 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 Class >= 0.5, predict class 1. Otherwise, class 0.

- Problems:
  -You can predict class > 1 and < 0, while that is not possible in reality.

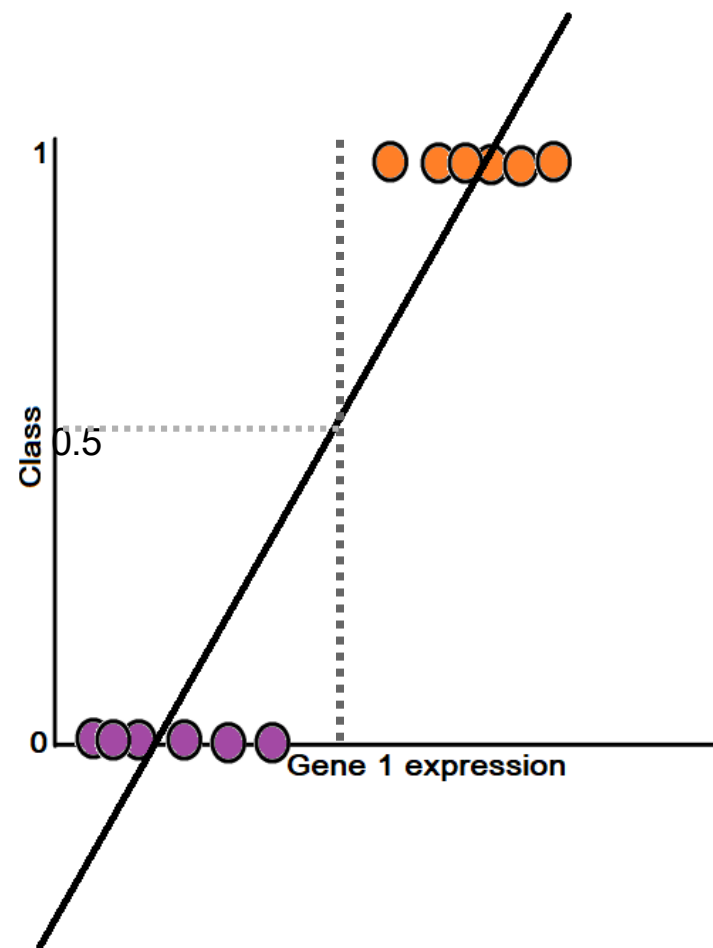# Logistic regression

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

- Problems:
  -You can predict class > 1 and < 0,
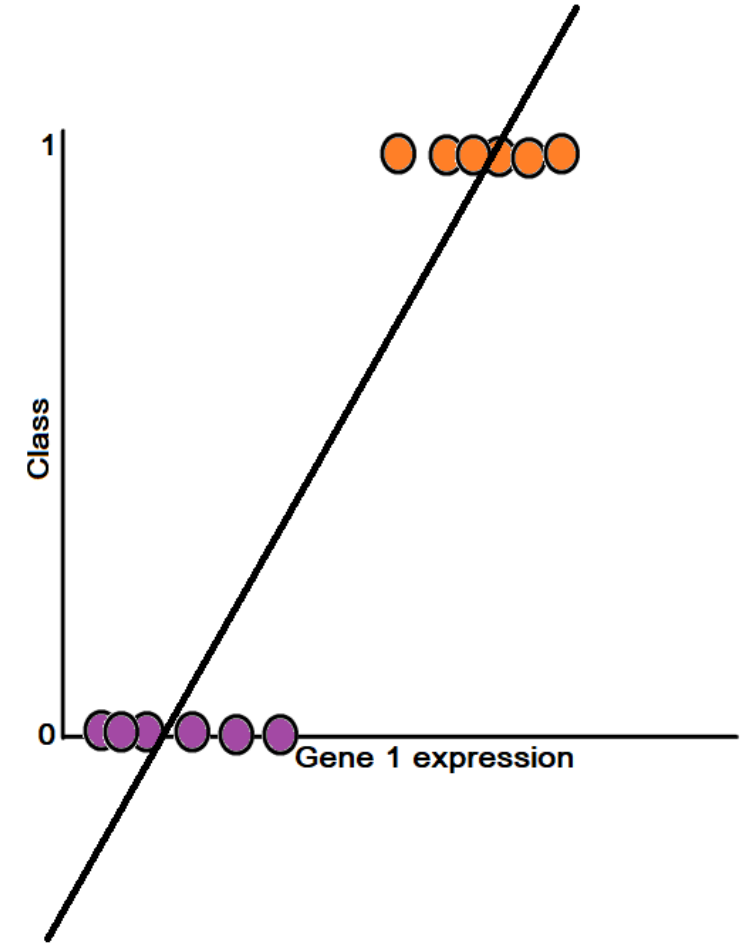  while that is not possible in reality.
  -This example seemed to work,
  but quickly breaks down →

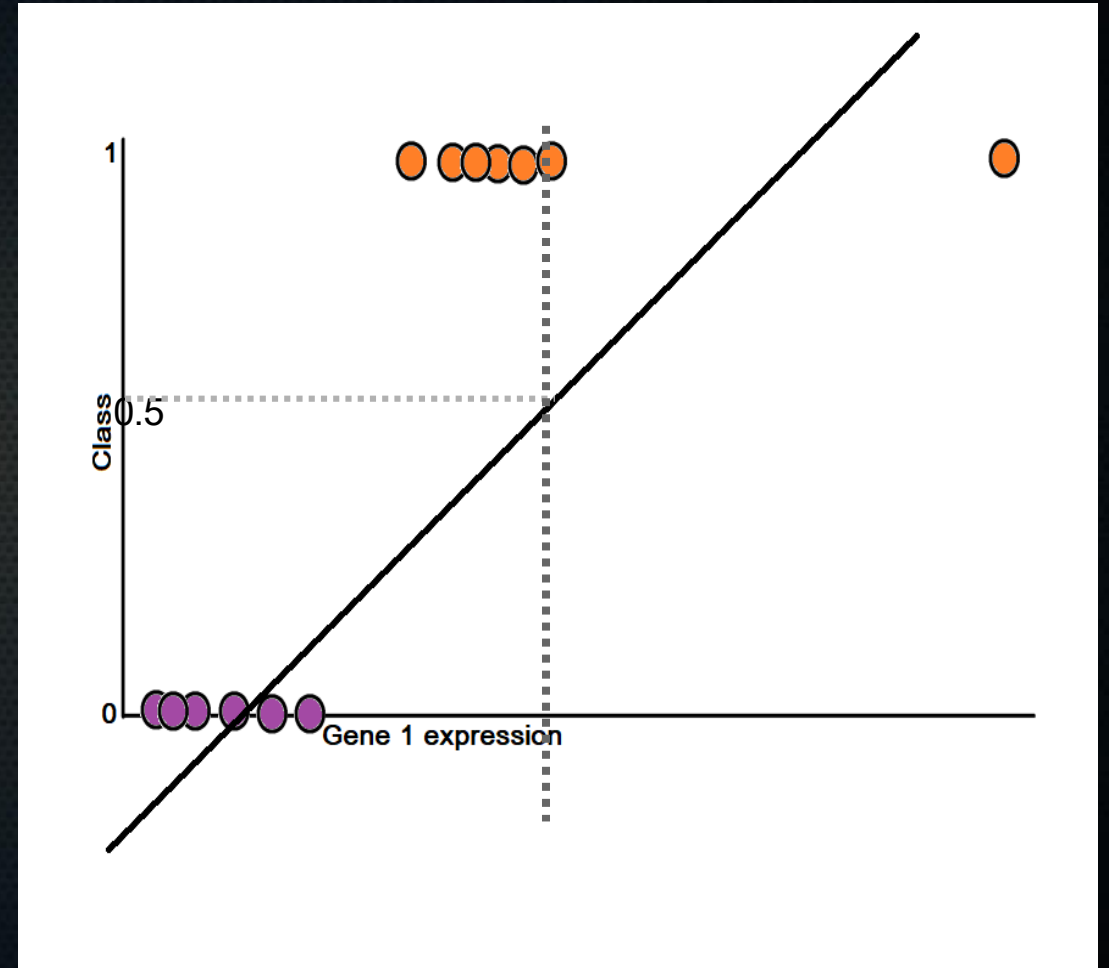# Logistic regression

- Naïve idea:
  Train a linear regression. If Class >= 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 →
  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:

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

→

$$[0.5 \quad 3 \quad -1.5] \cdot \begin{bmatrix} 1 \\ 3 \\ 8 \end{bmatrix}$$

Features for one sample (x0 = 1, intercept term, 2 data-derived features x1 and x2)

Learned parameters (theta 0 – theta 2)

# Sigmoid or logistic function

- Before, our hypothesis function was of the form:

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

# 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:
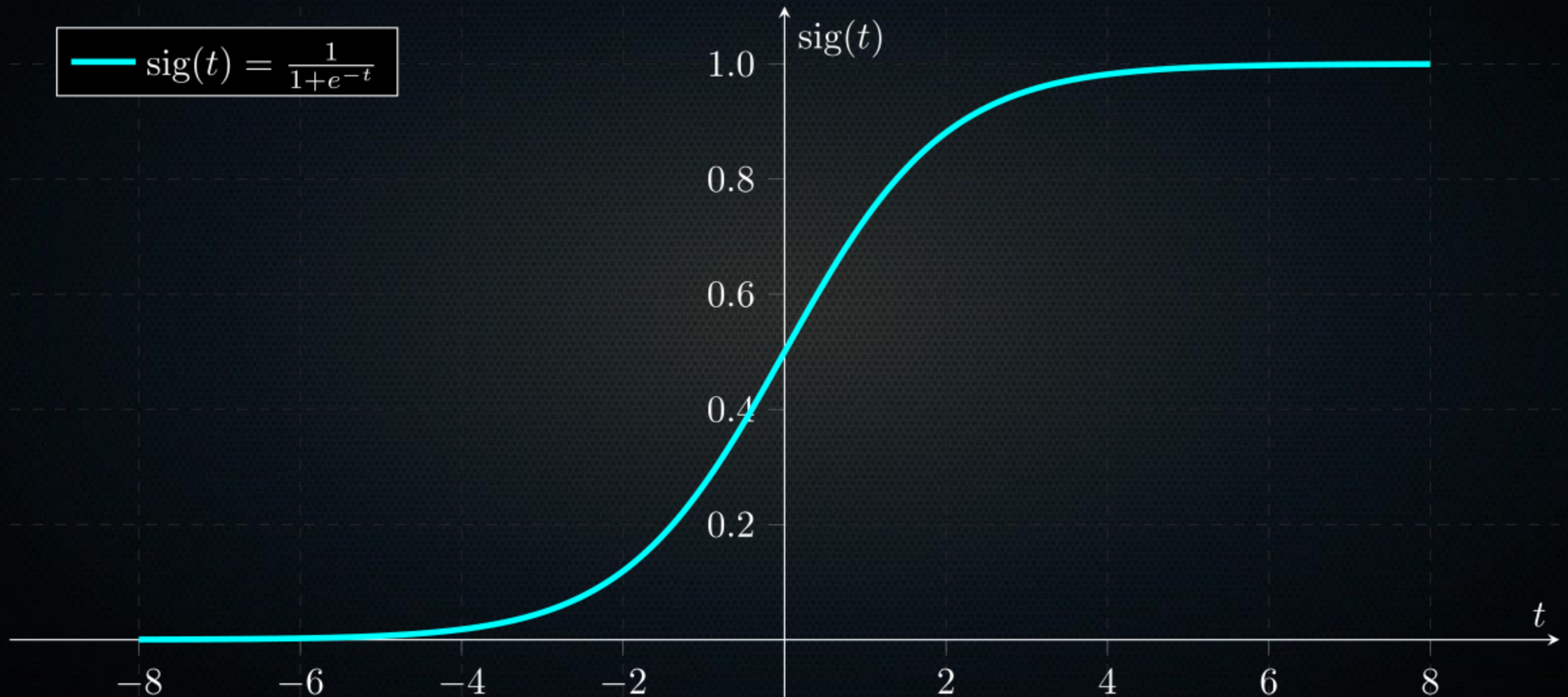
$$h_\theta(x) = \theta^T \cdot x$$

- Change that to the following:

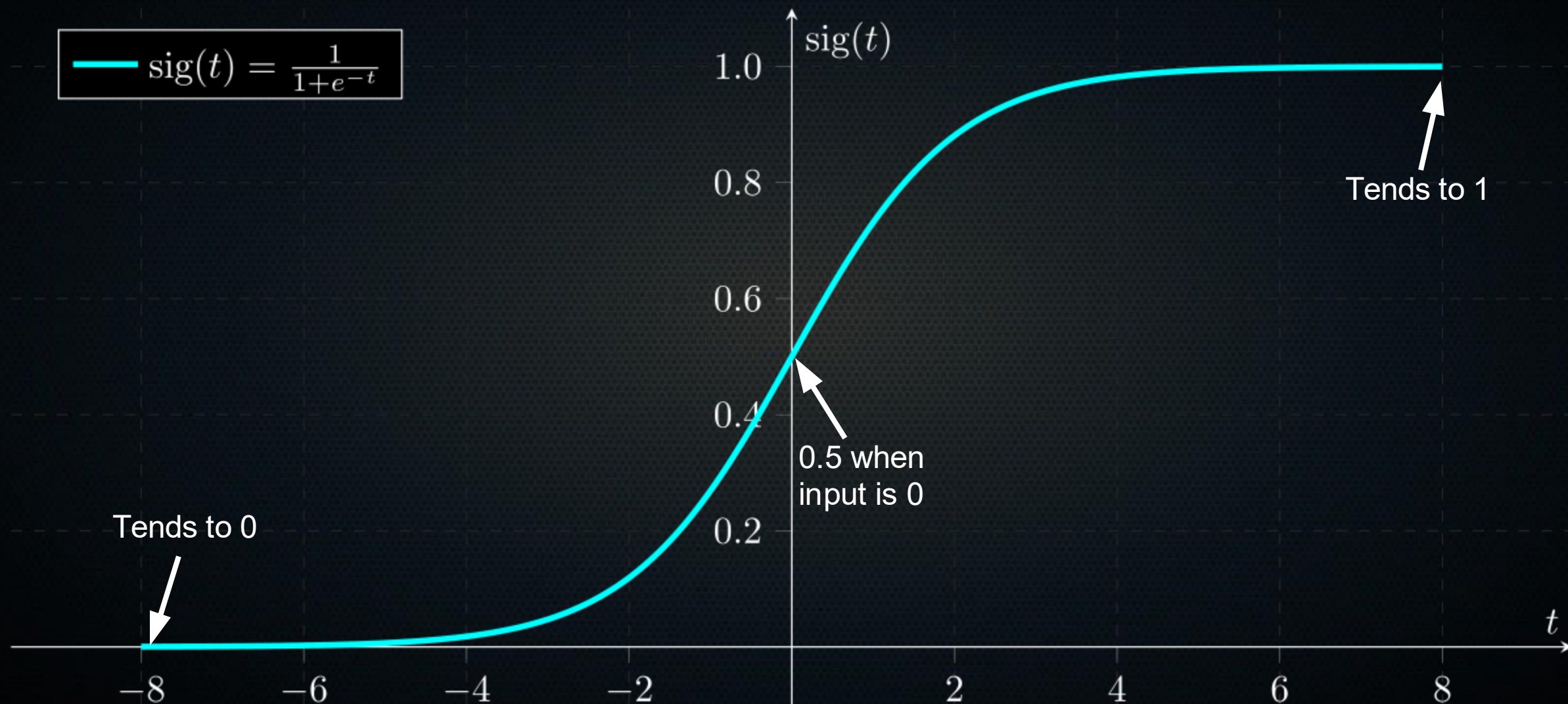$$h_\theta(x) = g(\theta^T \cdot x) \longrightarrow g(z) = \frac{1}{1+e^{-z}}$$

- What does that look like?

$$z \to \infty, e^{-z} \to 0$$

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

# What does the sigmoid function look like?



$$\text{sig}(t) = \frac{1}{1+e^{-t}}$$

# What does the sigmoid function look like?



$$\text{sig}(t) = \frac{1}{1+e^{-t}}$$

Tends to 1

0.5 when input is 0

Tends to 0

# Sigmoid or logistic function

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

# Sigmoid or logistic function

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

  - Interpret outcome of $h_\theta(x)$ as probability that class = 1 given the features.
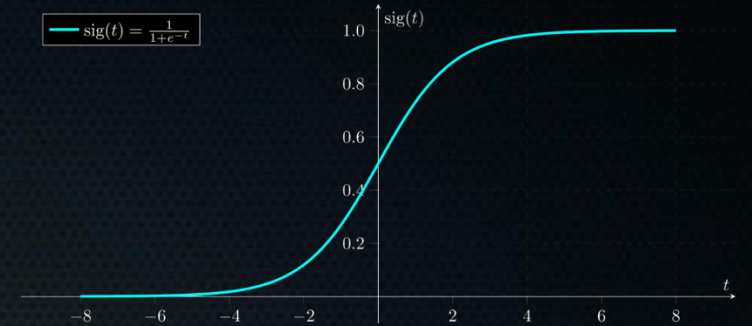
# Sigmoid or logistic function

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

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

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

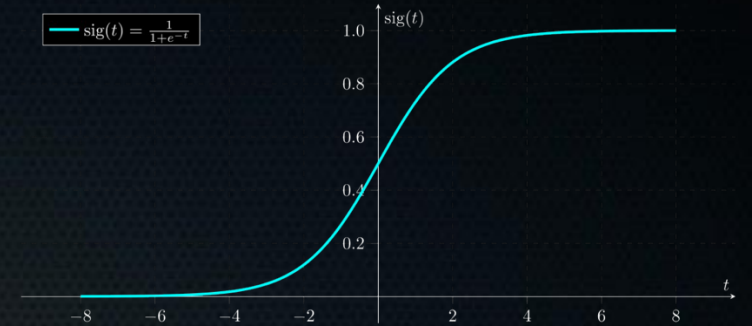$h_\theta(x) = 0.8$ $\longrightarrow$ 80% chance of tumor being malignant

# Sigmoid or logistic function

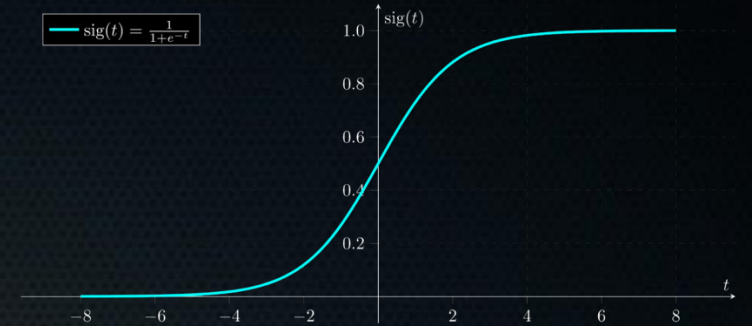- How do we work with this?  $h_\theta(x) = \dfrac{1}{1 + e^{-(\theta^T \cdot x)}}$

  - 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% → 20 % chance of being benign (class 0)

# Sigmoid or logistic function

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

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

  - Formally:

$$h_\theta(x) = \frac{1}{1 + e^{-(\theta^T \cdot x)}} = p(y=1 | x ; \theta)$$
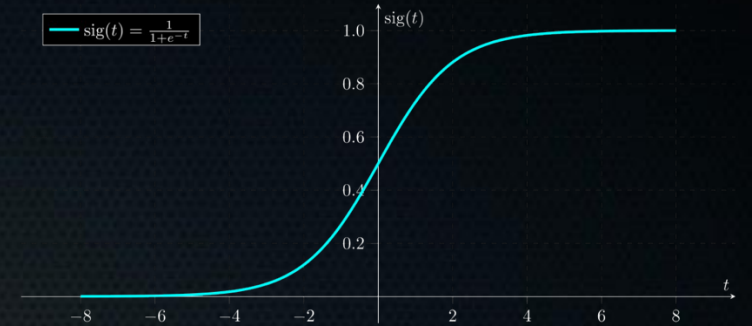
$$p(y=0 | x ; \theta) = 1 - h_\theta(x)$$

# Sigmoid or logistic function

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

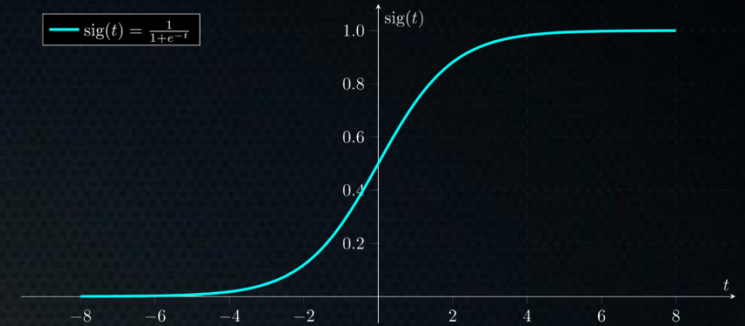  - Interpret outcome of $h_\theta(x)$ as probability that class = 1 given the features.

  - Formally:

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

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

$$\dfrac{p(y=1)}{1 - p(y=1)}$$

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 | . |

logOdds(diabetes) = -8.9 + (0.106*pregnant) + (0.037*glucose).... + (0.02*age)

# Decision boundary

- Threshold:



$$\mathrm{sig}(t) = \frac{1}{1+e^{-t}}$$

$$h_\theta(x) < 0.5$$

$$\theta^T \cdot x < 0$$

# Decision boundary

- Threshold:



$$h_\theta(x) \geq 0.5$$

$$h_\theta(x) < 0.5$$

$$\text{sig}(t) = \frac{1}{1+e^{-t}}$$

$$\theta^T \cdot x < 0 \qquad\qquad \theta^T \cdot x \geq 0$$
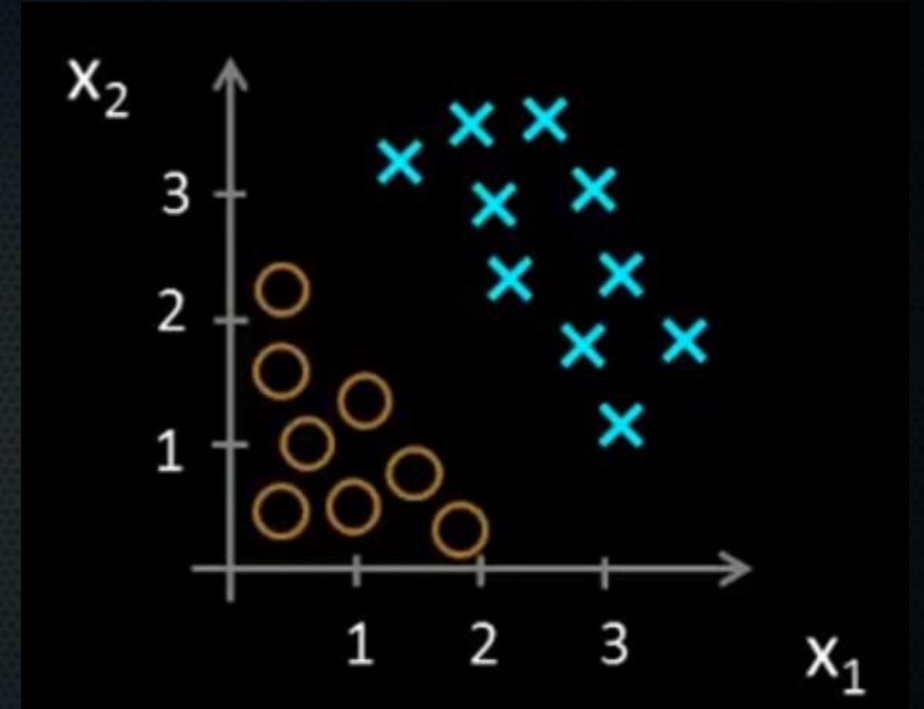
# Decision boundary

- How does it look?

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

$$h_\theta(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$$

# Decision boundary

- How does it look?

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

$$h_\theta(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$$

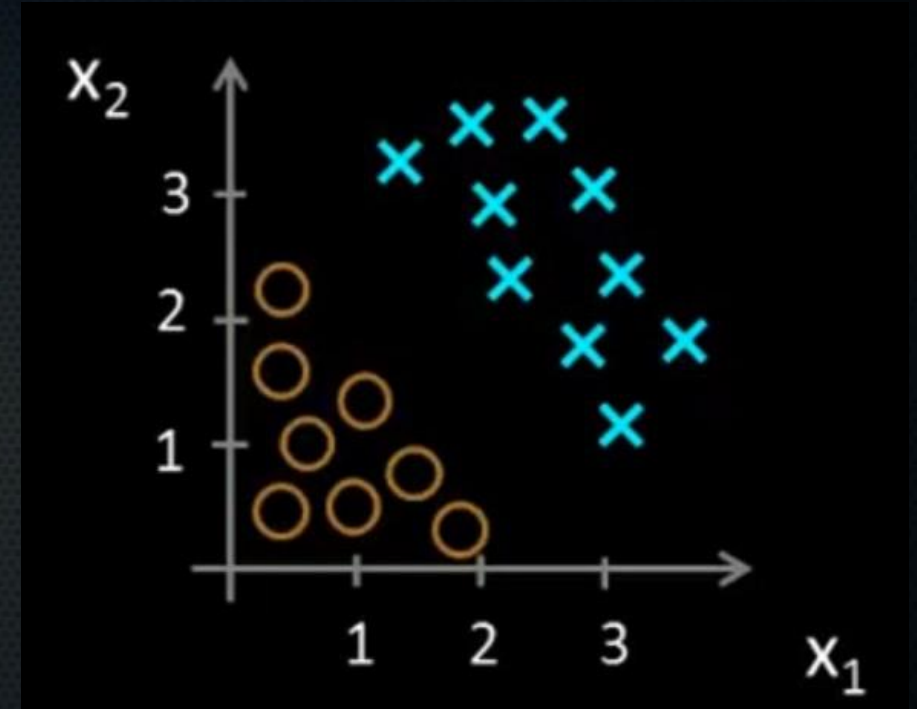$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

# Decision boundary

- How does it look?

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

$$h_\theta(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$$

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

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

# Decision boundary

- How does it look?

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



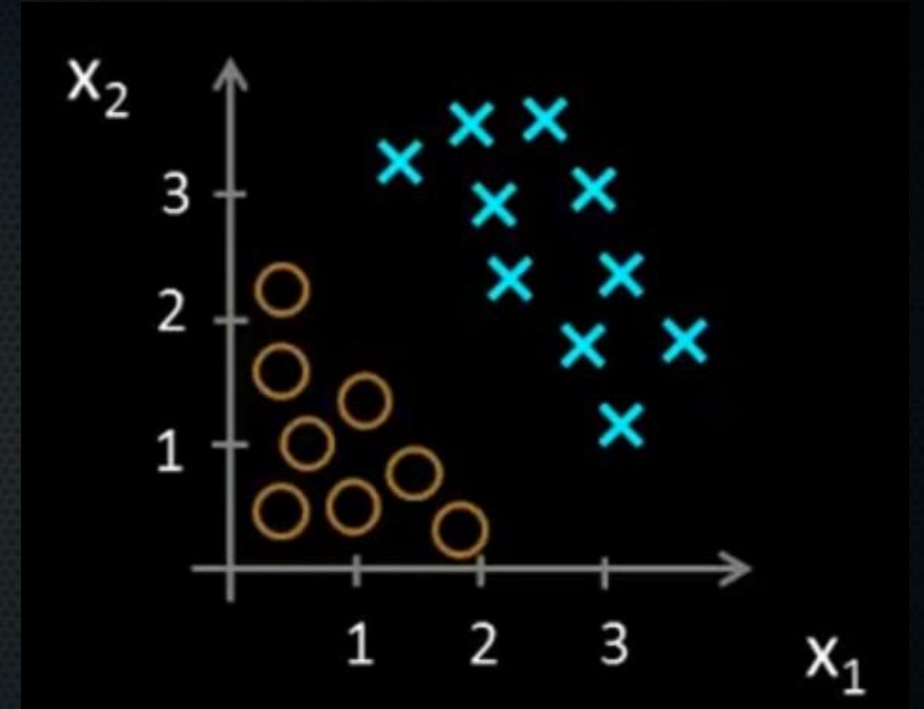$$h_\theta(x) \geq 0.5$$

$$\theta^T \cdot x \geq 0$$

# Decision boundary

- How does it look?

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

$$h_\theta(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$$

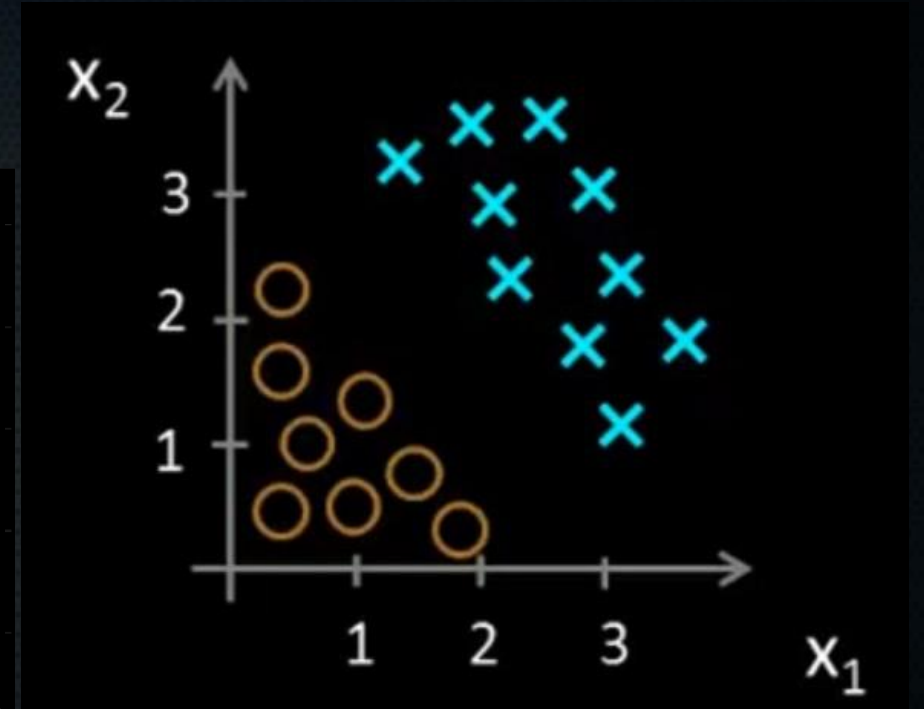$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

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

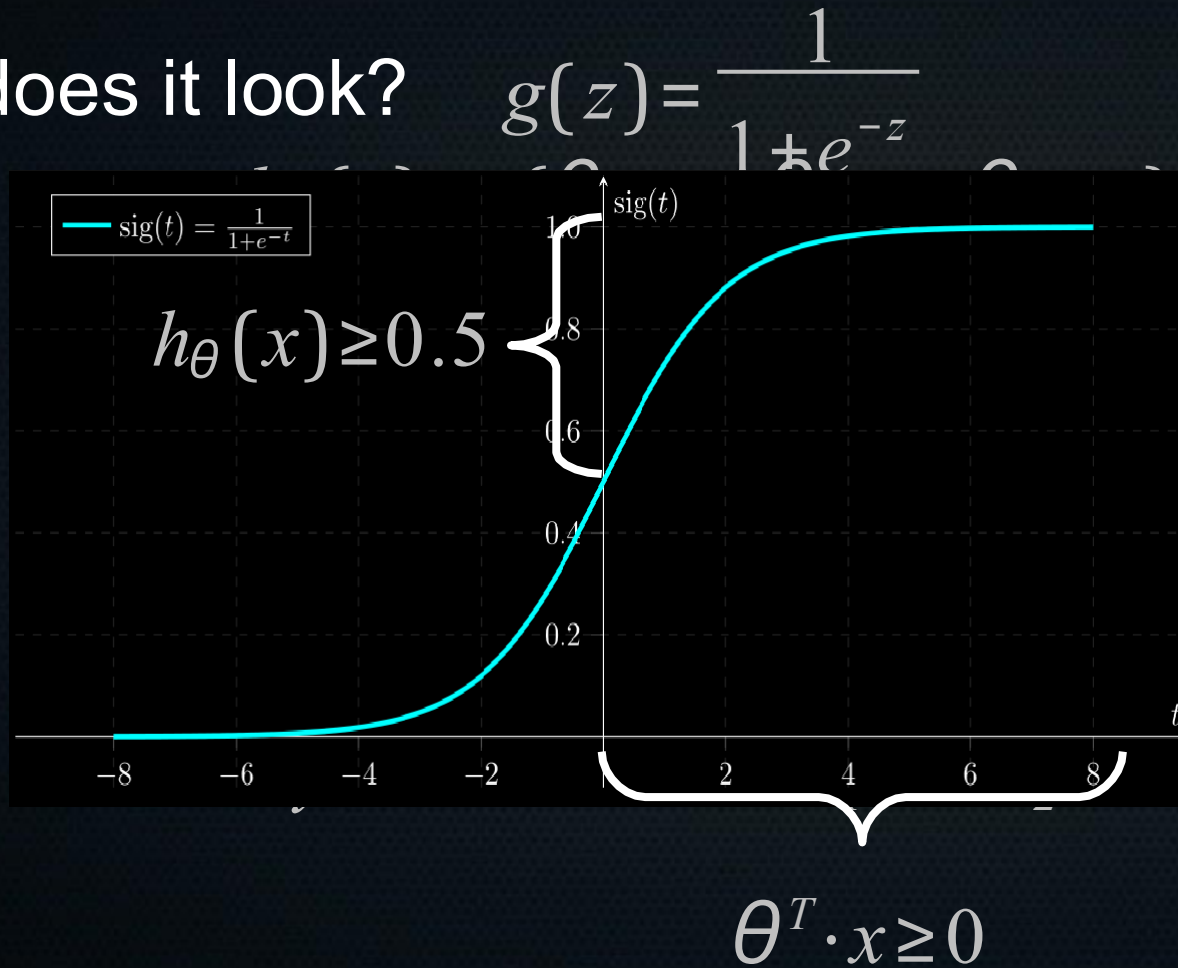$$-3 \cdot 1 + 1 \cdot x_1 + 1 \cdot x_2 \geq 0$$

$$x_1 + x_2 = 3$$
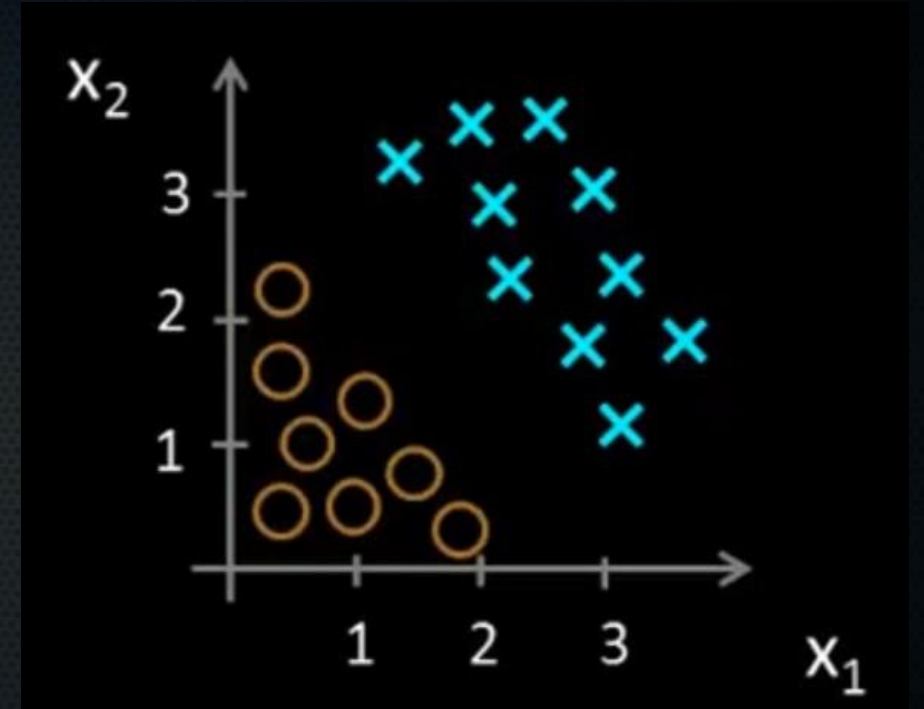
# Decision boundary

- How does it look?

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

$$h_\theta(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$$

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

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

$$-3 \cdot 1 + 1 \cdot x_1 + 1 \cdot x_2 \geq 0$$

$$x_1 + x_2 = 3$$

# Decision boundary

- How does it look?

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

$$h_\theta(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$$

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

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}}$$

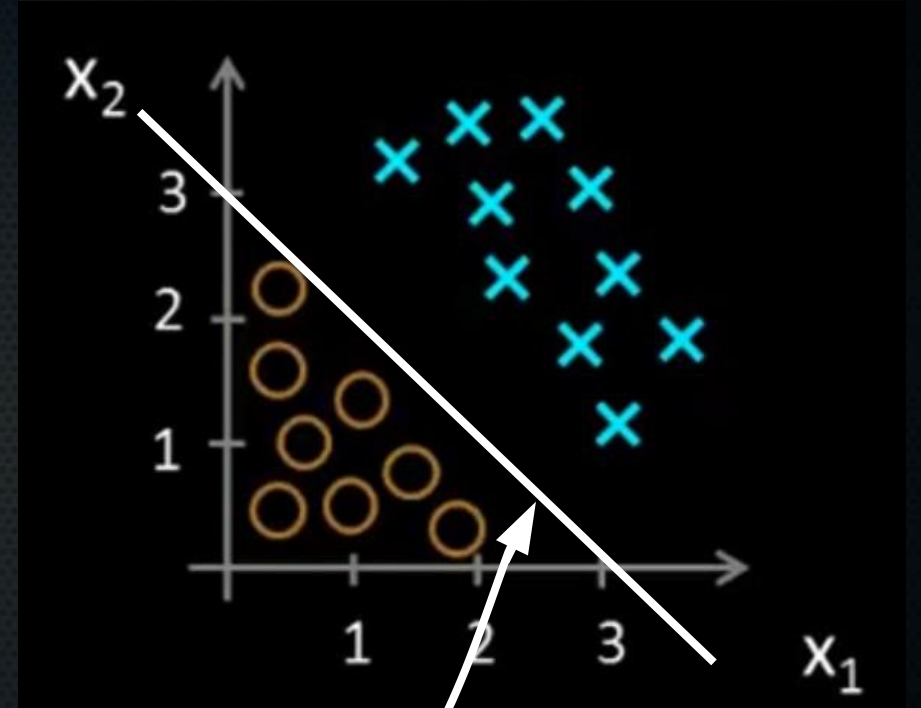$$y = 1 \text{ if } -3 \cdot 1 + 1 \cdot x_1 + 1 \cdot x_2 \geq 0$$

# Decision boundary

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

$$h_\theta(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$$

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

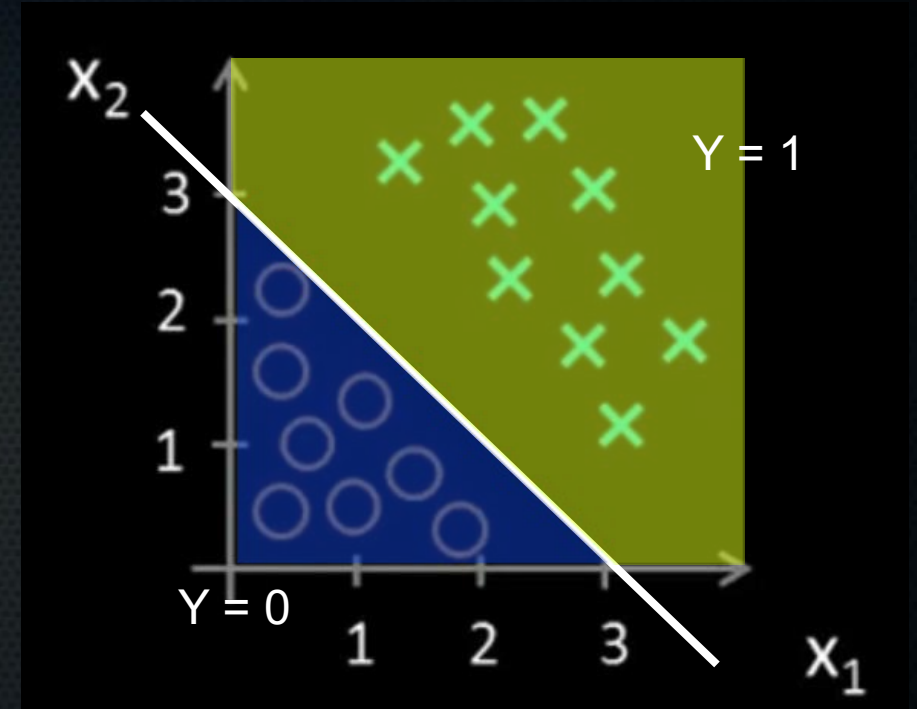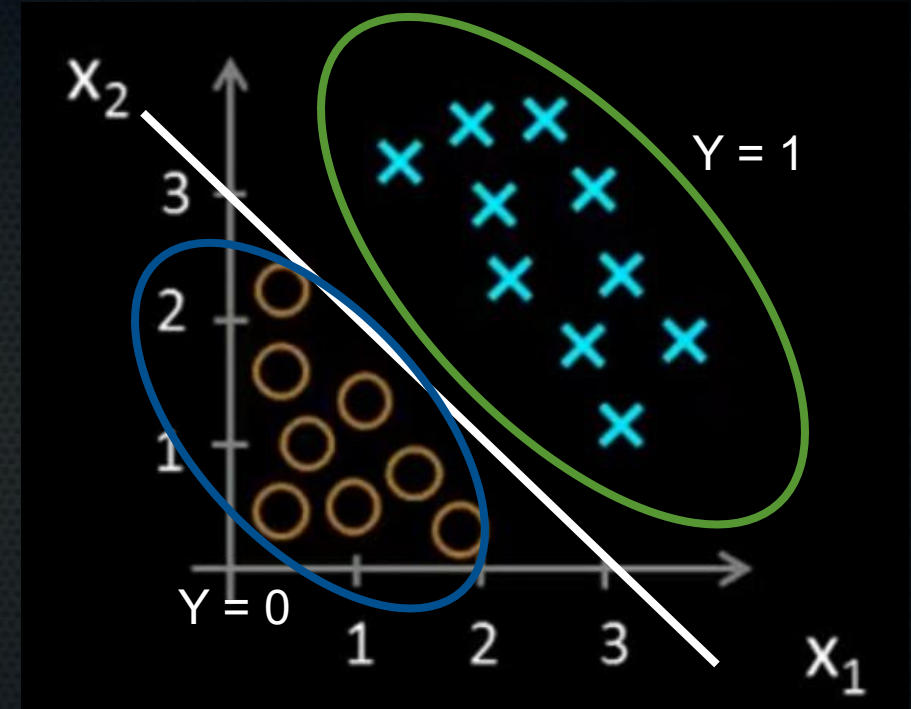$$y = 1 \ \text{if} -3 \cdot 1 + 1 \cdot x_1 + 1 \cdot x_2 \geq 0$$
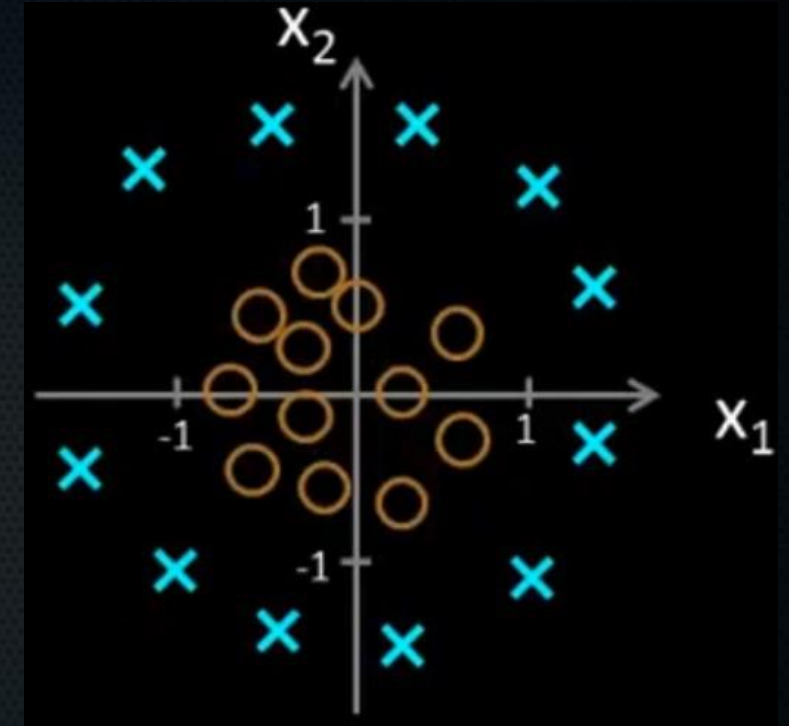
# Decision boundary

- How does it look?

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

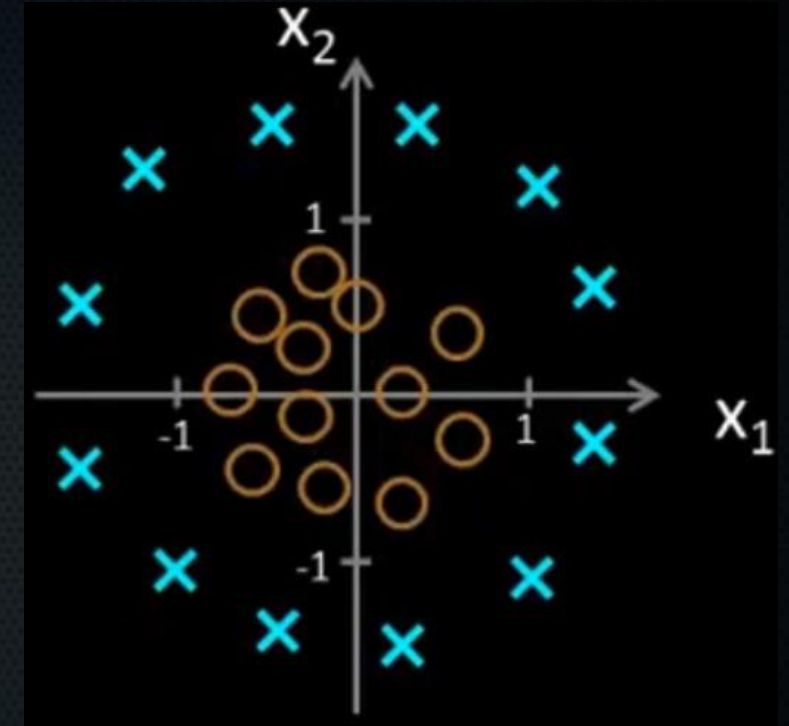$$h_\theta(x) = g(\theta_0 \cdot x_0, \theta_1 \cdot x_1, \theta_2 \cdot x_2)$$

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



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

# Non-linear decision boundary

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

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

- Add two polynomial features

# Non-linear decision boundary

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

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

- Add two polynomial features

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

# Non-linear decision boundary

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

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

- Add two polynomial features

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

Can you work out what the decision boundary will be?
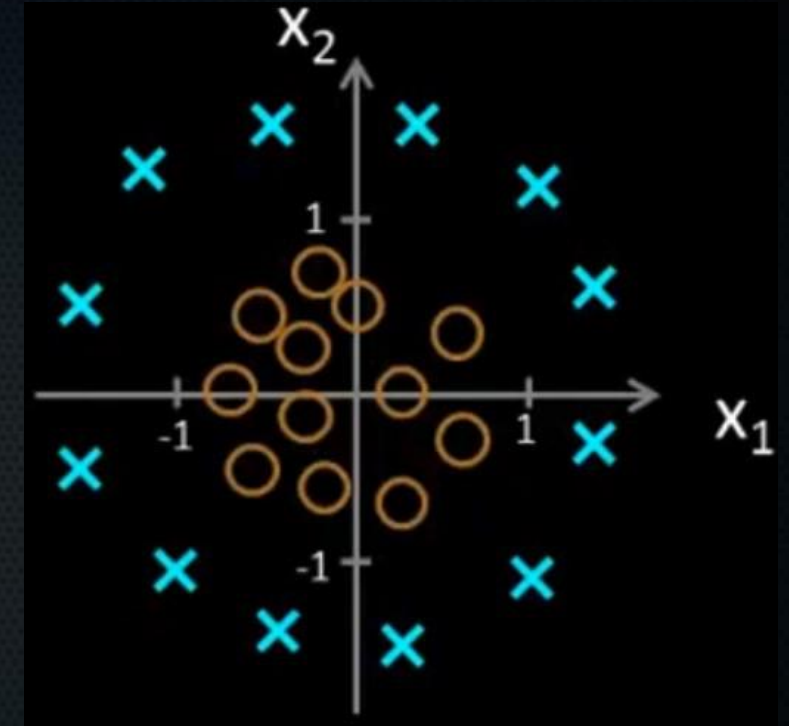
# Non-linear decision boundary

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

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

- Add two polynomial features

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \longrightarrow \begin{array}{l} y = 1 \ \text{if} \\[4pt] -1 + x_1^2 + x_2^2 \geq 0 \end{array}$$
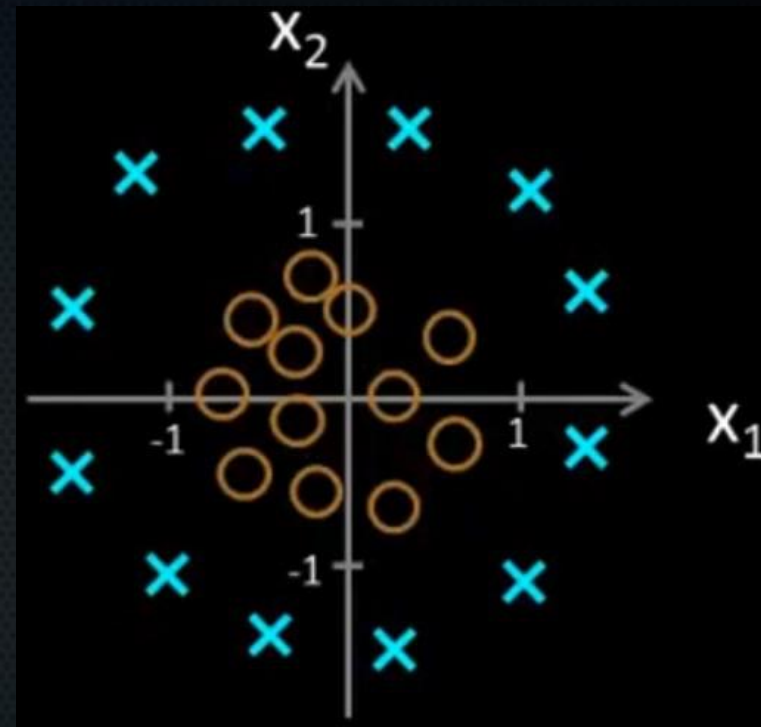
# Non-linear decision boundary



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

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

- Add two polynomial features

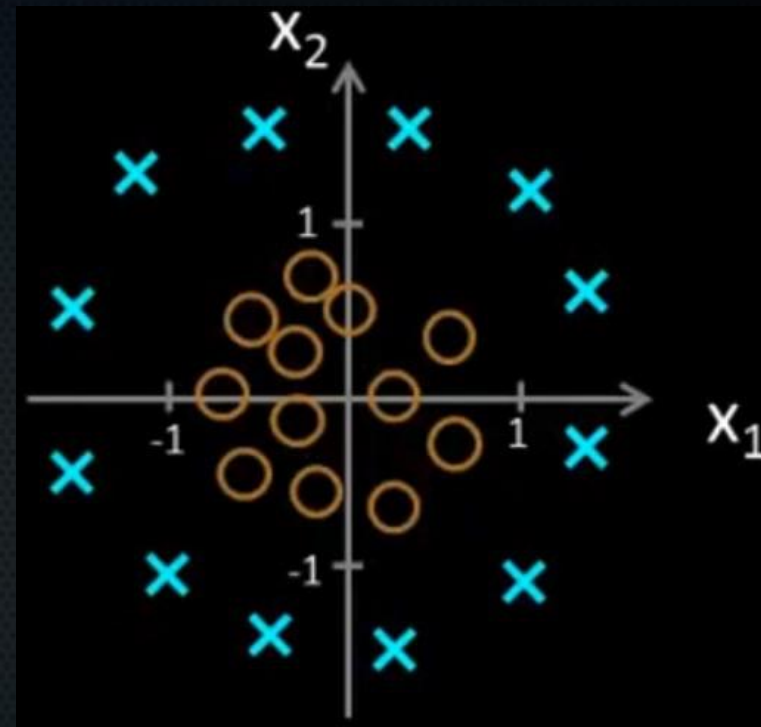$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \longrightarrow \quad -1 + x_1^2 + x_2^2 \geq 0$$

$$\downarrow$$

$$x_1^2 + x_2^2 \geq 1$$

# Non-linear decision boundary

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

$$h_\theta(x) = g(\theta_0 x_0, \theta_1 x_1, \theta_2 x_2, \boxed{\theta_3 x_1^2, \theta_4 x_2^2})$$
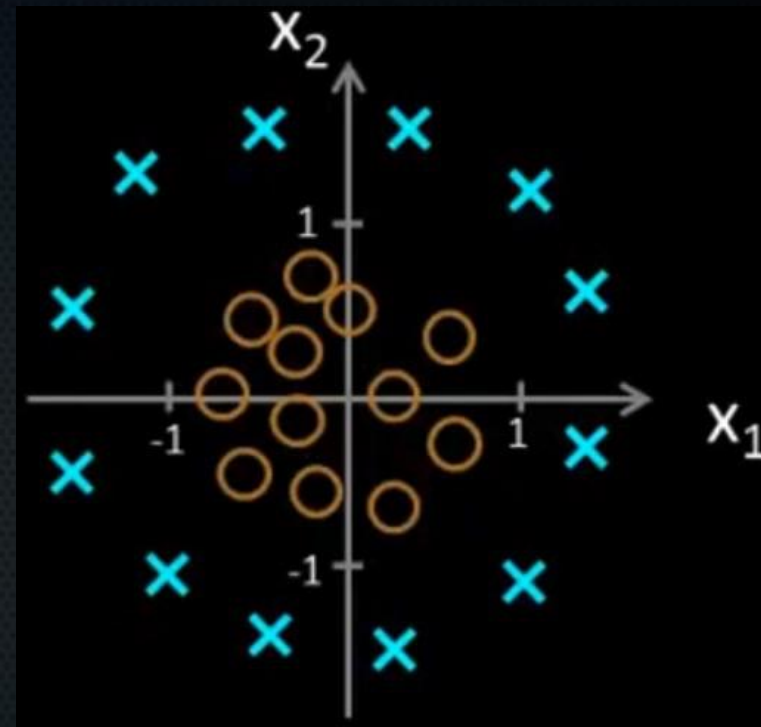
- Add two polynomial features

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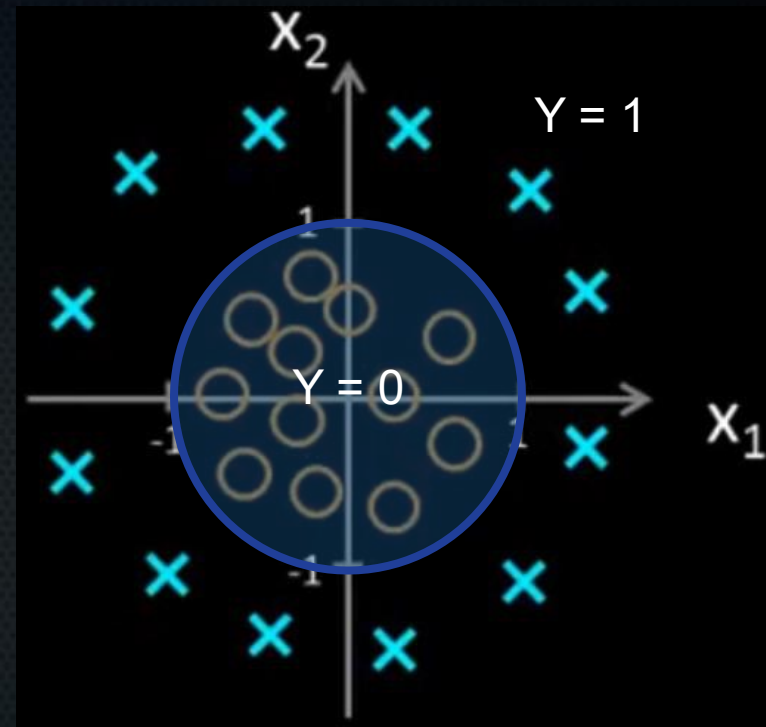
$$x_1^2 + x_2^2 \geq 1$$

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

$$\updownarrow$$

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

# So how do we get theta's?

- Need a cost function

- Before:

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

$$\updownarrow$$

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

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

# So how do we get theta's?

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

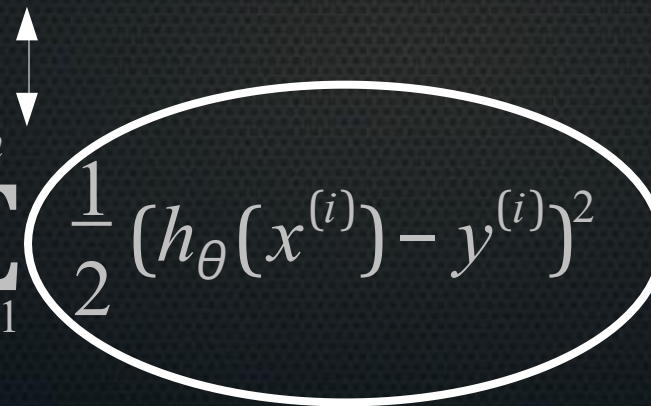$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \boxed{\frac{1}{2}(h_\theta(x^{(i)}) - y^{(i)})^2}$$

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

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

# So how do we get theta's?

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

# So how do we get theta's?

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

non-convex

$J(\theta)$

$\theta$

convex

$J(\theta)$

$\theta$

# So how do we get theta's?

- Need a cost function $J(\theta) = \dfrac{1}{m}\sum\limits_{i=1}^{m} \mathrm{Cost}(x^i)$ $\mathrm{Cost}(x) = \dfrac{1}{2}\left(h_\theta(x) - y\right)^2$
- Why not MSE? → not convex

non-convex

convex

$J(\theta)$

$J(\theta)$

$\theta$

$\theta$

# 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

$J(\theta)$

$\theta$

convex

$J(\theta)$

$\theta$

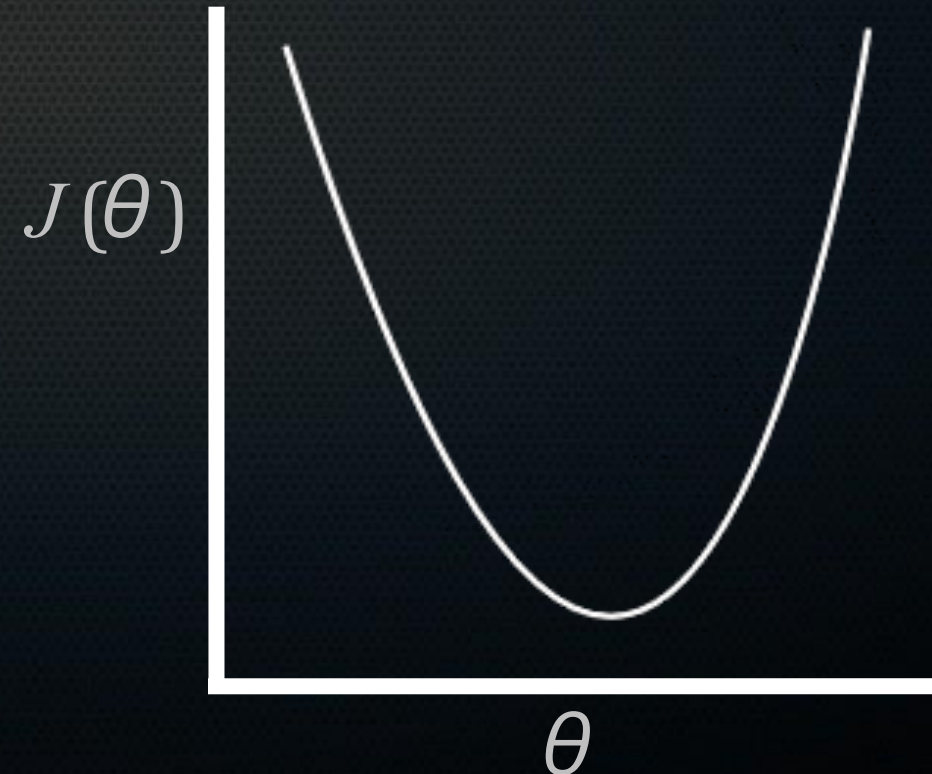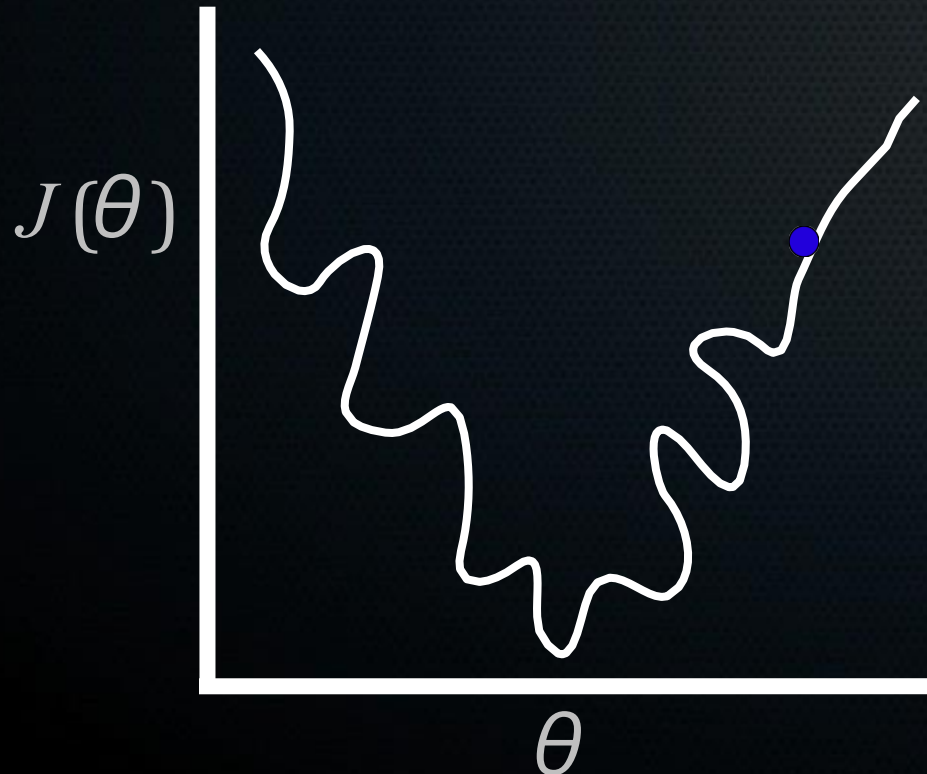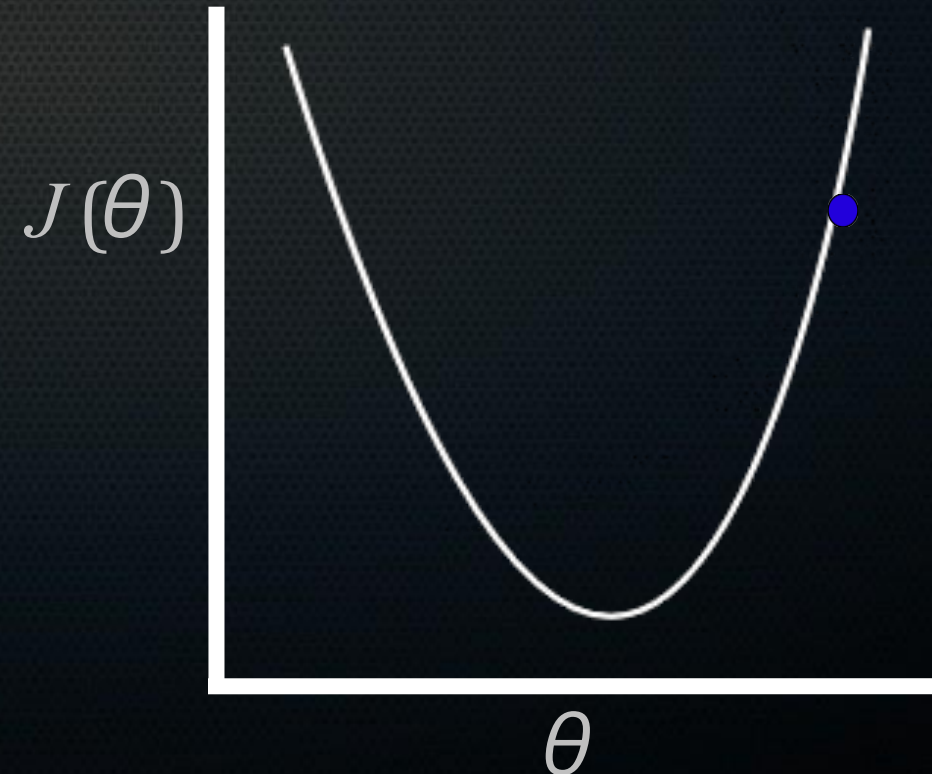# So how do we get theta's?

- 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

$J(\theta)$

$J(\theta)$

$\theta$

$\theta$

# So how do we get theta's?

- ▪ Need a cost function $J(\theta) = \dfrac{1}{m} \sum\limits_{i=1}^{m} \mathrm{Cost}(x^i)$ $\qquad \mathrm{Cost}(x) = \dfrac{1}{2}(h_\theta(x) - y)^2$
- ▪ Why not MSE? → not convex

non-convex

$J(\theta)$

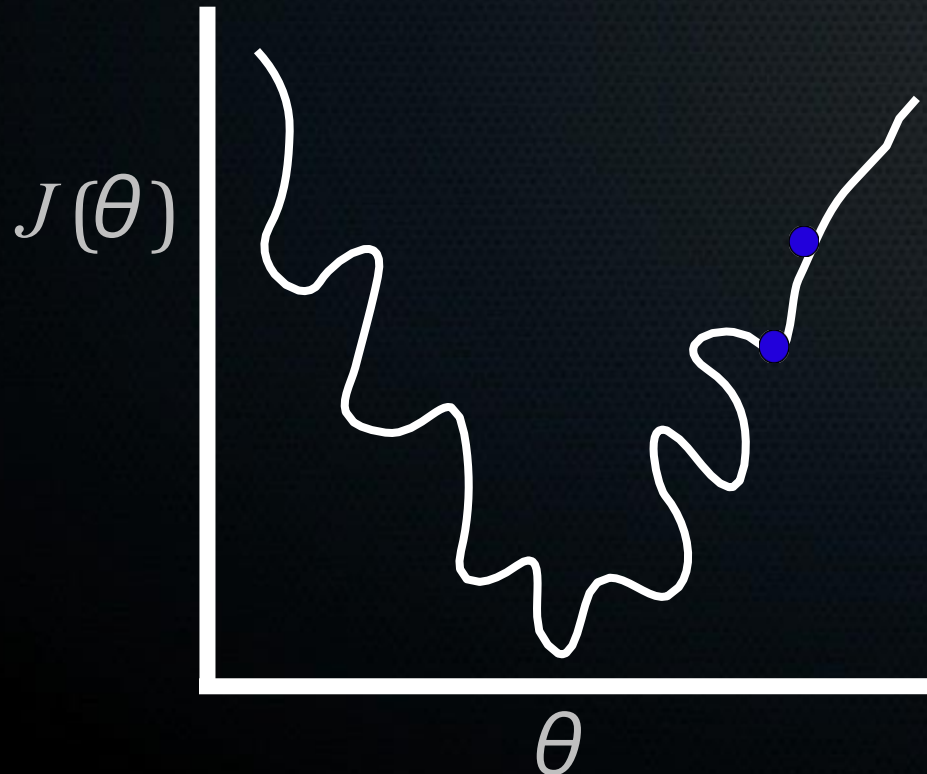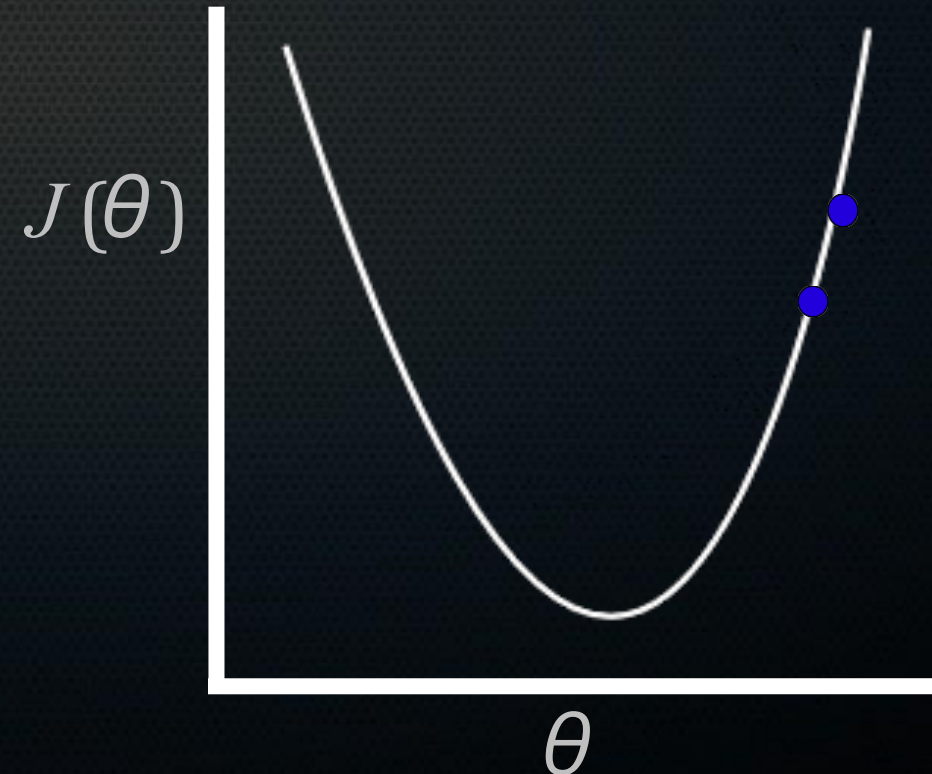convex

$J(\theta)$
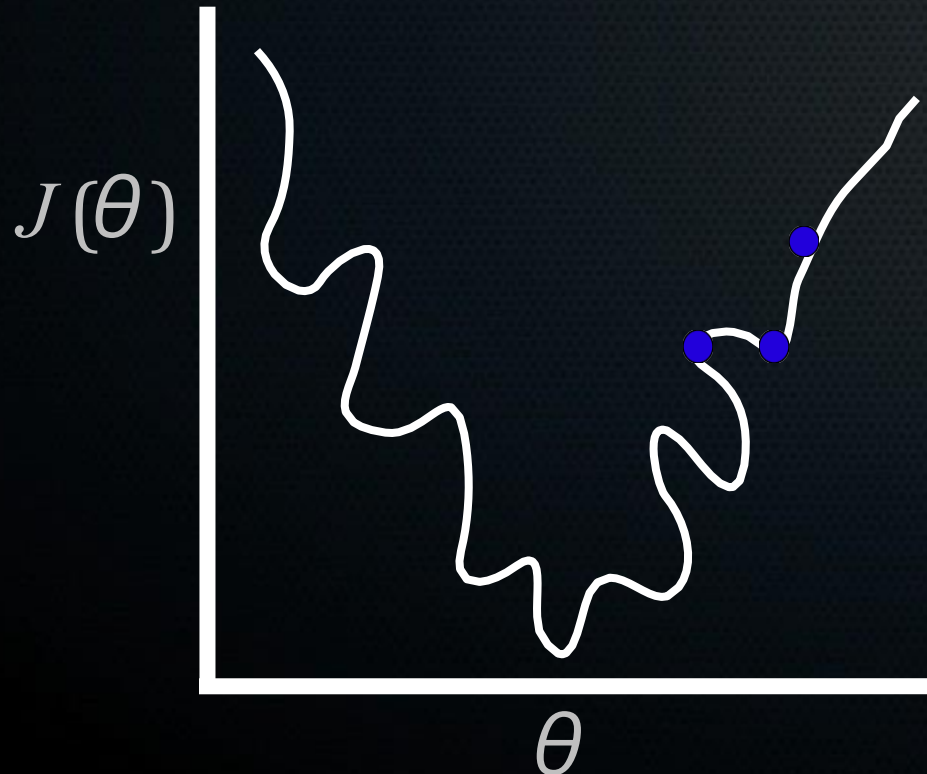
$\theta$

$\theta$

# So how do we get theta's?

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



non-convex

$J(\theta)$

$\theta$

convex

$J(\theta)$

$\theta$

# So how do we get theta's?

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

# So how do we get theta's?

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

$$\mathrm{Cost}(x) = \begin{cases} -\log(h_\theta(x)) \text{ if } y = 1 \\[2ex] -\log(1 - h_\theta(x)) \text{ if } y = 0 \end{cases}$$

# What does this function look like?

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

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

$\log(x)$

Tends to infinity

Tends to -infinity

# What does this function look like?

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

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



$-\log(x)$

# What does this function look like?

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

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



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

$h_\theta(x)$

# What does this function look like?

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

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

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

If y = 1
And our hypothesis predicts 1
The cost is 0

$h_\theta(x)$

# What does this function look like?

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

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

If y = 1
And our hypothesis predicts 0
The cost goes to  $\infty$

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

If y = 1
And our hypothesis predicts 1
The cost is 0

$h_\theta(x)$

# What does this function look like?

$$\text{Cost}(x) = \begin{cases} \boxed{-\log(h_\theta(x)) \text{ if } y=1} \\ -\log(1-h_\theta(x)) \text{ if } y=0 \end{cases}$$
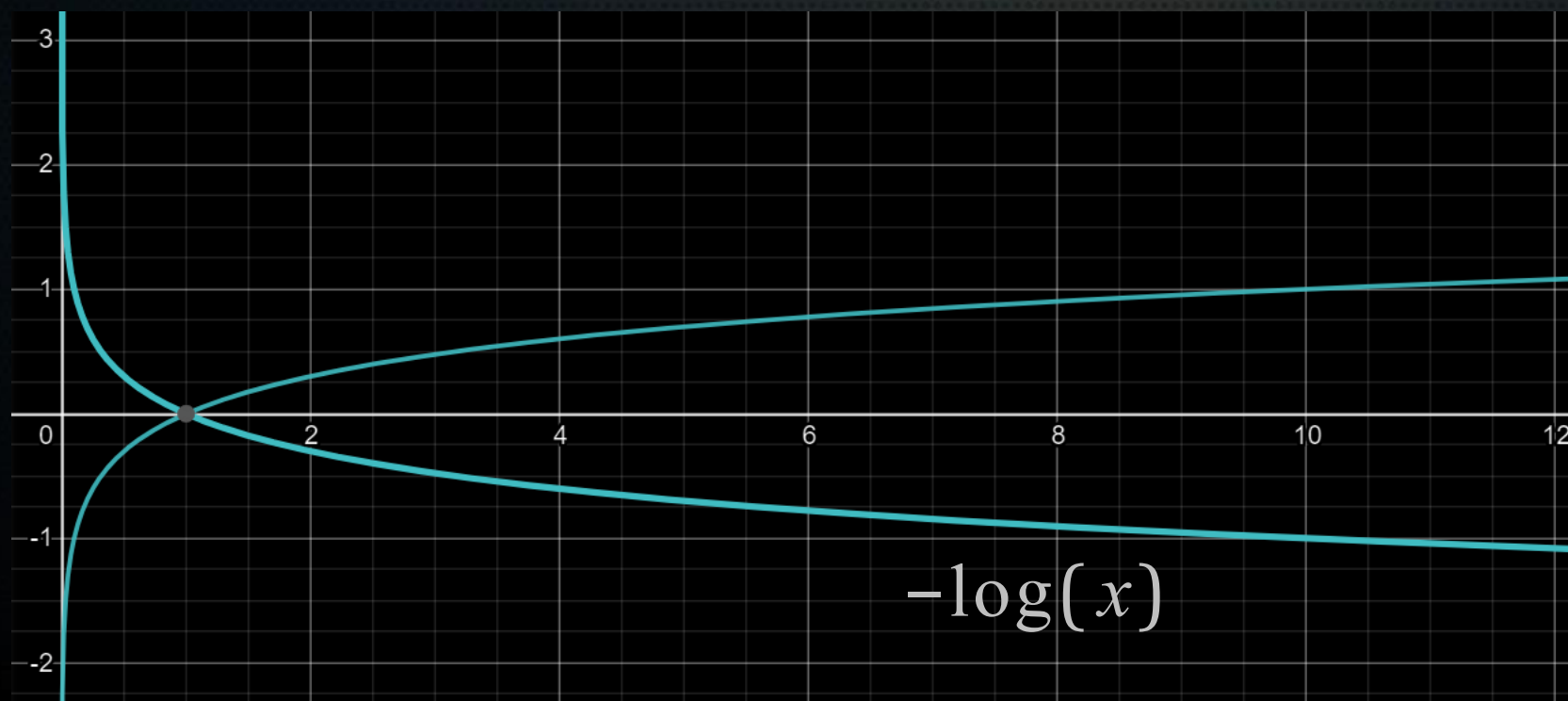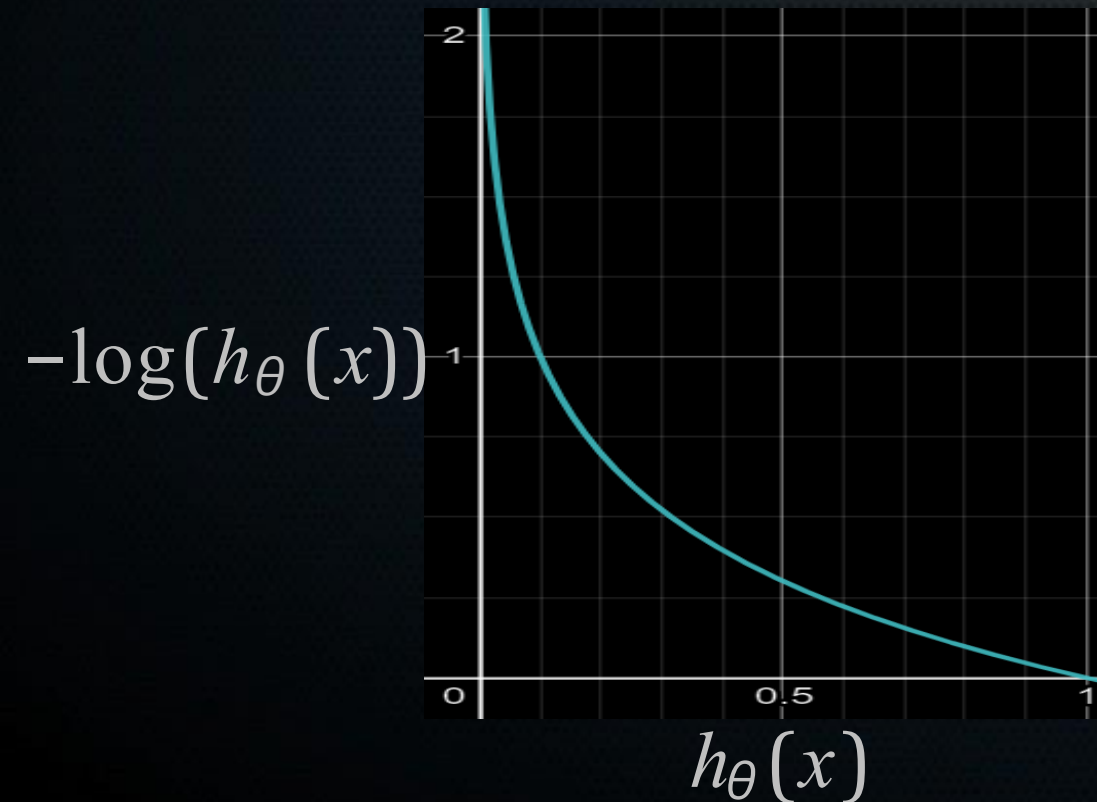
$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(x^i)$$

If y = 1
And our hypothesis predicts 0
The cost goes to $\infty$

If we had
used MSE

$(0-1)^2 = 1$

Penalised far less!

If y = 1
And our hypothesis predicts 1
The cost is 0

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

$h_\theta(x)$

# What does this function look like?

$$\text{Cost}(x) = \begin{cases} -\log(h_\theta(x)) \text{ if } y=1 \\ \boxed{-\log(1-h_\theta(x)) \text{ if } y=0} \end{cases}$$
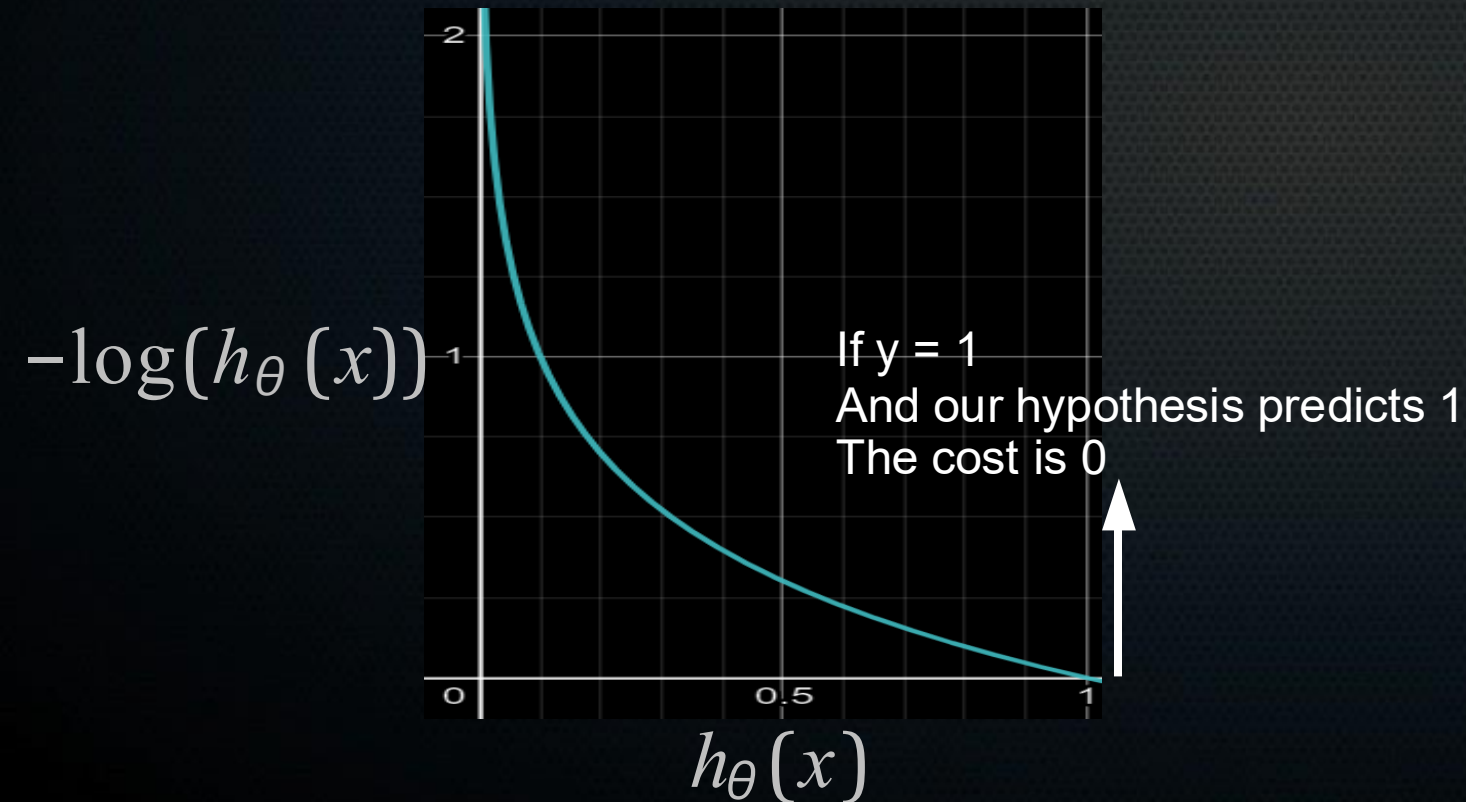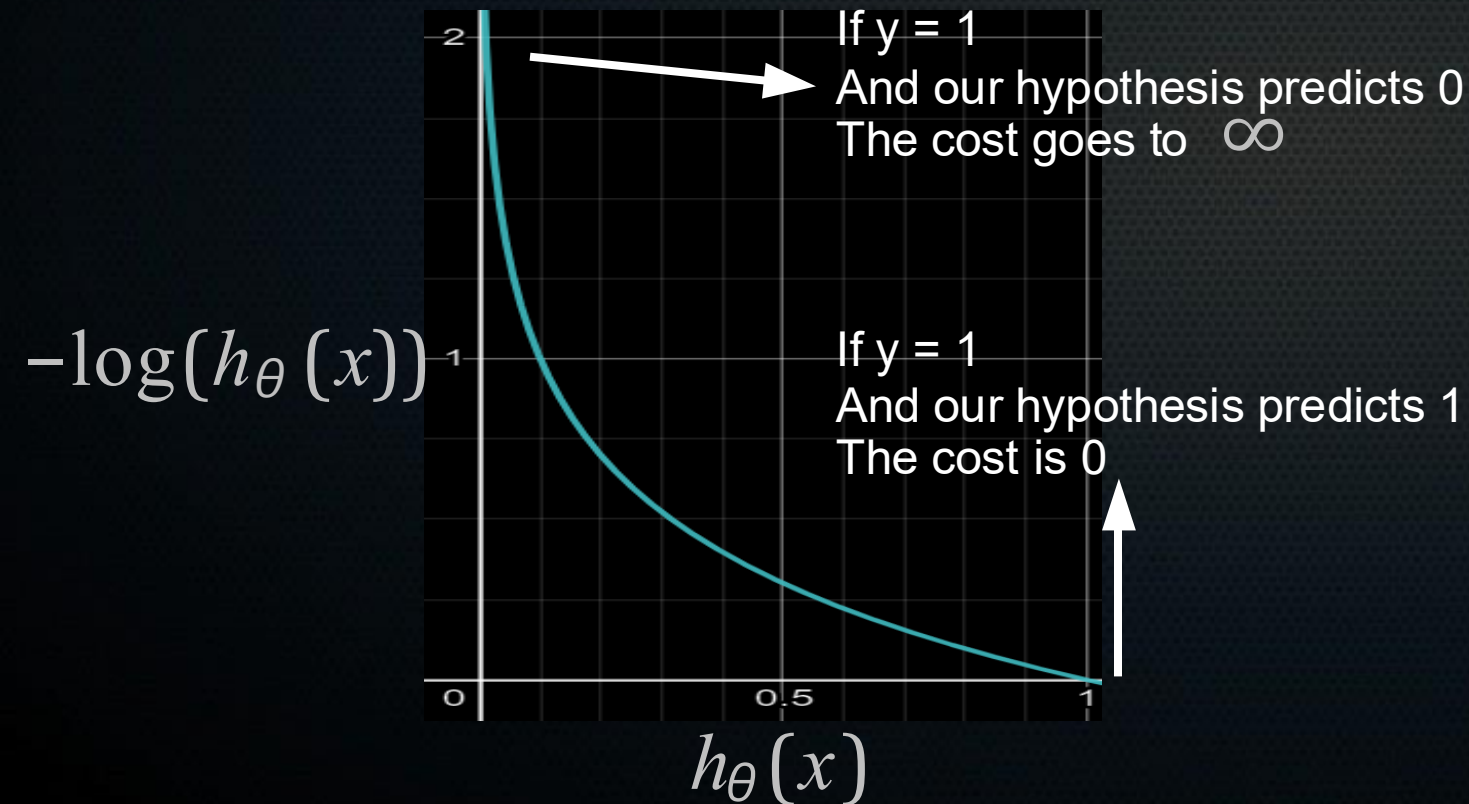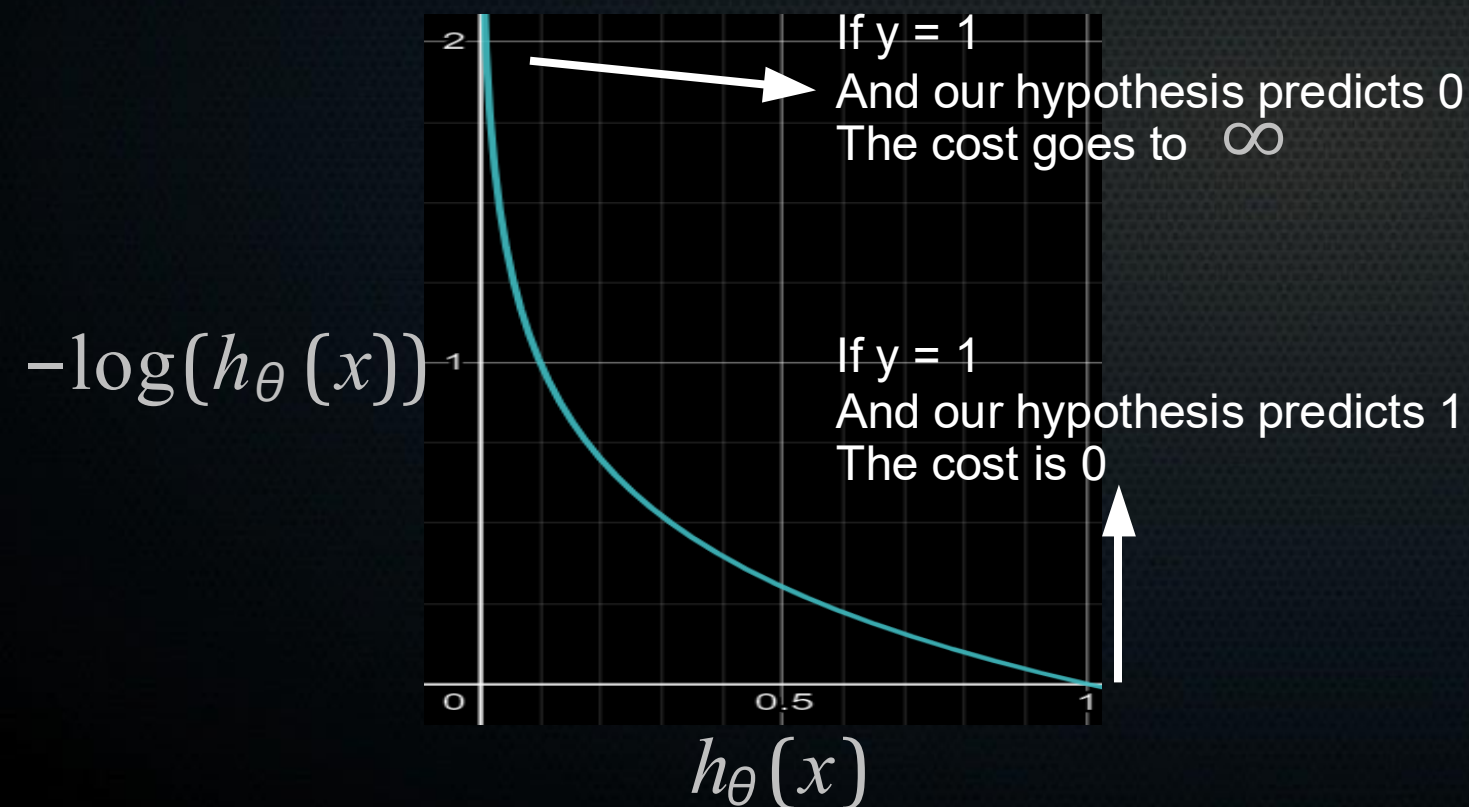
$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(x^i)$$

If y = 0
And our hypothesis predicts 1
The cost goes to ∞

If y = 0
And our hypothesis predicts 0
The cost is 0

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

$$h_\theta(x)$$

# Simplified notation

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

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

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

# Simplified notation

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

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

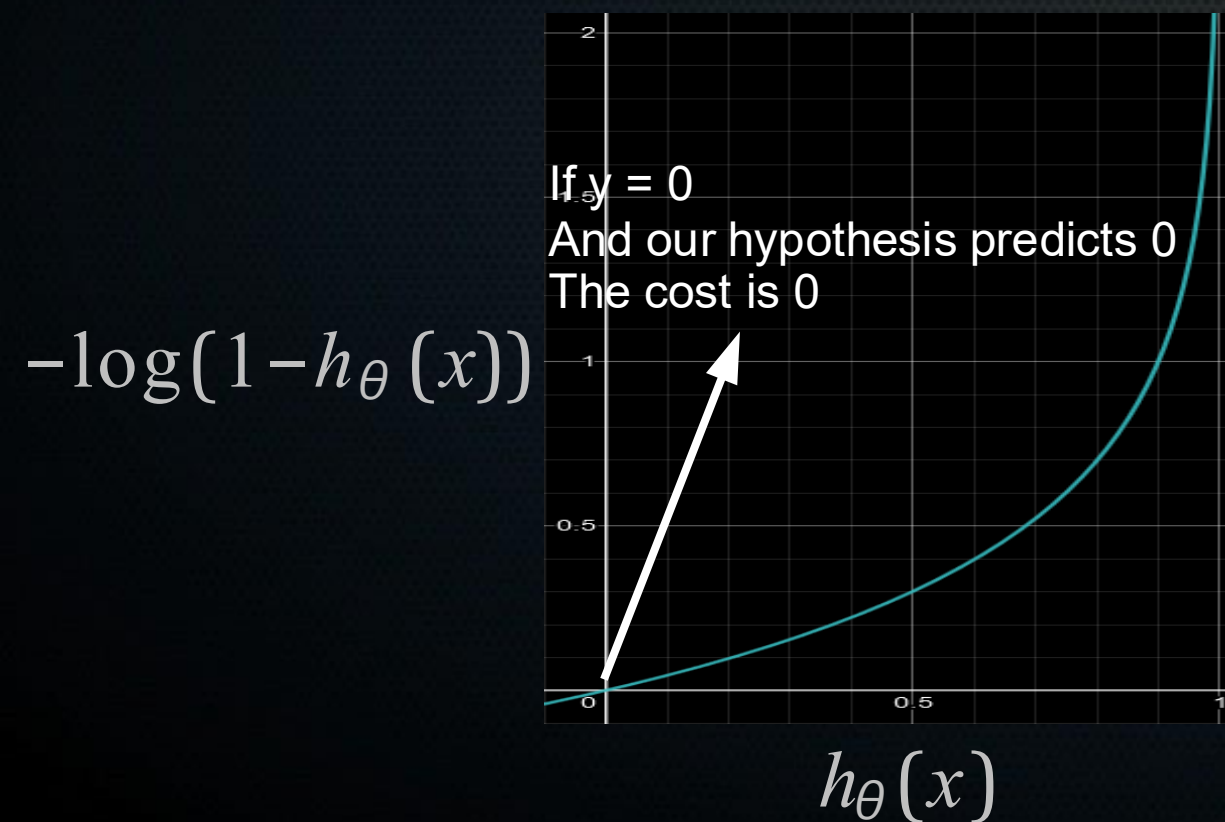$$\text{Cost}(x) = -y \cdot \log(h_\theta(x)) - (1-y) \cdot \log(1 - h_\theta(x))$$

$y = 0$

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

# Simplified notation

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

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

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

$y = 0$

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

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

# 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

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

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

# Putting it all together

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

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

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

# Optimising the cost function

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

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

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

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

# 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