

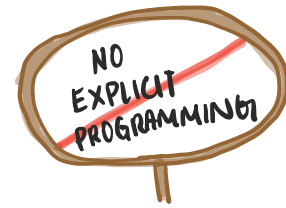
# Machine Learning Cartoons

Notes from Andrew Ng Coursera  
Stanford Machine Learning course



Illustrated by  
Patrick Tran

Q: What is machine learning?

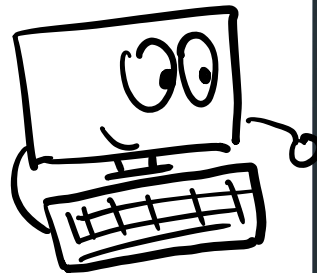


OLDER INFORMAL DEFINITION

The field of study that gives computers the ability to learn without being explicitly programmed

- Arthur Samuel

Example:  
Computer **explicitly** programmed to recognize a bear



```
while plugged_in:

    pixel = image[360][368]

    if pixel["color"] == "#8B4513" #brown
        return "brown bear"
    else
        return "a different color bear"
```



## Q: What is machine learning?

NEWER MORE MODERN DEFINITION

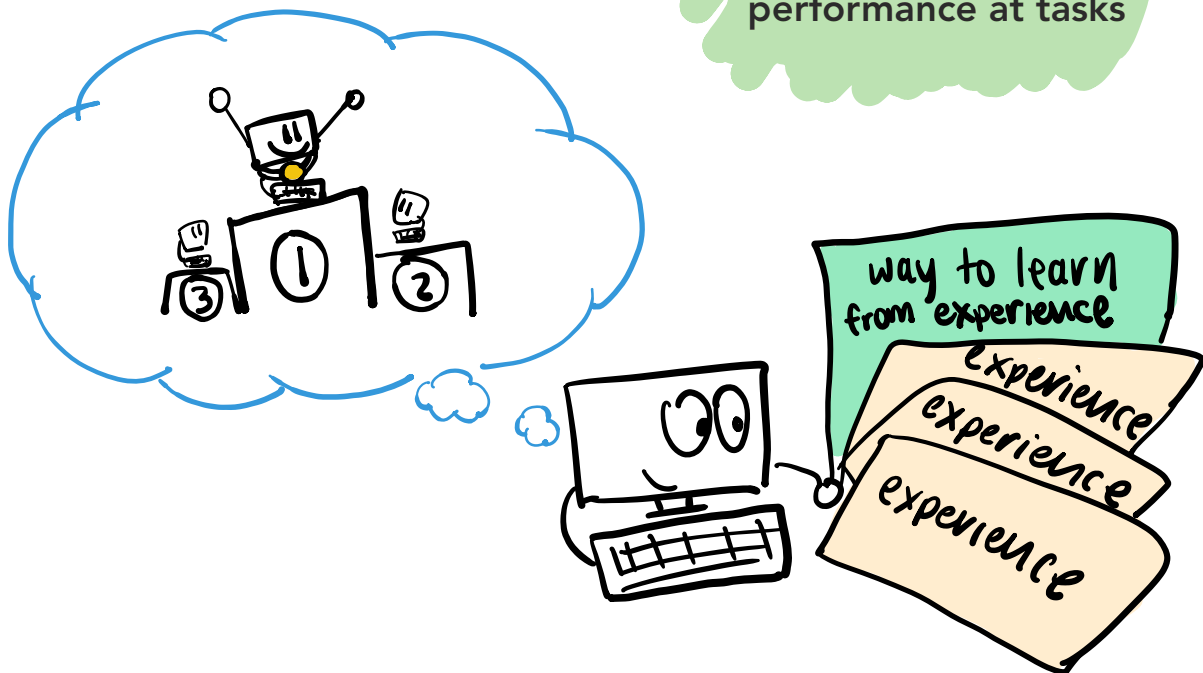
A computer program is said to learn  
from experience  $E$  with respect to some  
class of tasks  $T$  and performance  $P$

\_\_\_\_\_ IF \_\_\_\_\_  
its performance at tasks  $T$  as measure by  
 $P$  improves with experience  $E$

- Tom Mitchell

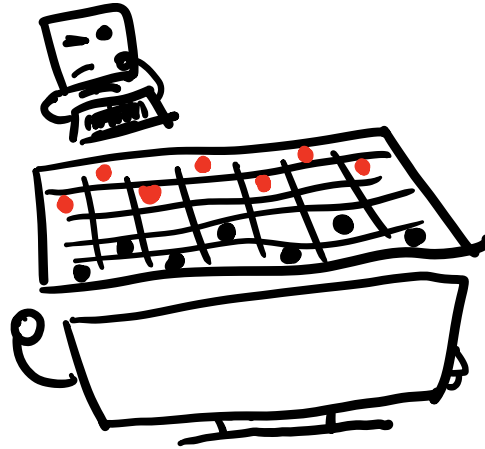
Example:  
Computer program learning

In plain English:  
More experience  
leads to higher  
performance at tasks



**Example:**

## Computer Program playing checkers



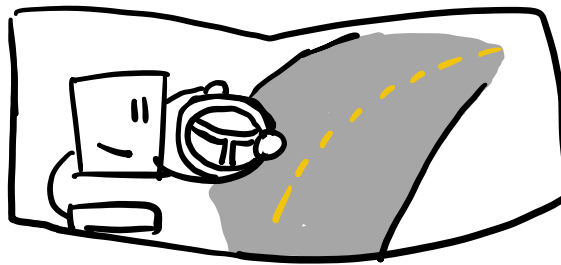
E (experience): playing many games of checkers (or data of games)

T (tasks) : task of playing checkers

p (performance): probability of winning

**Example:**

## Computer Program driving vehicles



E (experience): driving many roads (real and simulated)

T (tasks) : task of driving safely from A to B

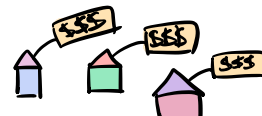
p (performance): number of accidents

This course covers supervised and unsupervised learning

## Supervised learning

We are given data with inputs and correct outputs. For example:

housing data (inputs)  
& house prices (outputs)



x1	x2	y
Size	Bedrooms	Price
1000	1	\$100k
1750	3	\$800k
1500	2	\$700k
⋮	⋮	⋮

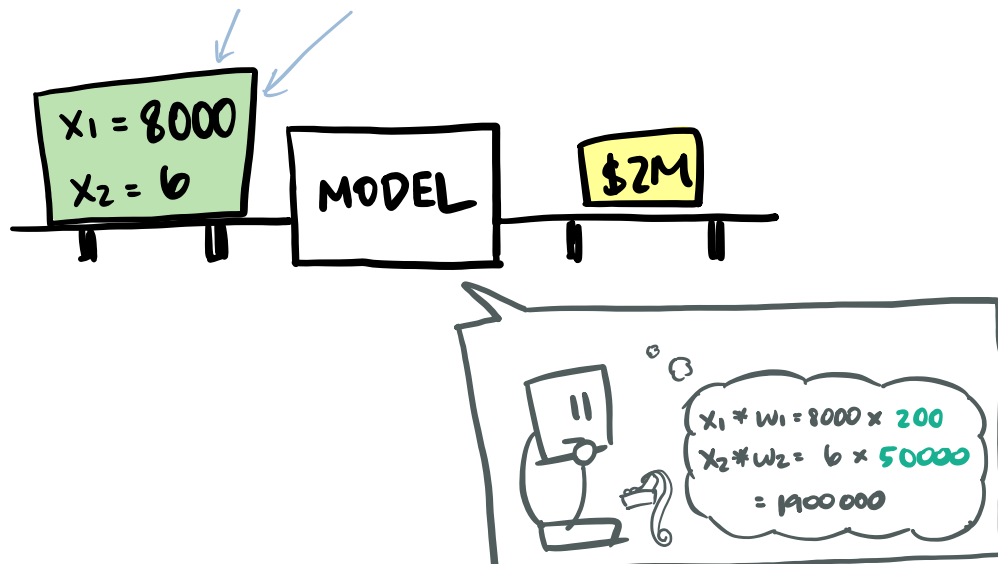
↑ these are inputs      ↑ output

This data is used to create a statistical model (or just "model") used to predict new data

①  $\text{train}(\text{inputs, output}) \Rightarrow \text{model}$

②  $\text{model}(\text{new housing data}) \Rightarrow \text{predicted house price}$

Say we wanted to predict the price for size=8000, bedrooms = 6

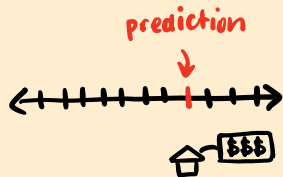


two main flavors 🍦 of


# supervised learning

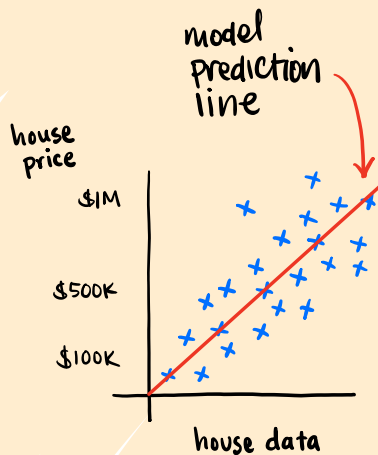
## linear regression

predict results in a continuous space



Example 1  
given: house data  
predicts: house price


Example 2  
given: photo of a person   
predicts: age of person  
Age: 43.287



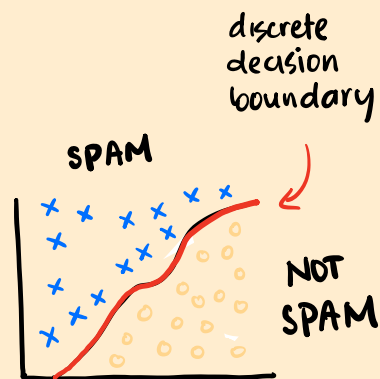
## classification

predicts results in a discrete space



Example 1  
given: email data   
predicts: spam or not spam

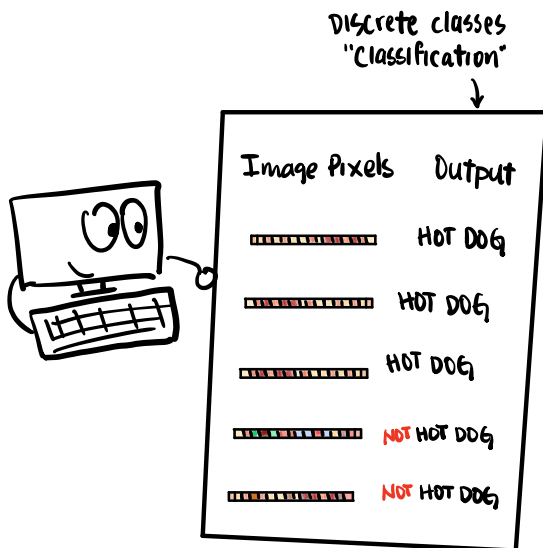
Example 2  
given: photo of tumor  
predicts: malignant or benign



# 🌸 ART IN LIFE 🌸

In an episode of HBO's Silicon Valley, the characters create an app called "Not Hot Dog" that detects whether an image is or isn't a hotdog.

The software engineers behind the show also created this app in real life using a total of 150,000 images to train their model to identify all types of hot dogs



This is a great example of logistic regression!  
Inputs are image pixels  
Output are the two discrete classes:  
1. Hot dog  
2. Not hot dog



two main flavors 🍦 of

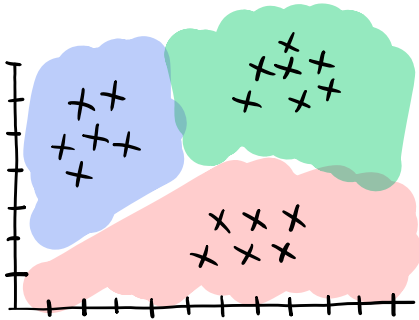
# Unsupervised Learning

( outputs are not in  
our provided data )



## clustering

"find clusters!"



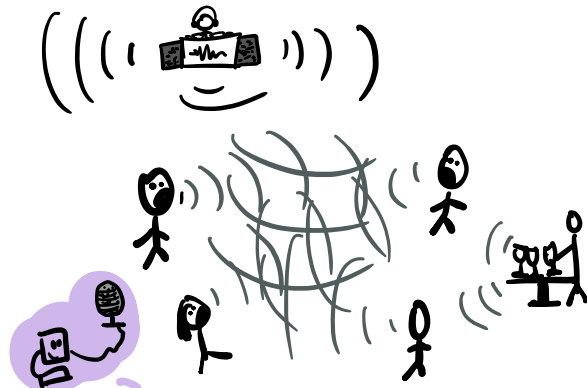
Given data, we use algorithms that find groupings or "clusters" to the data

id	location	lifespan	role
1	...	...	...
2	...	...	...
3	...	...	...
4	...	...	...
5	...	...	...
6	...	...	...
7	...	...	...
8	...	...	...
9	...	...	...
10	...	...	...



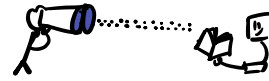
## not clustering

example: Cocktail Party Problem



This computer can listen to sounds and identify individual voices and music from a mesh of sounds

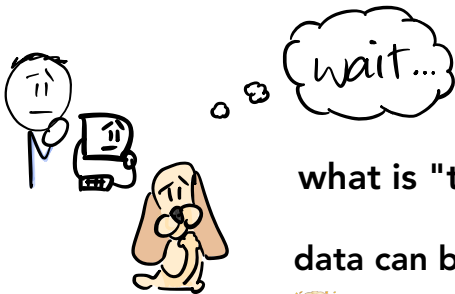
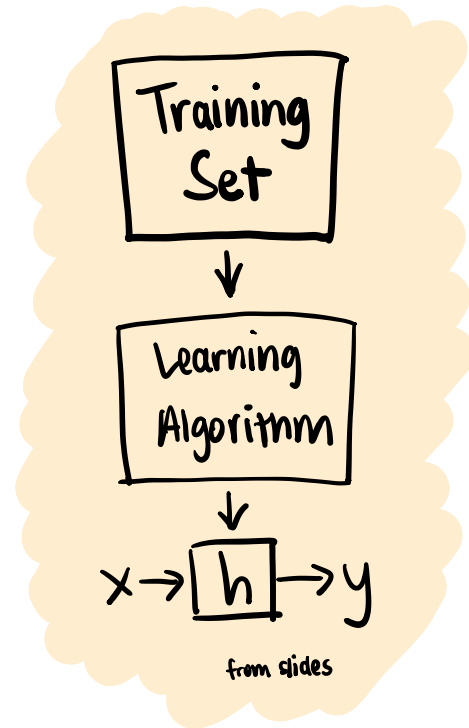
# Supervised Learning



In supervised learning, our goal is to learn a function given a set of data

$$h(x) \rightarrow y$$

Hypothesis function (h)  
given input (x)  
predicts output (y)



what is "training set"?

data can be used for training our model...

🤔 data we decide to use for training is "training data"

🤔 this set of data, our "dataset" can be referred to as our "training set"

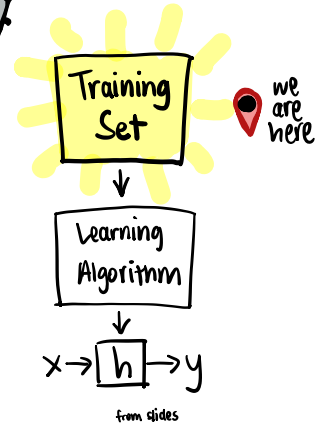
🤔 each individual data attribute of the data can be referred to as a "feature" (number of bedrooms and sq. feet are features of the data)

🤔 so sometimes the training dataset is also the feature dataset, feature data or simply our "features"

training set  
& feature set  
& training data  
& training dataset  
& feature set  
& feature data  
& features

# Training Set

Our data is our training set



Q: What makes up our data?

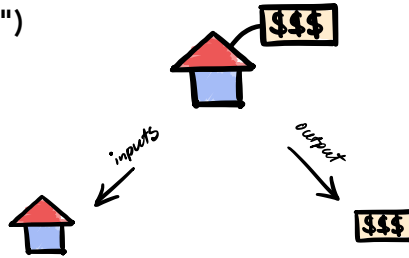
A: In supervised learning, our data has:

1. inputs (our "features")
2. outputs (our "labels")



For our housing example

inputs (features):  
 $x_1$  = square feet  
 $x_2$  = bedrooms



output (labels):  
 $y$  = house price

inputs		output
$x_1$ Size	$x_2$ Bedrooms	$y$ Price
1000	1	\$100k
1750	3	\$800k
1500	2	\$700k
⋮	⋮	⋮

There could be as many variables/features that your budget and computing power allows!

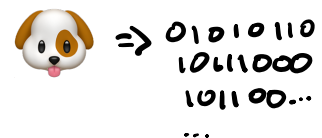
These outputs of price are in a continuous space so you can tell this data is fit for linear regression.

$(x^{(i)}, y^{(i)})$   
 This pair represents a single training example where the "i" represents the index

$m$  = number of training examples

Classification (logistic regression) would require "classes" as output ie. (can afford / cannot afford)

NOTE: all inputs can be converted to numeric, even images!

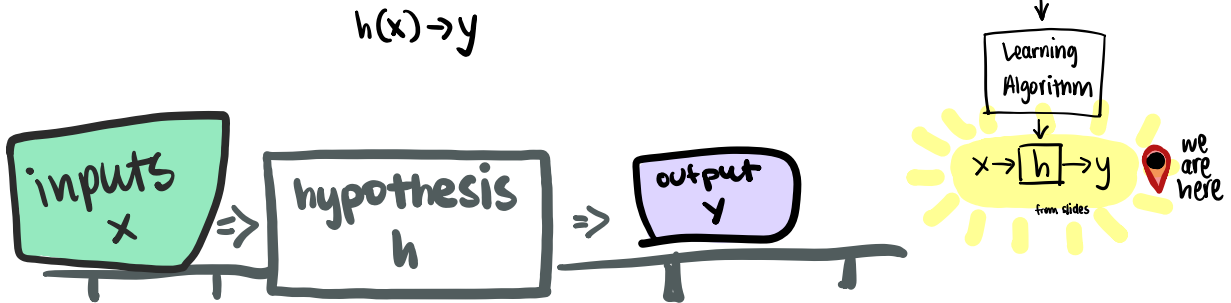




# HYPOTHESIS FUNCTION (h)

**Q:** what does the hypothesis function and look like?

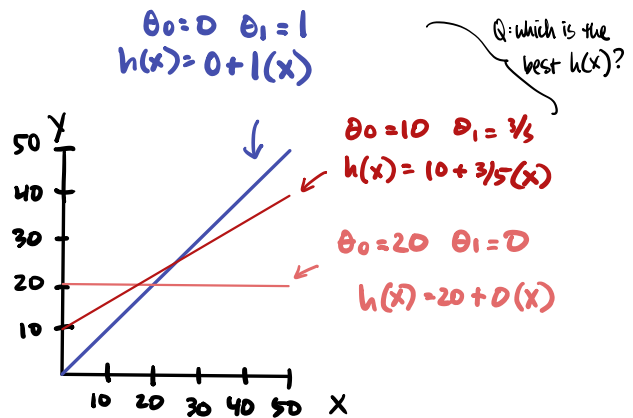
**A:** All data inputs are eventually converted to numbers and h is a function of these numbers.



Example h with one input (x) variable

$$h\theta = \theta_0 + \theta_1 x$$

h functions have a bias value ( $\theta_0$ ) we want to solve for  $\theta_0$  &  $\theta_1$  so that h(x) "fits" well with our training set



remember  $y=mx+b$  from high school?

Remember, our data can have one, two, three or even hundreds or thousands of inputs!

2  $\Rightarrow h\theta = \theta_0 + \theta_1 x_1 + \theta_2 x_2$

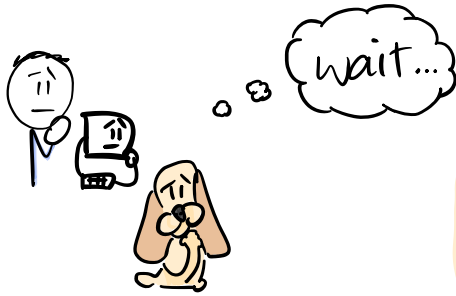
3  $\Rightarrow h\theta = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$

⋮

100  $\Rightarrow h\theta = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 + \theta_5 x_5$

↑ this many dimensions becomes hard to graph so this course usually visualizes 2 or 3 features





"we noticed  $h, h_\theta, h_\theta(x)$  being used"

These all refer to the hypothesis function  
 $h_\theta$  clarifies that  $h$  is parameterized by  $\theta$   
 $h_\theta(x)$  clarifies that it takes input  $(x)$

Note: hypothesis functions don't need to be simple functions

Depending on our domain knowledge about the data and the problem, we can create complex hypothesis functions

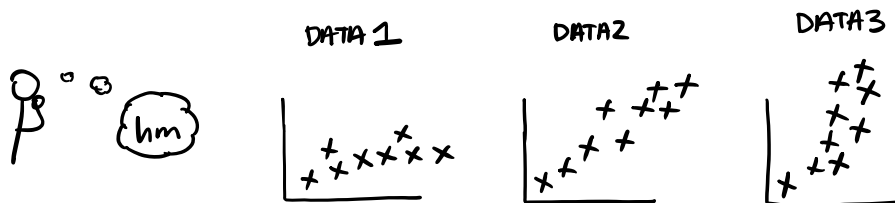
Example: two inputs with a special relationship

$$h_\theta = \theta_0 + \theta_1 x_1 + \theta_2 x_1 x_2 + \theta_3 x_2$$

Example: two inputs and higher order relationships

$$h_\theta = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$$

"whichever function we use will depend on our knowledge about the data"



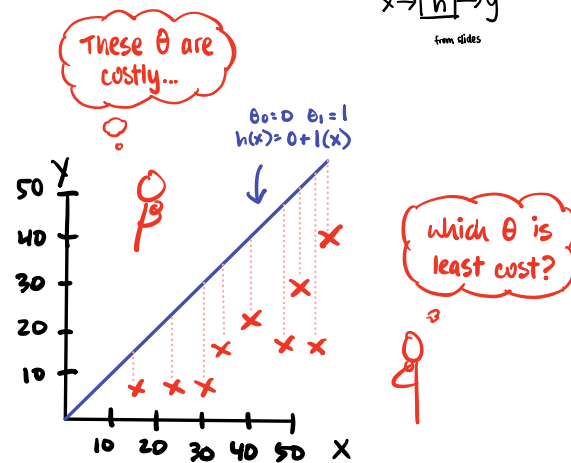
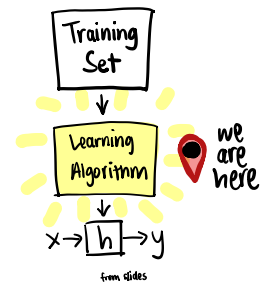
PROTIP: we can try different functions to see which one fits the best

# \$\$\$\$\$\$\$\$ COST function \$\$\$\$\$\$\$\$\$\$\$\$\$\$\$\$\$\$

Q: How do we evaluate our hypothesis function?

Q: First, why should we evaluate h?

A: We are using h for predictions and we should know how good or bad those predictions are?



A: We evaluate a hypothesis function with a “cost function”. Remember, in supervised learning our data includes the correct outputs we can use hypothesis function outputs. There are different cost functions but the most common one is “Average Squared Difference”\*

Our goal is to find theta ( $\theta$ ) that leads to the lowest cost ( $J(\theta)$ )

Linear Regression Cost Function

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

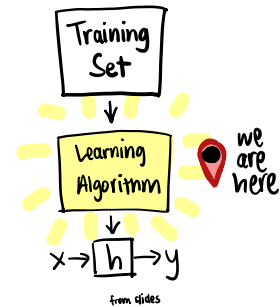
↑            ↑            ↑            ↑  
 averaging factor    sum    difference    square

Note: this function amplifies larger errors since an error of 1 is  $1^2=1$  while a diff of 5 is  $5^2=25$ !

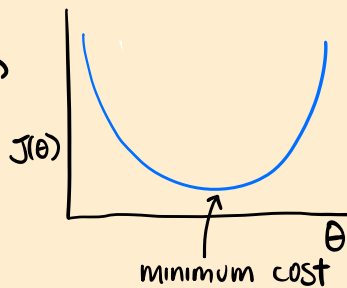
# Gradient Descent

Q: what is gradient descent (GD)?

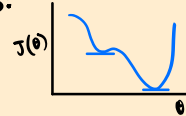
A: GD is an algorithm that allows us to find the  $\theta$  (theta) values that lead to the lowest  $J(\theta)$  (cost).



The cost function of linear regression  $\rightarrow$  could look like this



Note: linear regression will always be convex and doesn't have local optima like this:



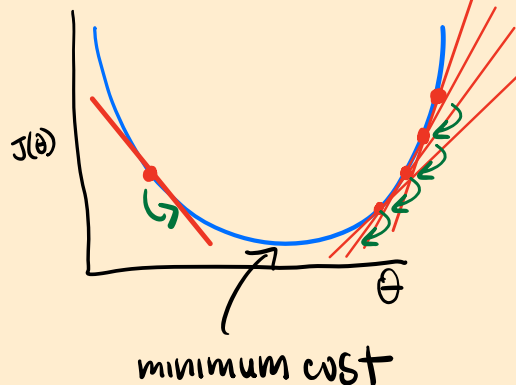
Let's use the derivative 

Q: How?

A: The derivative of the cost function informs us of the slope of the tangent line. with this we know which direction to "descend". The "learning rate" determines the step size.

Remember calculus?  
The derivative gives us a tangent line

red: tangent line  
green: learning rate



# Gradient Descent Algorithm

(one input example)

This is the condensed version with two thetas: one input theta and one bias theta

$$\text{repeat } \left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \end{array} \right. \dots$$

learning rate  $\downarrow$  partial derivative

$\uparrow$  perform simultaneously for all  $j$

Note: technically when there are more parameters we are using the partial derivative

This is the expanded partial derivative from calculus

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

Note: this a general version of the partial derivative

The final version of the algo has this expanded partial derivative and since  $x^{(0)}$  is always 1, we often see it disappear in the formula

## Gradient Descent Algorithm - one input

$$\text{repeat until convergence } \left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \\ \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \end{array} \right.$$

update simultaneously

week 2

WEEK2 WEEK2 WEEK2  
2 WEEK2 WEEK2 WE  
2 WEEK 2 WEEK 2 W  
K2 WEEK2 WEEK2 v  
EEK2 WEEK2 WEEK2  
WEEK2 WEEK2 WEEK  
WEEK2 WEEK2 WEEK  
2 WEEK2 WEEK2 W  
K2 WEEK2 WEEK2 W

# Multiple variables

Our data may contain

- (A) single feature or (B) multiple input features



here: size

X	Y
size	price
1000	\$100K
2000	\$200k
3000	\$300K
:	:

single feature

size?



you want me to predict price on this single feature???

here: size, #of bedrooms, #floors, age in years

X	x2	X3	x4	Y
size	bedroom	floors	age	price
1000	1	1	60	\$100K
2000	2	1	15	\$200k
3000	3	2	48	\$300K

$x^{(2)}$

$x_2^{(3)}$

multiple input features

size?

bedrooms?

floors?

age?



## Notation

$n$  = number of features

$x^{(i)}$  = features of  $i$ th training example

$x_j^{(i)}$  = value of feature  $j$  in  $i$ th training example

# Hypothesis Function

## Multi-variate + Vectorized

Single feature hypothesis function

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Multiple Variable hypothesis function

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 \quad \leftarrow \text{(four features)}$$

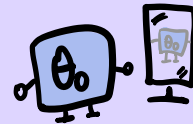
or generally

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

Note: we have more  $\theta$  variables than  $x$  because of the  $\theta_0$  bias value

In order to create a vectorized representation of  $h_{\theta}(x)$  we can assume an  $x_0$  value (which is always 1) Now our vectors both have length  $n + 1$

why is  $x_0 = 1$ ?



$x_0$  is just a placeholder and 1 multiplied by any value is that same value

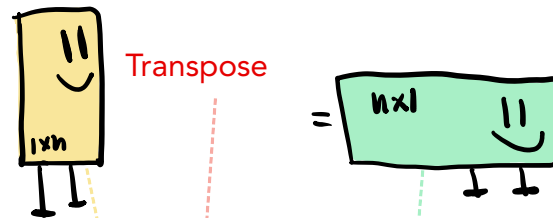
$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Multiplying these two vectors  $h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots$  can be represented as

Multiple feature hypothesis function

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

In order to multiply two column vectors we need to transpose a vector



for example

$$\theta^T = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}^T = [\theta_0 \ \theta_1 \ \theta_2 \ \dots \ \theta_n]$$



# Gradient Descent

Gradient Descent works similar going from one feature to multiple features

In essence, for a single feature, perform a simultaneously update for  $\theta$  in order to minimize cost.

Here, we only looked at  $\theta_0$  and  $\theta_1$

## Gradient Descent Algorithm - one input

repeat until convergence {

$$\begin{aligned} \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) \\ \theta_1 &:= \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) x^{(i)} \end{aligned}$$

update simultaneously

For multiple features, we are expanding to  $\theta_2, \theta_3, \dots$  all the way up to  $\theta_n$

## Gradient Descent Algorithm - multi input

repeat until convergence {

$$\begin{aligned} \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) \\ \theta_1 &:= \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) x^{(i)} \\ \theta_2 &:= \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) x^{(i)} \\ &\vdots \\ \theta_n &:= \theta_n - \alpha \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) x^{(i)} \end{aligned}$$

update simultaneously

we can more generally represent  $0, \dots, n$  as  $j$  and create this concise function

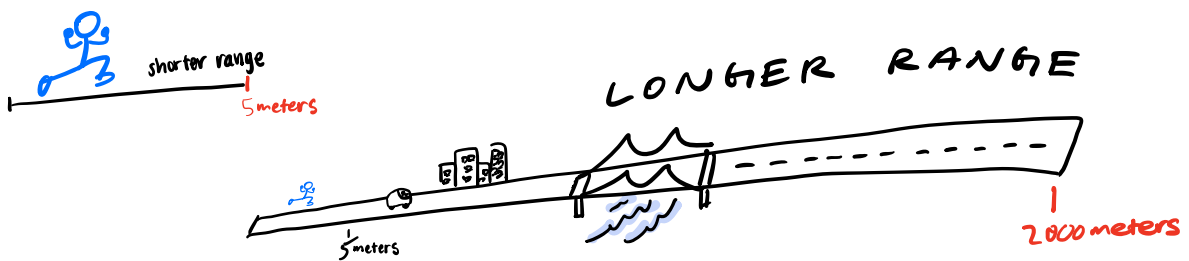
more generally

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

# Feature Scaling

## Gradient Descent in practice I

What: make sure our features are on a similar scale  
 Why: our goal is to make gradient descent run much faster as  $\theta$  will "descend" quickly on smaller ranges and slowly on long ranges

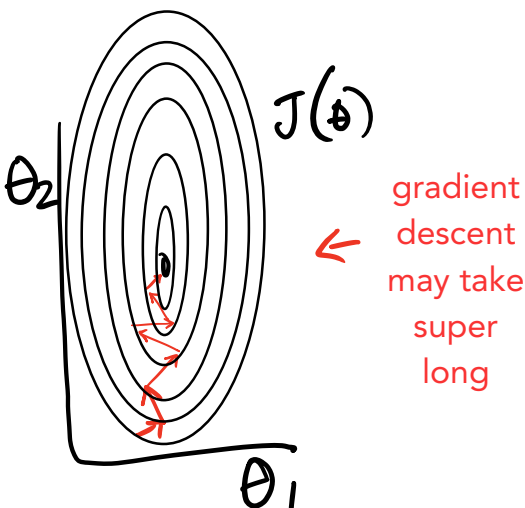


Say you have two features  
 $x_1 = \text{size (0-2000 feet)}$   
 $x_2 = \text{\#of bedrooms (1-5)}$

Feature scaling:  
 divides input values  
 by its range

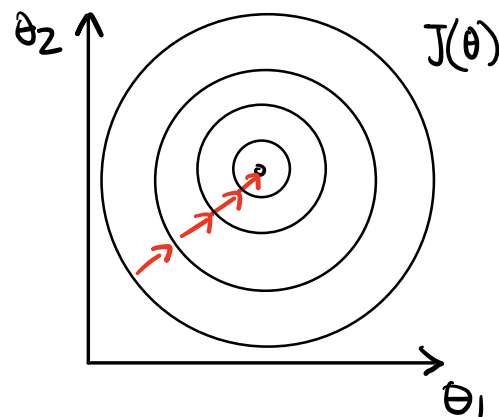
$$\frac{\text{value}}{\text{range}}$$

The difference in scale creates  
 a skewed cost curve graph



$$x_1 = \frac{\text{size}}{2000} \quad x_2 = \frac{\text{\# of bedrooms}}{2}$$

Which changes our cost curve to:



Two tricks are:

- ① feature scaling
- ② mean normalization

- ① Feature scaling involves dividing our feature value by its range in an attempt to shrink its range to

$$-1 \leq x \leq 1$$

Note: every researcher has their own rule of thumb for this range. Ng suggests  $-3 \leq x \leq 3$  and  $-1/3 \leq x \leq 1/3$  are also appropriate ranges

- ② Mean normalization is an additional option that replaces the feature value with feature value minus the mean so the new mean is roughly 0

How to **Mean Normalize**

In the above example if:  
Average size  $x_1 = 1000$   
Average bedrooms  $x_2 = 2$

$$\text{new (mean normalized) feature value} = \frac{\text{original feature value} - \text{mean}}{\text{range}}$$

$$x' \leftarrow \frac{x - \mu}{s}$$

$$x_1 = \frac{\text{size} - 1000}{2000}$$

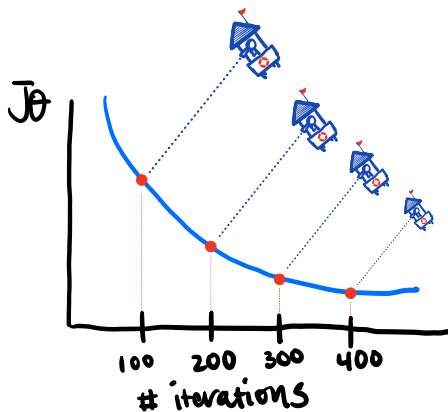
$$x_2 = \frac{\# \text{ bedrooms} - 2}{5}$$

Video: Gradient Descent in Practice II

Debugging: How do you know gradient descent is working correctly?

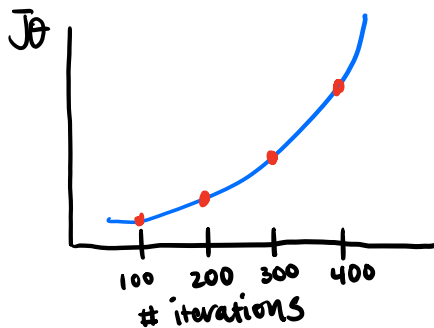
Sometimes GD never converges  
Sometimes GD has a slow convergence

Method I: plot of the cost function  $J(\theta)$  to its number of iterations



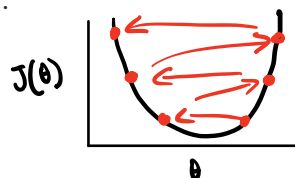
Use checkpoints at 100, 200, 300, 400 iterations to see if gradient descent is working properly.  
 $J(\theta)$  should decrease with every iteration

what if our graph looked like this?



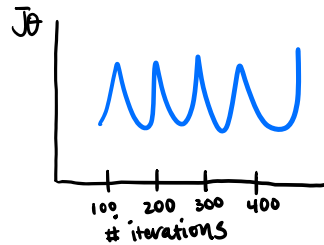
Solution:  
Use a smaller learning rate alpha

This means each iteration is making our  $J(\theta)$  cost larger. This could mean that in our convex curve...



the learning rate alpha is causing our us to overshoot the minimum and actually increasing  $J(\theta)$

what if our graph looked like this?



Solution:  
Use a smaller  
learning rate  
alpha

## Learning Rate ( $\alpha$ )

small  $\alpha$

- for sufficiently small  $\alpha$ ,  $J(\theta)$  should decrease on every iteration
- BUT if  $\alpha$  is too small, gradient descent can be slow to converge

# LARGE $\alpha$

- if alpha is too large,  $J(\theta)$  may not decrease on every iteration AND it may never converge

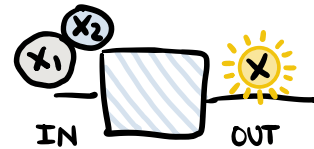
To choose an alpha try  
different values:

..., 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, ...

$\nearrow$   
 $\approx \times 3$       $\nearrow$   
 $\times 3$       $\nearrow$   
 $\times 3$       $\nearrow$   
 $\times 3$

# FEATURES

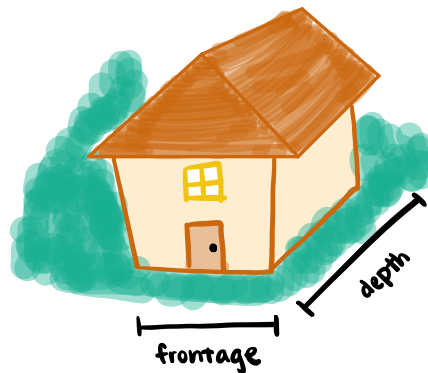
Depending on your domain knowledge about the problem, sometimes defining new features may lead to a better model



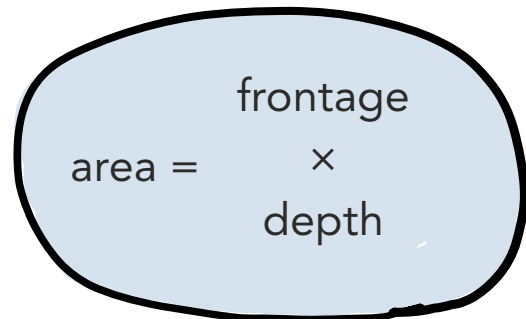
Given two features

$x_1 = \text{frontage}$

$x_2 = \text{depth}$

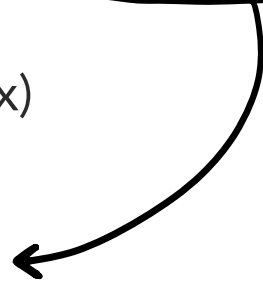


Based on your real estate knowledge you know that "area" is a better predictor of price



Now we can use a better predicting single feature  $h\theta(x)$

$$h\theta(x) = \theta_0 + \theta_1 x_1$$

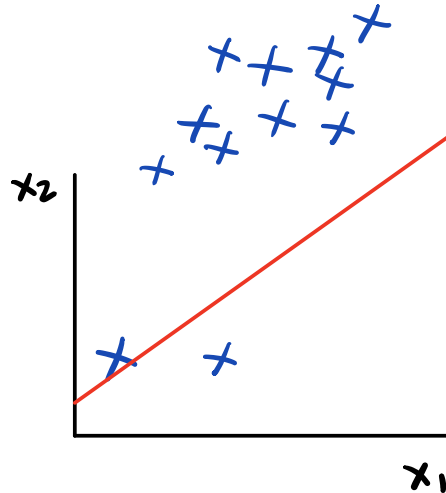


# POLYNOMIAL REGRESSION

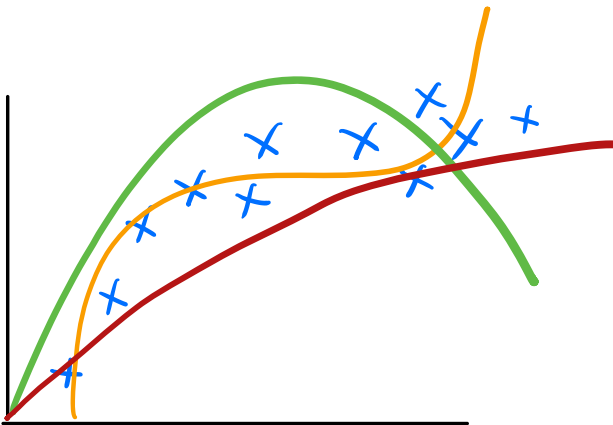
we have been working with straight lines so far

$$h\theta = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

no matter what theta is this will represent a straight fitting line through the data



Depending on the data we may want to use a polynomial function



quadratic

$$\theta_0 + \theta_1 x + \theta_2 x^2$$

OR

cubic

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

OR

square root!

$$\theta_0 + \theta_1 (\text{size}) + \theta_2 \sqrt{\text{size}}$$

Note that feature scaling becomes very important now since these ranges get exponentially large!

for scale, if:

$$\text{size} = 1-1,000$$

$$\text{size}^2 = 1-1,000,000$$

$$\text{size}^3 = 1-1,000,000,000$$

Parting words:

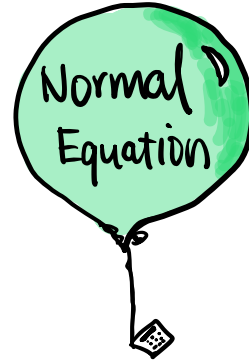
“ This may be bewildering!  
Which features do I use?  
which equation for the  
hypothesis function?  
Later in the course we  
talk about algorithms  
that will choose features.  
For now just know that you  
can choose different features  
and equations when your  
data calls for it

”



## Lecture: Normal Equation

we have been using gradient descent so far, alternatively we can use the "normal equation" to solve for theta analytically.



Note: this lecture does not prove why the normal eq works, just how to use it and when to use it

HOW:

$$\text{Normal Equation: } (X^T X)^{-1} X^T y$$

where  $X$  = inputs

$Y$  = outputs

$m$  = # training examples

$n$  = # features

10:51 AM Sat Dec 7

Normal Equation

Examples:  $m = 4$ .

	Size (feet <sup>2</sup> )	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
$x_0$	$x_1$	$x_2$	$x_3$	$x_4$	$y$
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$   
 $m \times (n+1)$

$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$   
 $m$ -dimensional vector

$\theta = (X^T X)^{-1} X^T y$

Andrew Ng

where  $(X^T X)^{-1}$  is the inverse of  $X^T X$

NORMAL EQUATION IN OCTAVE:

$$\text{pinv}(X' * X) * X' * y$$

When should we use either?

	Gradient Descent	Normal Equation
Advantage	<ul style="list-style-type: none"><li>- Works well even when n is large</li></ul>	<ul style="list-style-type: none"><li>- No need to choose learning rate <math>\alpha</math></li><li>- No need to iterate</li></ul>
Disadvantage	<ul style="list-style-type: none"><li>- Need to choose learning rate <math>\alpha</math></li><li>- Needs many iterations</li></ul>	<ul style="list-style-type: none"><li>- Need to compute <math>(X^T X)^{-1}</math></li><li>- Slow if n is very large</li></ul>

## Normal Equation and Non-invertibility

Q: What if  $X^T X$  is non-invertible?

A: This should happen very rarely...  
But this may be possible because

① Redundant features

Example: you have one feature size (in feet)  
and another feature size (in meters)

Solution: delete one of these features

② Too many features

solution: delete features

week 3

# Logistic Regression (Classification)

Classification (binary) Examples:

Email: spam, not spam  
 Online transaction, fraud?: Yes, no  
 Tumor: malignant, benign  
 Photo: hot dog, not hot dog

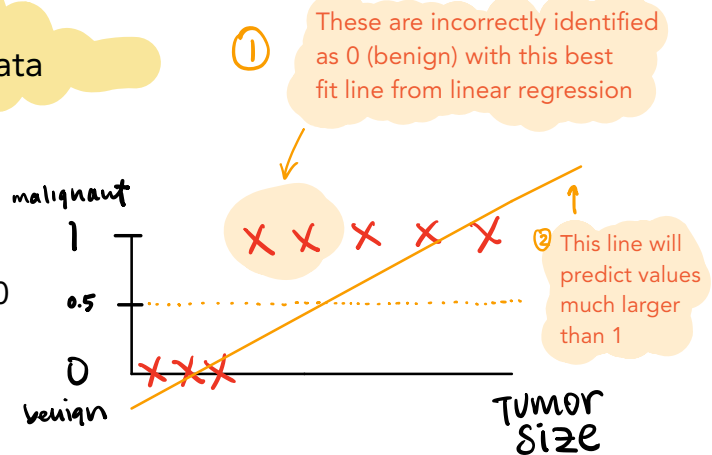
More concisely, 0 or 1  
 Given features  $X$ ,  $h_{\theta}(x) = \{1, 0\}$



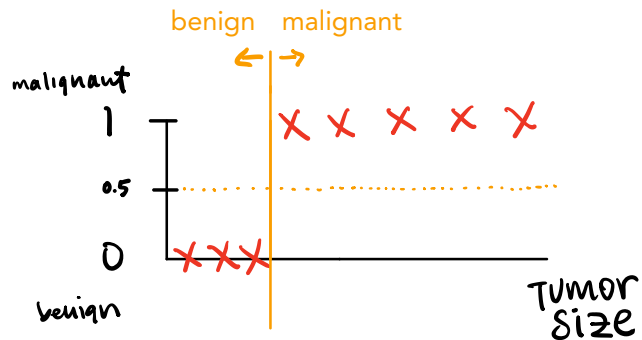
Note:  $\{1, 0\}$  is binary and can be interpreted as {Yes, No}, {SPAM, NOT SPAM},  $\{X, O\}$  {HOT DOG, NOT HOT DOG}

## Looking at a line through data

Here we use the best fit line from linear regression:  
 we may be tempted to interpret  $h_{\theta}(x) > 0.5 = 1$  and  $h_{\theta}(x) < 0.5 = 0$



we want to go from best fit to best split



Intuition:

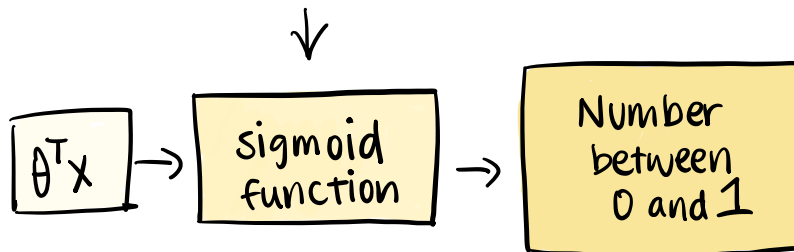
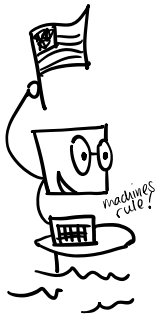
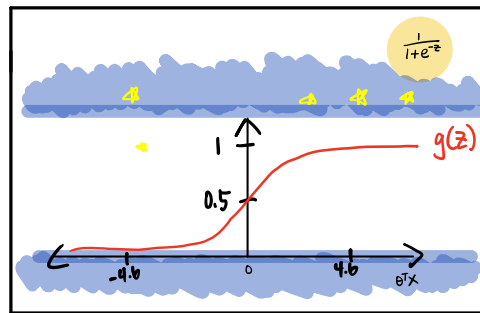
This was our linear regression hypothesis function

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

However for **logistic regression**

$h_{\theta}(x)$  should only return  $\{0,1\}$

so we use a "sigmoid" function ( $g$ ) otherwise known as a "logistic" function to take  $\theta^T x$  and fit it to an "S" curve



### Logistic Regression Hypothesis Function

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

where

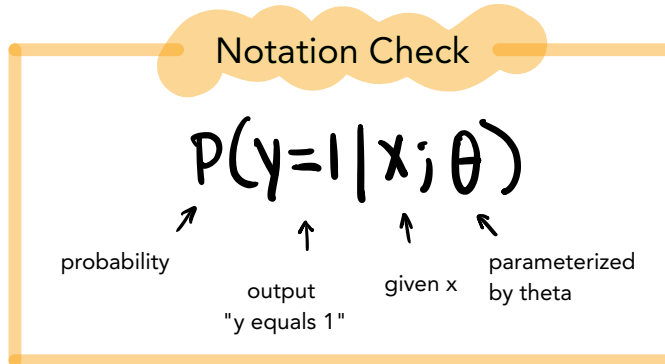
$$g(z) = \frac{1}{1 + e^{-z}}$$

one liner:

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




In logistic regression we can think of  $h\theta(x)$  as the probability that  $y = 1$  so if  $h\theta(x) = 0.7$  then there is a probability of 70% that  $y=1$



### More about probability algebra

There is a 100% probability that  $y = 1$  or  $y = 0$   
Therefore:

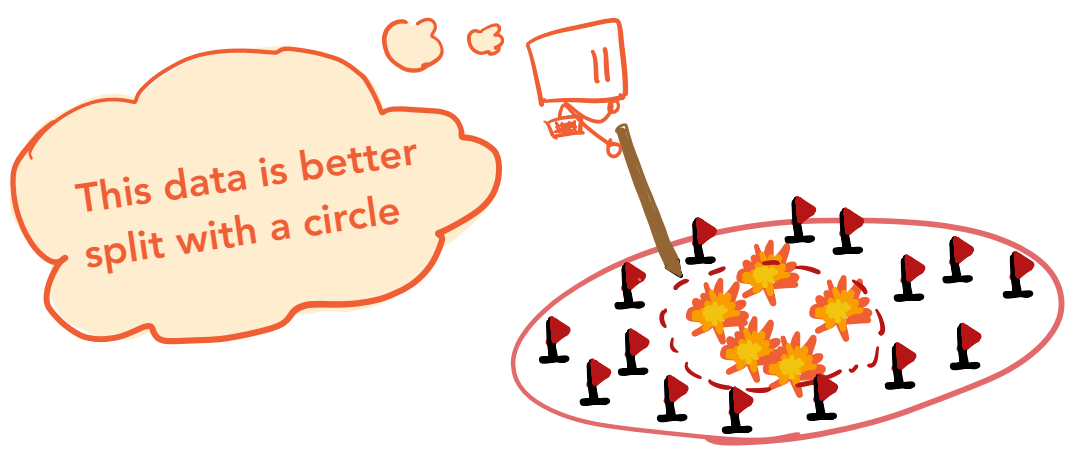
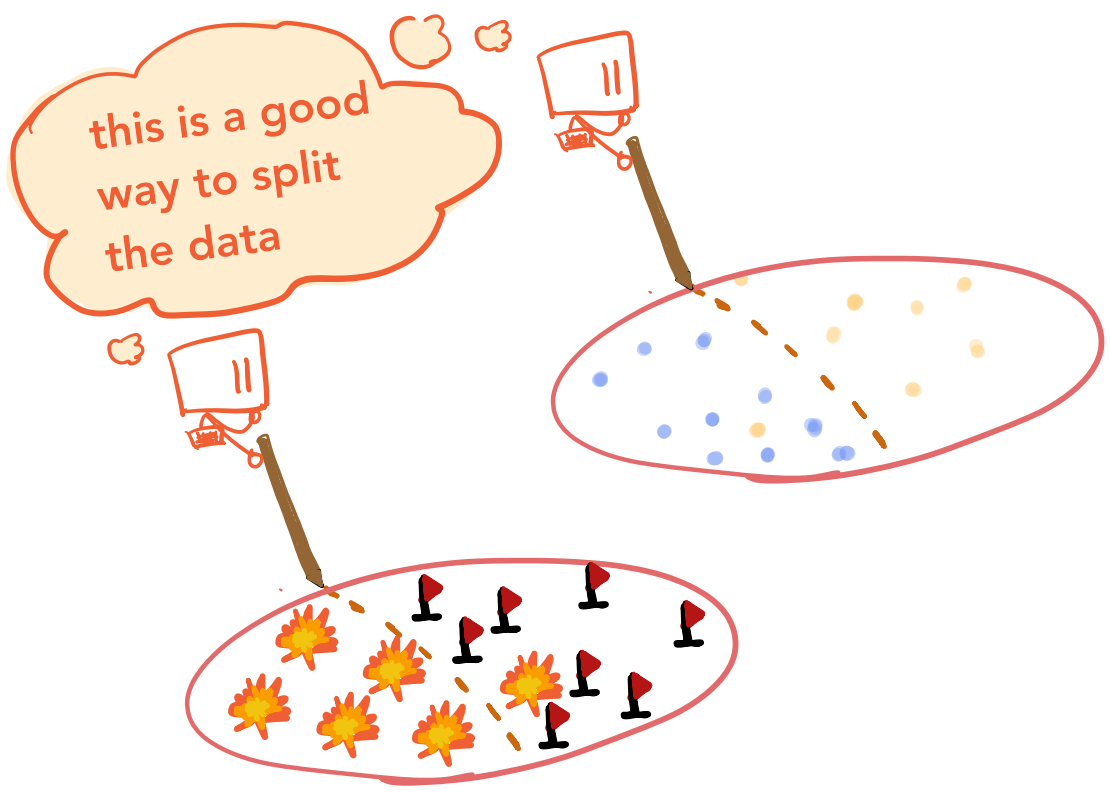
$$P(y=1) + P(y=0) = 1$$


$$P(\text{DO}) + P(\text{DO NOT}) = 1$$
$$P(\text{SPAM}) + P(\text{NOT SPAM}) = 1$$
$$P(\text{🌯}) + P(\text{NOT 🌯}) = 1$$
$$P(\text{NOT 🌯}) = 1 - P(\text{🌯})$$

Note: once you know  $P(y=1)$  or  $P(y=0)$  you can derive the other through algebra!

# Decision Boundary

Finding the boundary to best split the data





## Q: How does $h_{\theta}(x)$ represent a decision boundary?

Say we have...

this hypothesis function and this data

Hypothesis Fn.

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

Data

$x_1$	$x_2$	$y$
1	0	0
0.5	1	0
2.5	2.5	1
1.5	1	0
1	3	1
3	0.5	1
0	2	0
...	...	...



we need to solve for  $\theta$  (theta) that creates the best decision boundary

★ we learn how to solve for theta in a later lecture

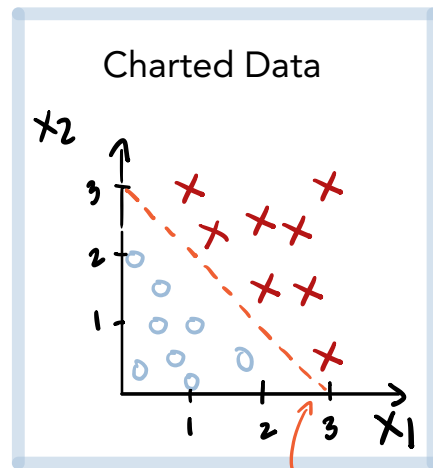
Imagine we solved for these theta values:

$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \text{ so } \begin{matrix} \theta_0 = -3 \\ \theta_1 = 1 \\ \theta_2 = 1 \end{matrix}$$

Hypothesis Fn.

$$h_{\theta}(x) = -3 + 1(x_1) + 1(x_2)$$

and predict  $y=1$  if  $h_{\theta}(x) \geq 0$



This is the decision boundary solved by finding theta

# LOGISTIC REGRESSION COST FUNCTION

Q: Can we reuse the linear regression cost function as the logistic regression cost function?

A: No

Q: Why not?

First lets look at our linear regression cost function:

$$\text{linear regression } J\theta = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h\theta(x^{(i)}) - y^{(i)})^2$$

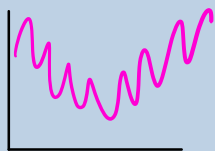
A: Since we have a higher order (non-linear) hypothesis function

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

IF we plugged it into this function our cost graph will be non-convex

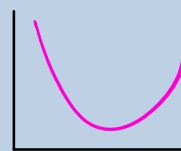
## NON CONVEX

This shape means many local optima; gradient descent will struggle to find the best option



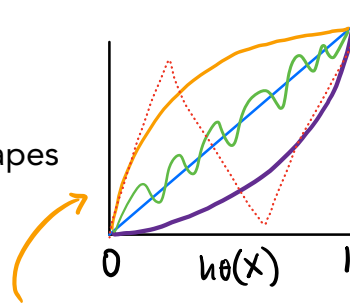
## CONVEX

Our goal is to get this pretty looking convex shape that gradient descent can help with



**Q: what should the shape of our cost function be?**

if  $y = 0$   
our cost graph for  $h\theta(x)$   
could be ya of these shapes



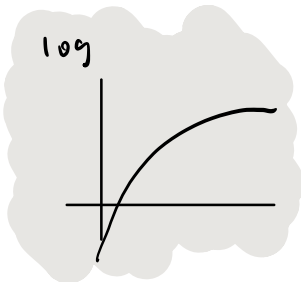
② We know that our cost should go up if  $y=0$  and  $h\theta(x)$  goes toward 1. Our cost should go up as our prediction is more incorrect. If  $y=0$ ,  $h\theta(x) = 0.5$  should "cost" more than  $h\theta(x) = 0.75$

① we know if  $h\theta(x)$  and  $y=0$  we want the cost to be very very low or zero

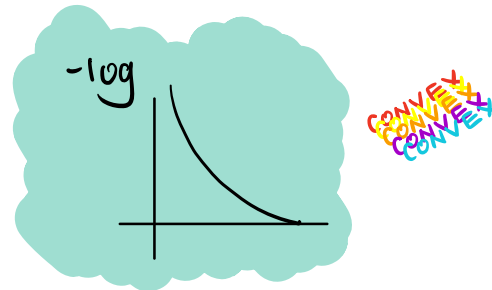
③ We know that this cost curve SHAPE needs to be convex. so gradient descent can work its magic.

**A: The log and -log shapes are great shapes**

This is the log function graph



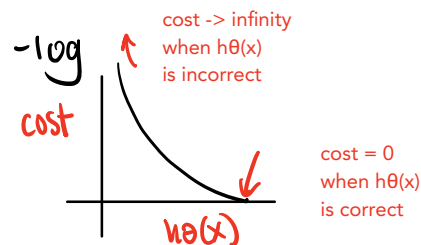
This is the negative log function graph



Why are these shapes great?

Why is this function great for modeling our cost curve?

- They are convex
- Can have zero cost
- We can model cost to approach infinity (helps to highly penalize wrong predictions)



## Logistic Regression Cost Function

TLDR →

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

### Intuition

if  $y=0$ ,  
we want the correct prediction to "cost" less  
than incorrect predictions

so...

$h\theta(x) = 1$  to "cost" less than  $h\theta(x) = 0$

### A TALE OF TWO HYPOTHESIS PREDICTIONS

①

Different parameters lead to different predictions

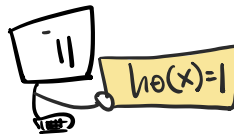


ANSWER



②

We set up the cost function to make mistakes more costly



ANSWER



③

NICELY DONE

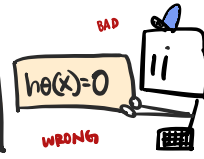
NO COST

ANSWER

OH NOEZ... MISTAKE

cost=0

GOOD



cost → ∞ !!!

WRONG

BAD

In notation

cost = 0 if  $y=1$  &  $h_{\theta}(x)=1$

but as  $h_{\theta}(x) \rightarrow 0$

cost  $\rightarrow \infty$

$h_{\theta}(x)$  approaches zero

( $h_{\theta}(x)$  approaches zero)

(cost approaches infinity)

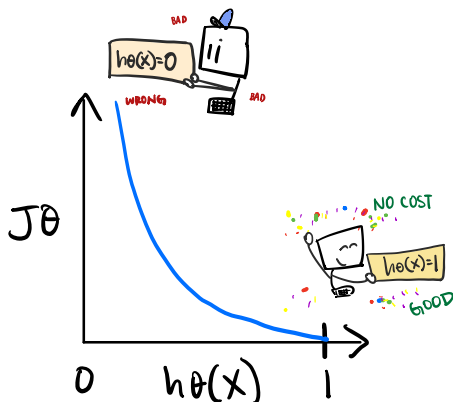
## Visualizing the cost graph

### Logistic Regression Cost Function

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

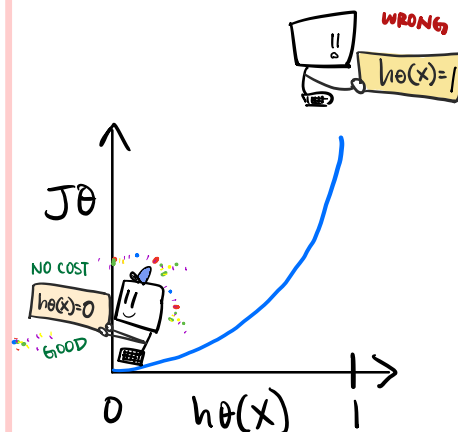
if  $y=1$

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



if  $y=0$

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



★ Note that  $1-h_{\theta}(x)$  flips the graph because  $h_{\theta}(x)$  will be between 0 and 1



## Converting our cost function to one line for $J(\theta)$

logistic Regression cost Function

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

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

$$\begin{aligned} \text{cost}(h\theta(x), y) &= \overbrace{(y * -\log(h\theta(x)))}^{\text{if } y=0 \text{ this portion} = 0} + \overbrace{(1-y) * (-\log(1-h\theta(x)))}^{\text{if } y=1 \text{ this portion} = 0} \\ &= \begin{matrix} \text{rearrange} \\ \text{(-) sign} \end{matrix} \quad \begin{matrix} \text{rearrange} \\ \text{(-) sign} \end{matrix} \\ &= - (y * \log(h\theta(x))) + - ((1-y) * (\log(1-h\theta(x)))) \\ &= \begin{matrix} \text{extract} \\ \text{(-) sign} \end{matrix} \quad \begin{matrix} \text{remember} \\ \text{transitive property} \\ -1(x+y) = -x + (-y) \end{matrix} \\ &= - \left( y * \log(h\theta(x)) + (1-y) * (\log(1-h\theta(x))) \right) \end{aligned}$$

one line!

We use this in  $J\theta$  to measure average cost

$$J\theta = \frac{1}{m} \sum_{i=1}^m \text{cost}(h\theta(x^{(i)}), y^{(i)})$$

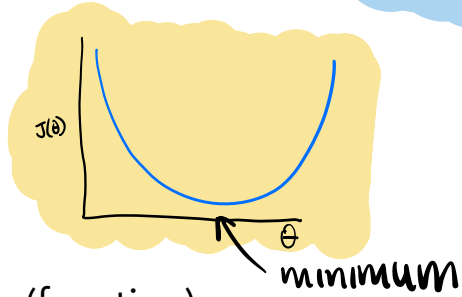
$$J\theta = \frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \log h\theta(x^{(i)}) + (1-y^{(i)}) \log(1-h\theta(x^{(i)})) \right]$$

# Minimizing Cost $J(\theta)$

Like always, given parameters theta ( $\theta$ ) we want to minimize cost



$$\min_{\theta} J(\theta)$$

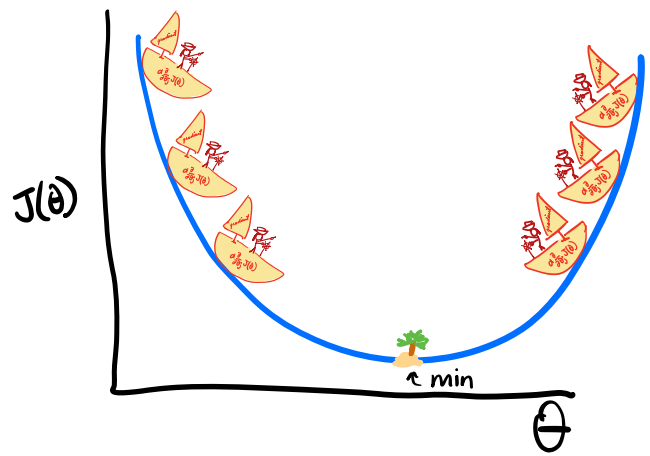
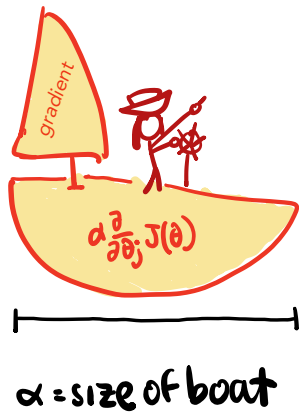


we have a convex shape (function) and we want to find the min value

## Welcome Back Gradient Descent

$$\text{Repeat } \left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \end{array} \right. \leftarrow \text{until convergence}$$

If the gradient step were a boat...



the gradient boat leads us to the minimum cost!

Expanding  $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$  to something usable

Think back to calculus class

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

Note:  $\theta_j$  is parameterized by "j" and we have the partial derivative

$$\frac{\partial}{\partial \theta_j}$$

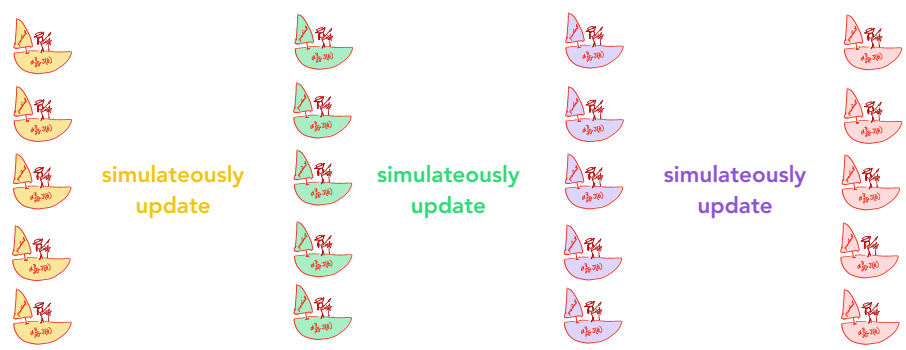
The chain rule from calculus says  $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$  can be rewritten as

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

### Gradient Descent formula

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (\text{Simultaneously update all } \theta_j)$$

If gradient descent were a group of gradient step boat...





## Lecture: Advanced optimization

### ★ gradient descent alternatives

1. conjugate descent
  2. BFGS
  3. L-BFGS
- \* the details of these three are outside the scope of this course

#### Advantages

- no need to manually pick a learning rate ( $\alpha$ ) (these algo's will choose  $\alpha$  for you)
- often faster than gradient descent

#### Disadvantages

- more complex

Recommendation:

Do not write these algos yourself  
you can use these algo's without fully  
understanding the implementation

# Multiclass Classification

Example classes

multiple classes

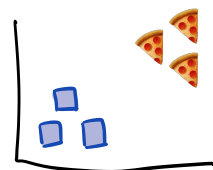
email tagging: work, friends, family  
y=1      y=2      y=3

\* indices can start at y=0 or y=1

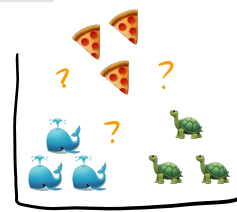
weather: sunny, cloudy, snowy



Instead of binary classification where we only have two classes (true or false, 🍕 or 🐳) ...

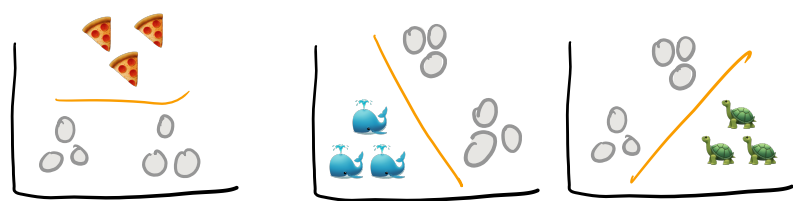


In multi-class classification we can have many classes (🍕, 🐳, 🐢)



We can solve this by creating multiple binary classification problems.

- 1. 🍕 and NOT 🍕
- 2. 🐳 and NOT 🐳
- 3. 🐢 and NOT 🐢



We now have 3 classifiers and for predictions, we run all 3 classifiers and pick the highest score

# Overfitting & Underfitting

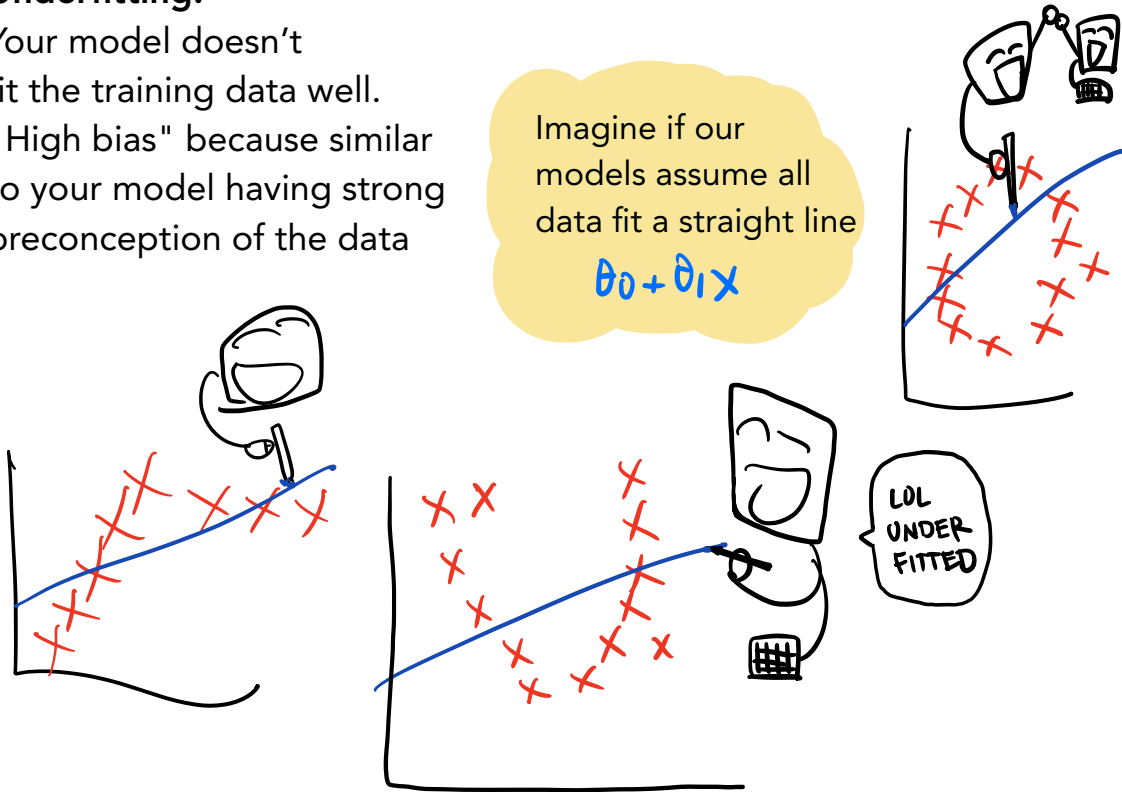
Your model may be improperly fit to the data

## ① Underfitting:

Your model doesn't fit the training data well. "High bias" because similar to your model having strong preconception of the data

Imagine if our models assume all data fit a straight line

$$\theta_0 + \theta_1 x$$



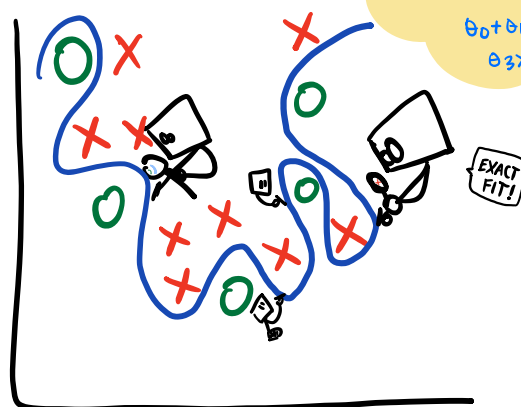
## ② Overfitting

Your model fits the training data too well. It tries too hard to fit the data and fails to generalize to new data

"high variance" ↗

Imagine if the model essentially memorized all the data points of the training set

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \dots$$

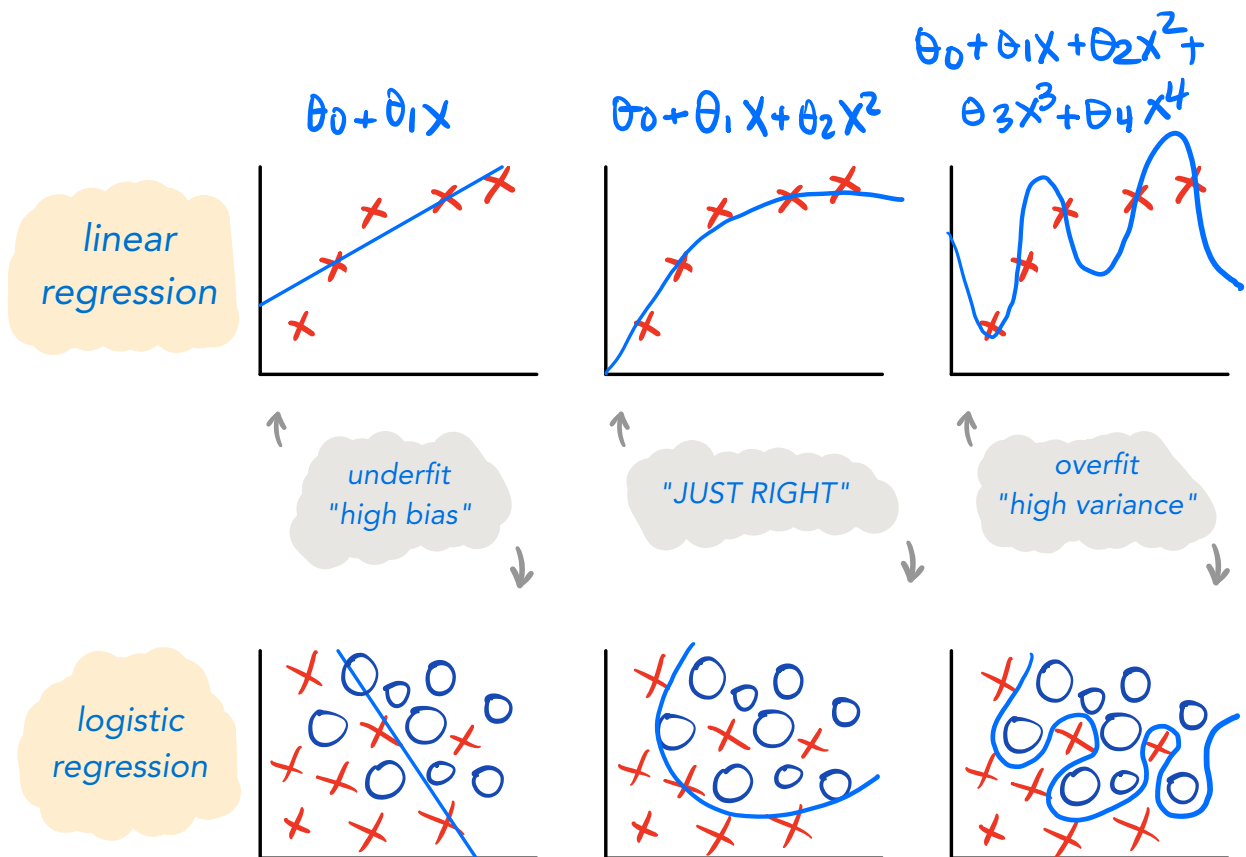



# OVERfitting

Q: What is over fitting?

A: If we have too many features, the learned hypothesis may fit the training data so well that it fails to "generalize" to new example inputs

Note: Given this data we can fit a linear, quadratic, even a higher order function



## Addressing Overfitting

Its easy to plot data with only a couple features  
but we will encounter data with many features

$X_1$	size of house
$X_2$	color of house
$X_3$	# of bedrooms
$\vdots$	
$X_{100}$	cardinal direction of 2nd bedrooms 3rd window

### Options

- ① Reduce number of features
  - Ⓐ manually remove features
  - Ⓑ choose a feature selection algo
- ② Regularization

### Q: What is regularization?

The idea behind regularization is that having smaller values for our  $\theta$  parameters creates a "simpler" hypothesis, smoother functions and is less prone to overfitting

Suppose we want to penalize and make  $\theta$  values really small. We do this in our cost function by adding an additional term that magnifies the affect to  $\theta$

$$J(\theta) = \text{NORMAL COST FUNCTION} + \text{REGULARIZATION TERM}$$

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

This lambda term ( $\lambda$ ) is a parameter we can adjust

this regularization term will reduce the  $\theta$  values

LINEAR REGRESSION  
W/REGULARIZATION:

$$\text{LINEAR REGRESSION} + \text{REGULARIZATION}$$

LOGISTIC REGRESSION  
W/REGULARIZATION:

$$\text{LOGISTIC REGRESSION} + \text{REGULARIZATION}$$



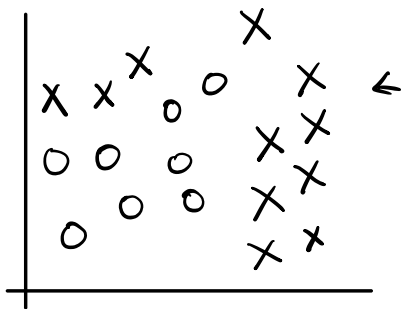
# Neural Networks

## Intuition

This week we learn about neural networks.

We already have linear regression and logistic regression, so why do we need another learning algorithm?

There are situations we want to learn **complex nonlinear hypothesis**. Consider you have this data.

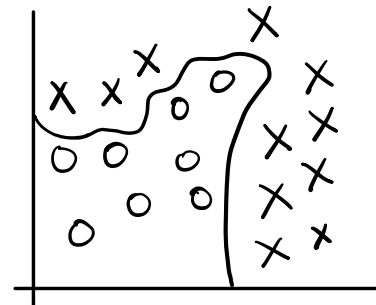


we could add non-linear terms

$$x_1 x_2 + x_1^2 x_2 + x_1 x_2^2 + \dots$$



When you have only two features we can afford to add all these terms... **but** often you have many more than two features and this becomes computationally expensive **so** we need a better learning algorithm



## How do we solve this?

If we used logistic regression, we would add non-linear terms that are complex enough to fit interesting datasets.

Quadratic such as  $x_1 x_2$ ,  $x_1 x_3$  etc...

Or cubic such as  $x_1 x_2 x_3$ ,  $x_1 x_2^2$

However, that is a lot of features which leads to

- overfitting,
- computationally expensive

Complex non-linear hypothesis are hard to learn when  $n$  is large



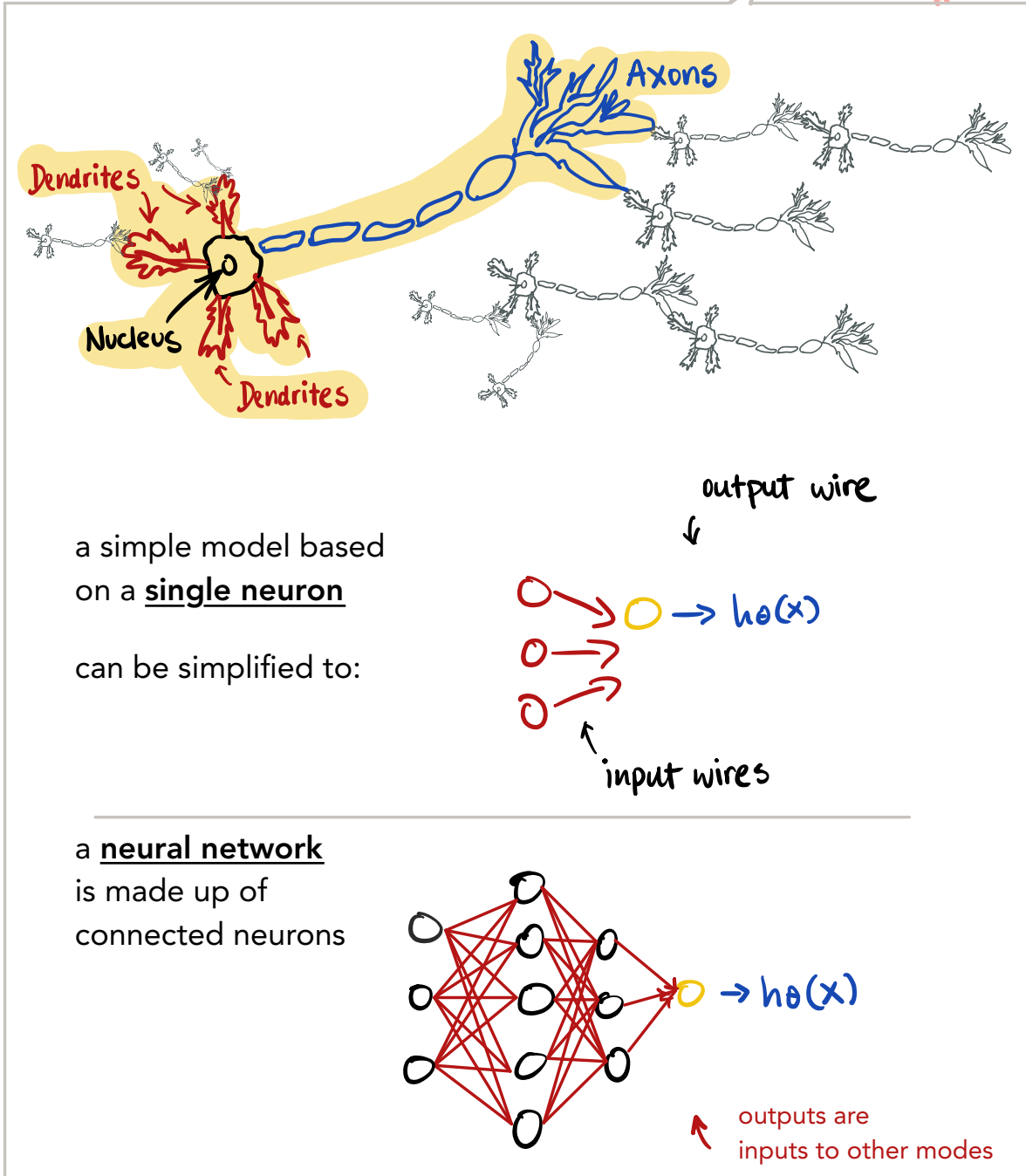
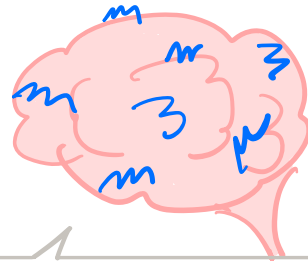
# Neurons in the brain

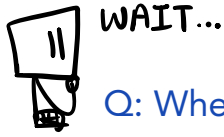


The neurons that make up our brain have

Dendrites - receivers of input

Axons - broadcasters of output





WAIT...

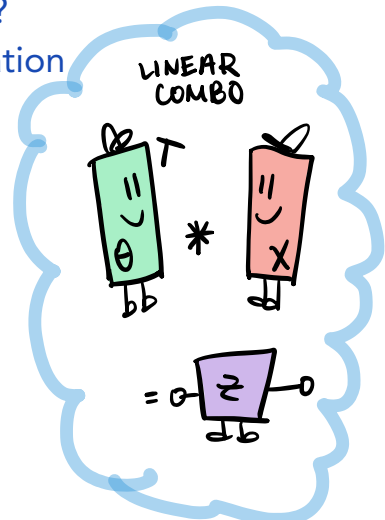
Q: Where does the nonlinearity come from?

A: At each node,  $\theta^T X$  is the linear combination of theta (weights) and x (inputs).

We call  $\theta^T X$ , "z"

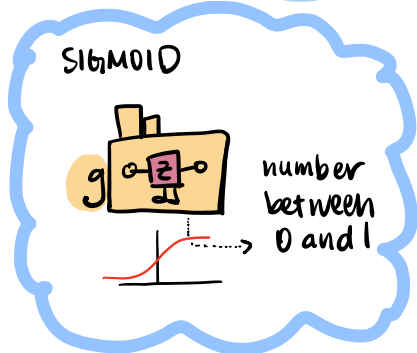
$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} \quad X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z = \theta^T X$$

$$= [\theta_0 \ \theta_1 \ \theta_2 \ \theta_3] * \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$



The difference is that we put z through an "activation function" which is a nonlinear function "g"

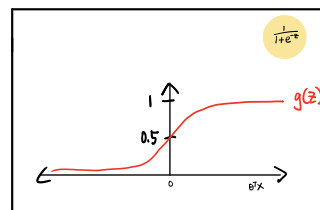
In this course, g will be the sigmoid function.



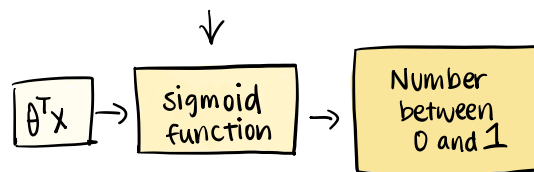
$$h_{\theta}(x) = g(z) \quad \text{sigmoid function!} \quad \frac{1}{1+e^{-z}}$$

A hand-drawn diagram showing three red arrows pointing towards a yellow circle. To the right of the circle is the text  $h_{\theta}(x)$ . Further right is the equation  $so \ h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$ .

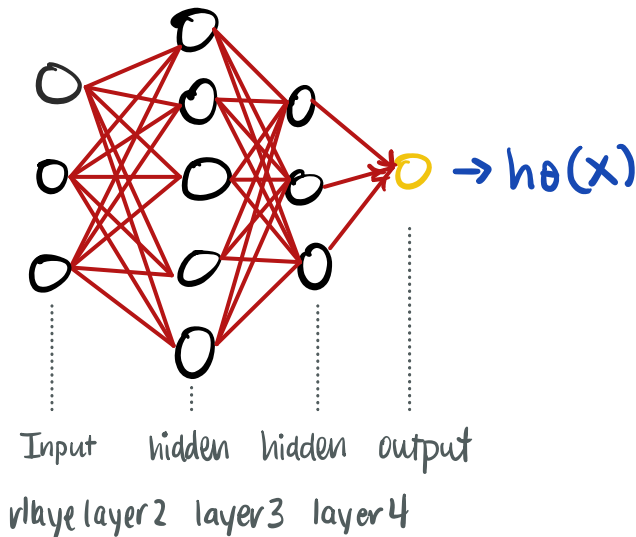
Here "g" is the same sigmoid function from our logistic regression lecture



\*In fact, without the network, this should look similar to logistic regression



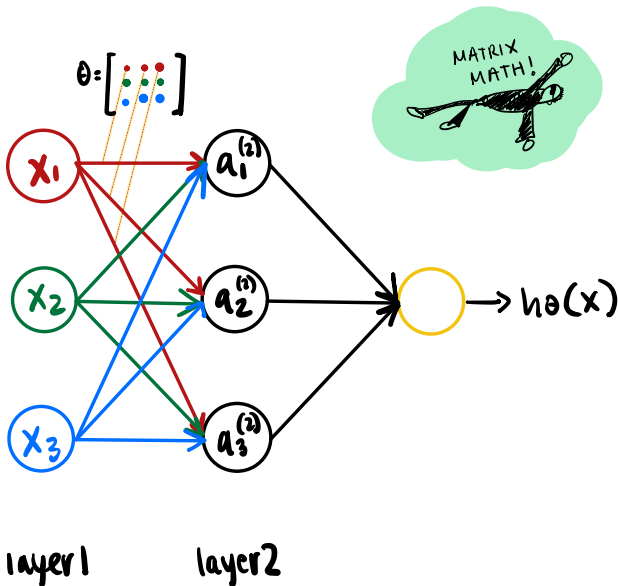
A neural network is a collection of these neuron calculations.



**Terminology**  
 The **first** layer is called the **input** layer  
 The **last** layer is the **output** layer  
 All layers in between are "**hidden**" layers

$$\theta = \begin{bmatrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix}$$

**Note:** Theta is now stored in a matrix  
 In the single neuron example (like in logistic regression), output to a single node means theta is a 1D vector.  
 In the network,, each layer can output to 1+ nodes for the next layer so theta is a matrix of paramters (or "weights")



**Notation**

- $a_i^{(j)}$  activation unit of unit  $i$  in layer  $j$
- $\theta^{(j)}$  matrix of weights going from layer  $j$  to layer  $j + 1$
- $x_n$  inputs from the data

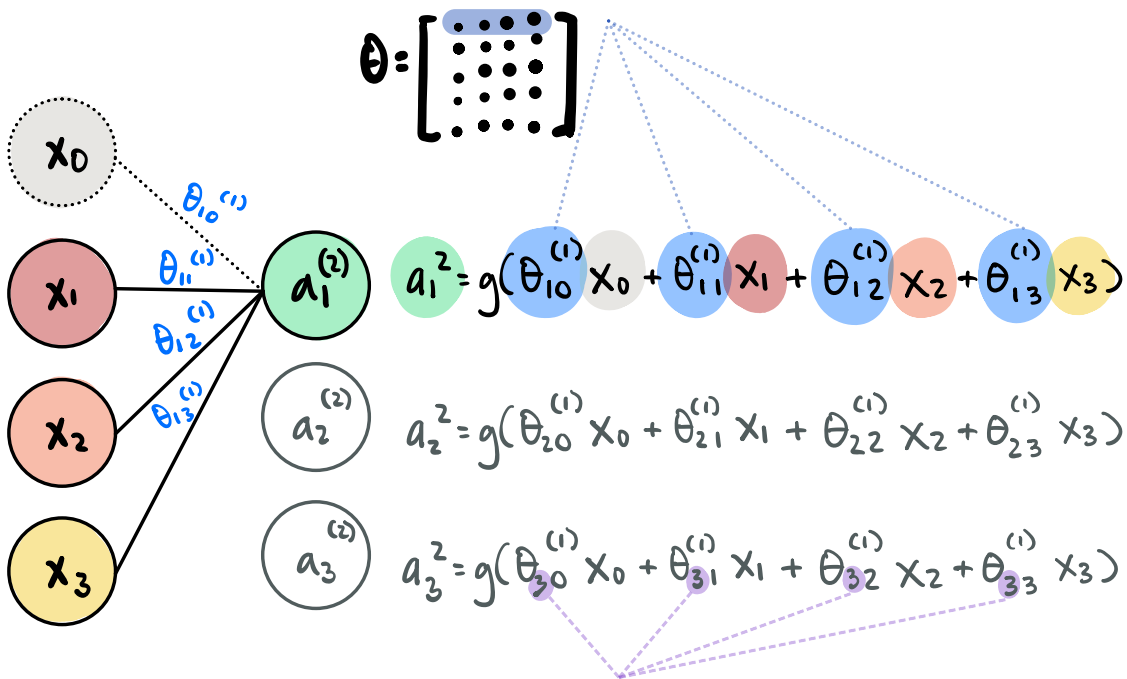
## Notation

This is a single  
theta value

$\theta$ 
  
 from layer
   
 to node
   
 from node

$\theta$ 
  
 $(1)$ 
  
 $30$ 
  
 for the 0th node of the layer (1)
   
 for the 3rd node of the next layer (2)

In matrix  $\theta$ , each row  
will be the theta weights  
for a single neuron to multiply



### Notice

That these values determine  
what row in the theta matrix  
is used to calculate this neuron

**Q: What are the dimensions of theta of a single layer?**

A: If a network has

- k units in layer j and
- y units in layer j+1,
- then  $\Theta(j)$  has dimensions  $(y * k+1)$ .
- Where 1 is added to k to adjust for the bias term

$$\Theta = \begin{bmatrix} ??? \\ ??? \\ ??? \end{bmatrix}$$

$$\Theta = \begin{bmatrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix}$$

**Example:**

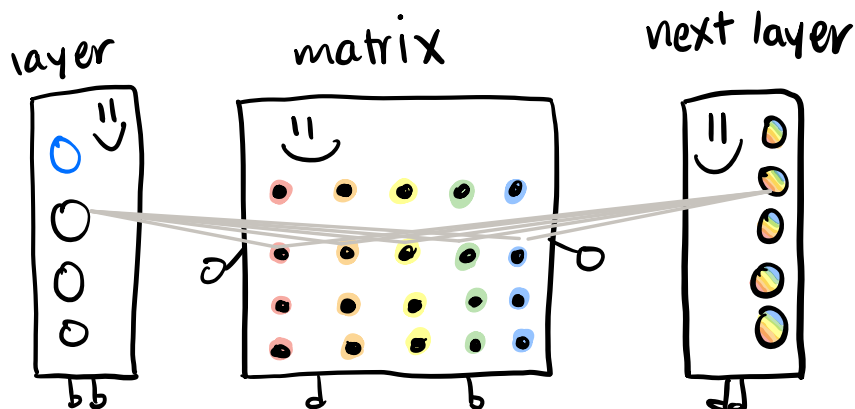
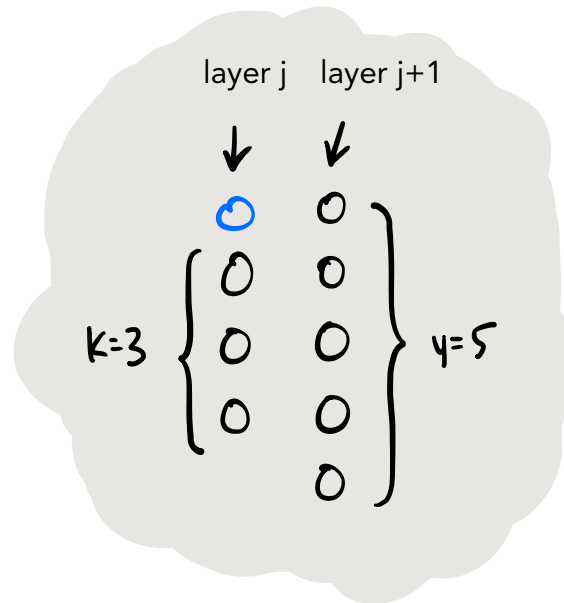
layer 1 has 3 units and  
layer 2 has 5 units

$\theta^{(1)}$  will have dimensions:

$$(y * k+1) = (5 * 4)$$

dimensions:

$$\left( \begin{array}{c} \text{units in} \\ \text{next layer} \end{array} \times \begin{array}{c} \text{units in current} \\ \text{layer plus bias unit} \end{array} \right)$$

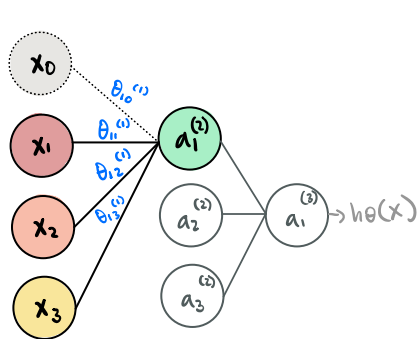


## Model Representation 2

Q: How do we **calculate** the layers efficiently?

A: Computers are extremely fast at vectorized and matrix multiplication so we solved with vector math.

"Forward Propagation": Vectorized Implementation

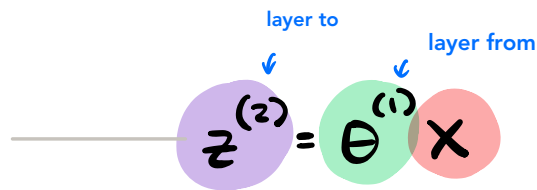
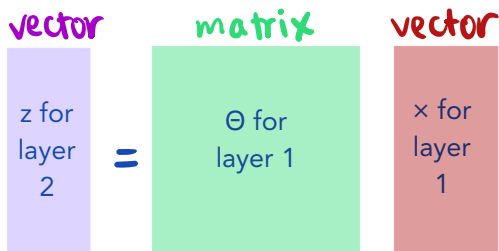


$$a_1^{(2)} = g(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3)$$

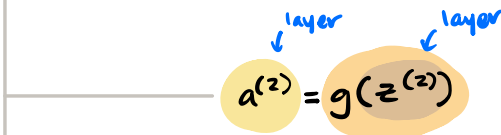
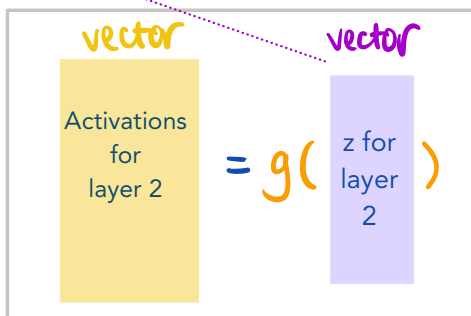
$$a_2^{(2)} = g(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3)$$

$$a_3^{(2)} = g(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3)$$

$$h_{\theta}(x) = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)})$$



Z is the linear combination of theta and X values



$a^{(2)}$  is the activation layer of layer 2 that is input to layer 3

otherwise known as "layer 1"  
or "activation layer 1"

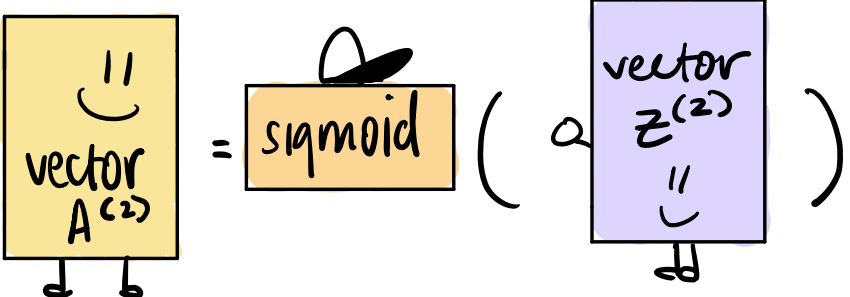
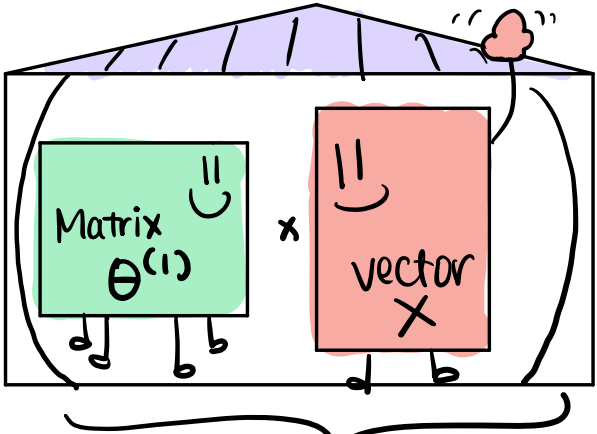
Say we have input vector  $x$   
Our first goal is to find activation layer 2.  
Here we represent these vectors in  
shorthand

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \quad a^{(2)} = \begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \end{bmatrix}$$

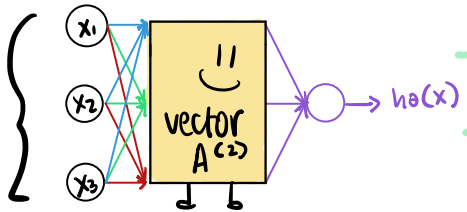
Solution:  
calculate the  
linear combination  
and take the sigmoid

$$z^{(2)} = \Theta^{(1)} x$$

$$a^{(2)} = g(z^{(2)})$$

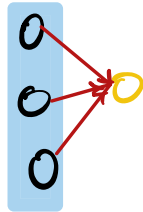


The input layer can  
also be called  $a^{(1)}$  or  
activation layer 1

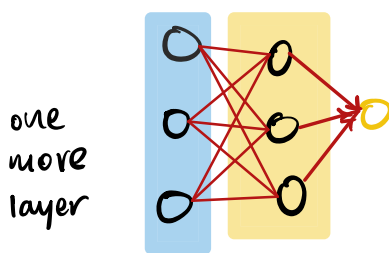


this process  
"feeds forward"  
into all further  
layers

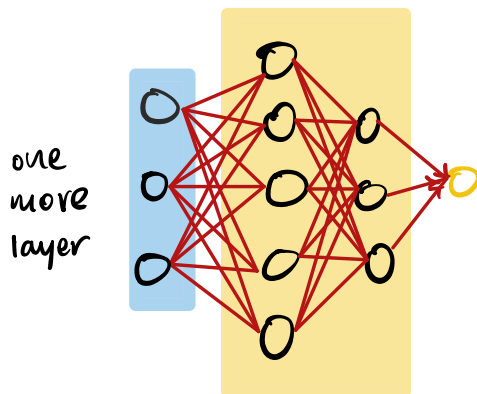
# # NN learning its own features



In logistic regression these are the features that contribute to output

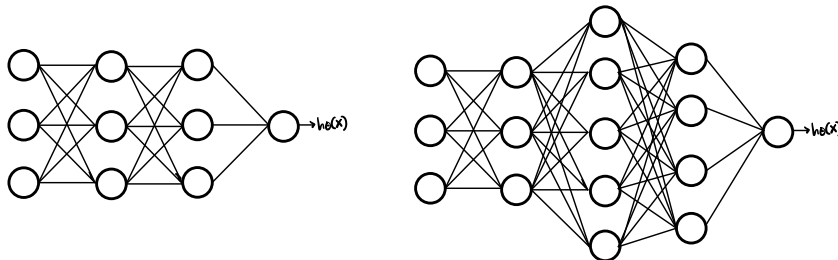


With neural networks layers are added that serve as inputs to the next layer



Its like NN are learning its own features layer by layer until the output

Note: we can choose different architectures with a different number of layers and nodes





# Neural Network: Examples and Intuitions

NN's allow us to model complex relationships that cannot be easily modeled as linear combinations

Two complex logical functions are

- ① XOR - exclusive OR
- ② XNOR - the inverse of XOR

## What is XOR (exclusive OR)?

XOR returns true (T)  
 if  $x_1$  OR  $x_2$  are true  
 but not if  
 $x_1$  AND  $x_2$  are true

$x_1$	$x_2$	$y$
0	0	0
1	0	1
0	1	1
1	1	0

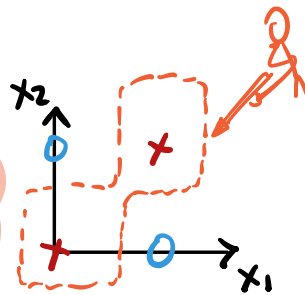
note the inverse

## What is XNOR?

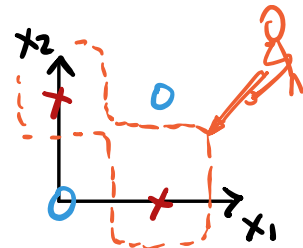
XNOR returns true  
 if  $x_1$  AND  $x_2$  are true  
 or if  $x_1$  AND  $x_2$  are false

$x_1$	$x_2$	$y$
0	0	1
1	0	0
0	1	0
1	1	1

Imagine the difficulty in drawing a decision boundary



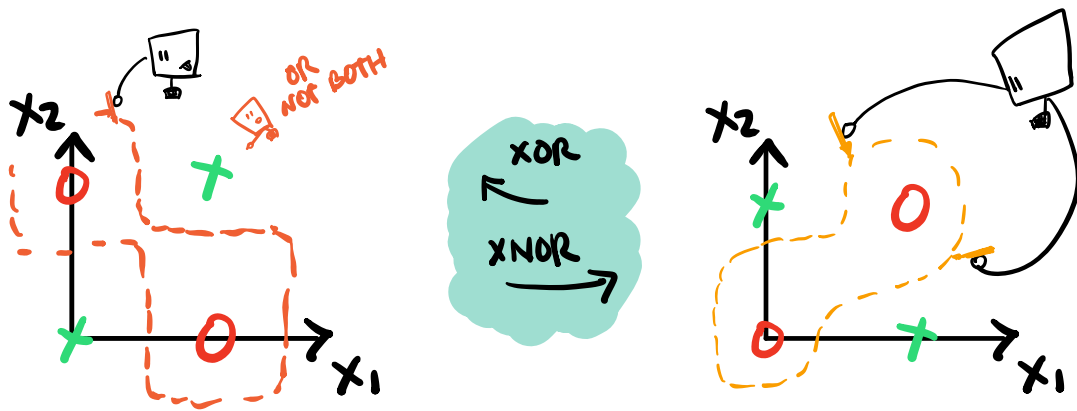
XNOR  
 XOR



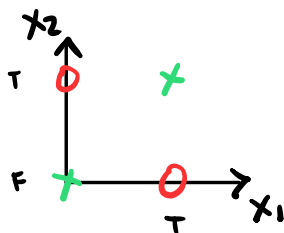
# Neural Network: Examples and Intuitions

NN's allow us to model complex relationships that cannot be easily modeled as linear combinations

Imagine how difficult it would be for a *linear function* to represent **XOR** (exclusive OR) or **XNOR** (inverse of XOR)



What is XOR (exclusive or)?

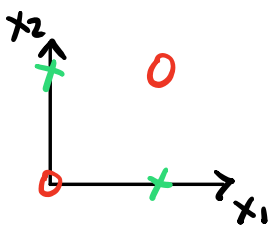


XOR returns true if  $x_1$  OR  $x_2$  are true but not if  $x_1$  AND  $x_2$  are true

$x_1$	$x_2$	$y$
0	0	0
1	0	1
0	1	1
1	1	0

note the inverse

What is XNOR?



XNOR returns true if  $x_1$  AND  $x_2$  are true or if  $x_1$  AND  $x_2$  are false

$x_1$	$x_2$	$y$
0	0	1
1	0	0
0	1	0
1	1	1

Q: How can we build XNOR?



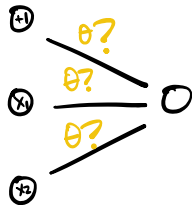
we can use simpler building blocks!



**AND**

Q: how do you build AND?

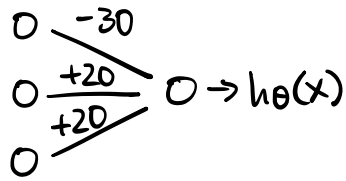
negative(-4.6) is .01



Q: What parameters ( $\theta$ ) would you need to create the AND function?

context: In our sigmoid function  
 positive(+4.6) is .99  $\leftarrow$  nearly 1  
 negative(-4.6) is 0.01  $\leftarrow$  nearly 0

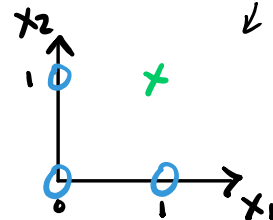
Say we used parameters: -30, +20, +20



$$= -30(1) + 20x_1 + 20x_2$$

Truth Table

$x_1$	$x_2$	$h(\theta)$
0	0	0
1	0	0
0	1	0
1	1	1



we created this chart!

AND here!  
I've made an AND here!

TO DO LIST

- AND
- OR
- NOR

## OR

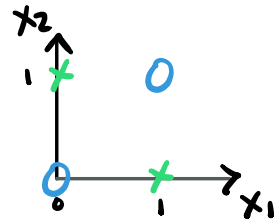
returns true  
 $x_1$  is true or  $x_2$  is true

$$\begin{array}{l} \text{O} \begin{array}{l} \swarrow -10 \\ \searrow +20 \end{array} \\ \text{O} \begin{array}{l} \swarrow +20 \\ \searrow +20 \end{array} \\ \text{O} \end{array} \rightarrow h_{\theta}(x)$$

$$= -10(\cdot) + 20x_1 + 20x_2$$

Truth Table

$x_1$	$x_2$	$y$
0	0	0
1	0	1
0	1	1
1	1	1



## NOT

returns true  
if  $x_1$  is false

$$\begin{array}{l} \text{O} \begin{array}{l} \swarrow +10 \\ \searrow -20 \end{array} \\ \text{O} \begin{array}{l} \swarrow -20 \\ \searrow -20 \end{array} \\ \text{O} \end{array} \rightarrow h_{\theta}(x)$$

$$= 10(\cdot) - 20x_1$$

Truth Table

$x$	$y$
0	1
1	0



## NOR

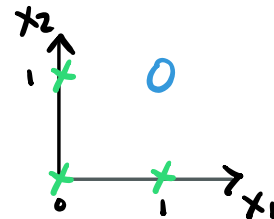
returns true  
 $x_1$  is false OR  $x_2$  is false

$$\begin{array}{l} \text{O} \begin{array}{l} \swarrow +10 \\ \searrow -20 \end{array} \\ \text{O} \begin{array}{l} \swarrow -20 \\ \searrow -20 \end{array} \\ \text{O} \end{array} \rightarrow h_{\theta}(x)$$

$$= 10(\cdot) - 20x_1 - 20x_2$$

Truth Table

$x_1$	$x_2$	$y$
0	0	1
1	0	1
0	1	1
1	1	0



XNOR made of

- ① AND:  $x_1$  and  $x_2$
- ② NOR: NOT  $x_1$  and  $x_2$
- ③ OR:  $x_1$  OR  $x_2$

**AND**

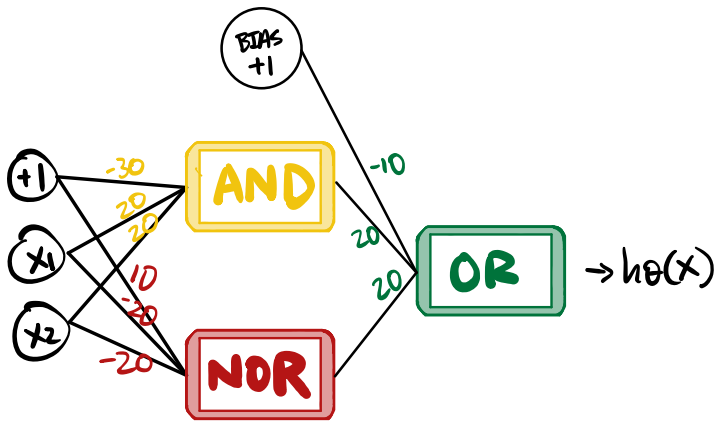
**NOR**

**OR**

$$\theta^{(1)} = [-30, 20, 20]$$

$$\theta^{(1)} = [10, -20, -20]$$

$$\theta^{(1)} = [-10, 20, 20]$$



"Building" XNOR with AND, OR and NOT

$$\text{XNOR} = (x_1 \text{ OR } x_2) \text{ AND } (\text{NOT } (x_1 \text{ AND } x_2))$$

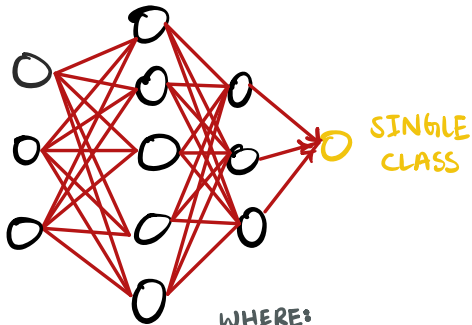
NOR

NN's can output to many nodes to represent multiple classes

SINGLE CLASS

vs

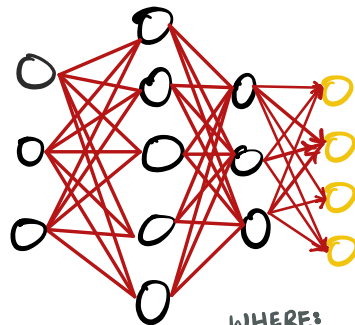
MULTIPLE CLASSES



WHERE:

[1] POSITIVE

[0] NEGATIVE



WHERE:

$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$  class 1

$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$  class 3

$\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$  class 2

$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$  class 4

Multiple classes can be represented by one dimensional vectors where all values are 0 except for a single 1 value to represent the class

$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$  class 1  $\Rightarrow$  🐢

$\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$  class 2  $\Rightarrow$  🐳

$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$  class 3  $\Rightarrow$  🍕

$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$  class 4  $\Rightarrow$  🐶

week 5

week 5

week 5

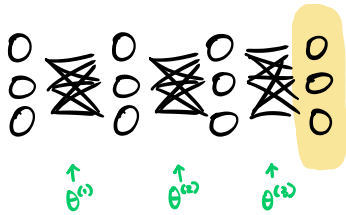
week 5

week 5

wee  
w

## Neural Network Cost Function

Just like linear & logistic regression, we need a cost function, the derivative of which will allow us to fit parameters that will minimize cost

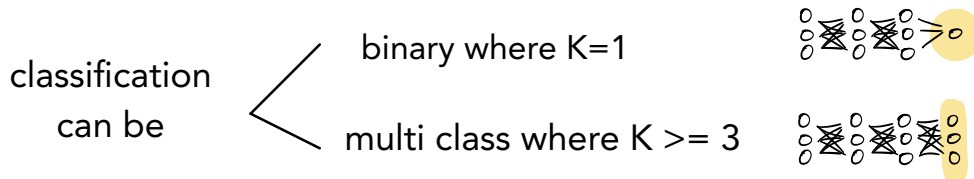


### Terminology

$L$  : total no. of layers in network

$S_L$  : no. of units in layer  $L$

$K$  : no. of output units



The NN cost function is a generalized version of the logistic regression cost function adjusted for multiple output nodes and theta dimensions of a matrix



cost function = cost of predictions + regularization

### Logistic regression cost function

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

with NN's there are  $K$  output nodes so we have to sum over every  $K$  output

regularization now accounts for layers and multi dimensions of theta matrix

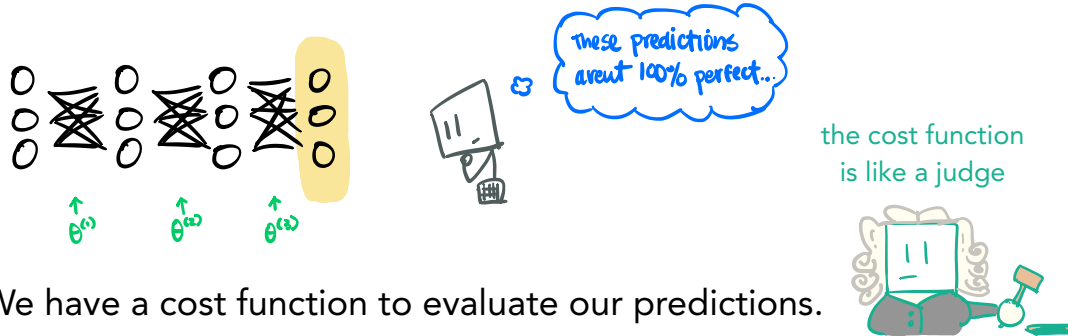
### Neural network cost function

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log h\theta(x^{(i)})_k + (1 - y_k^{(i)}) \log(1 - (h\theta(x^{(i)}))_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{ji}^{(l)})^2$$



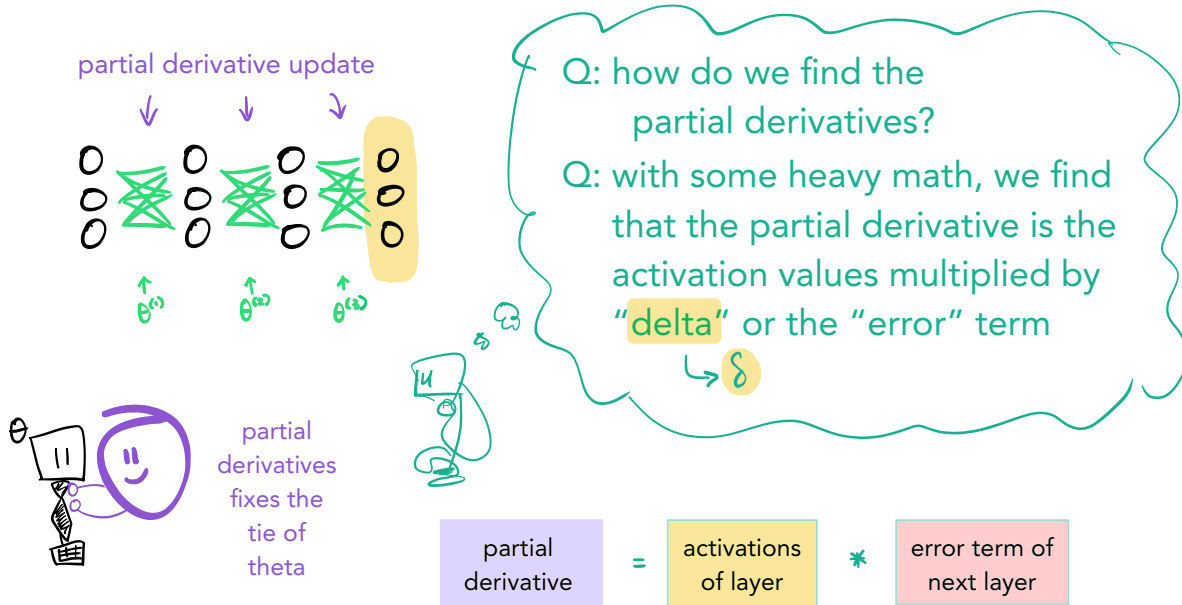
# Introducing backpropagation and why

Through the feed forward mechanism our NN creates an output prediction layer



We have a cost function to evaluate our predictions.

**Backpropagation:** we calculate partial derivatives so we can nudge our theta (parameters/weights) by tiny amounts to minimize our cost "J(θ)"



partial derivative for a single training example

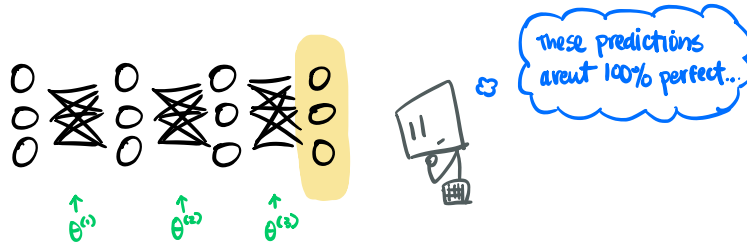
$$\frac{\partial J(\theta)}{\partial \theta^{(n)}} = a^{(n)} \delta^{(n+1)}$$

where  $\delta$  is "delta" or the "error"  
 $a$  is the activation values  
 $n$  is the layer

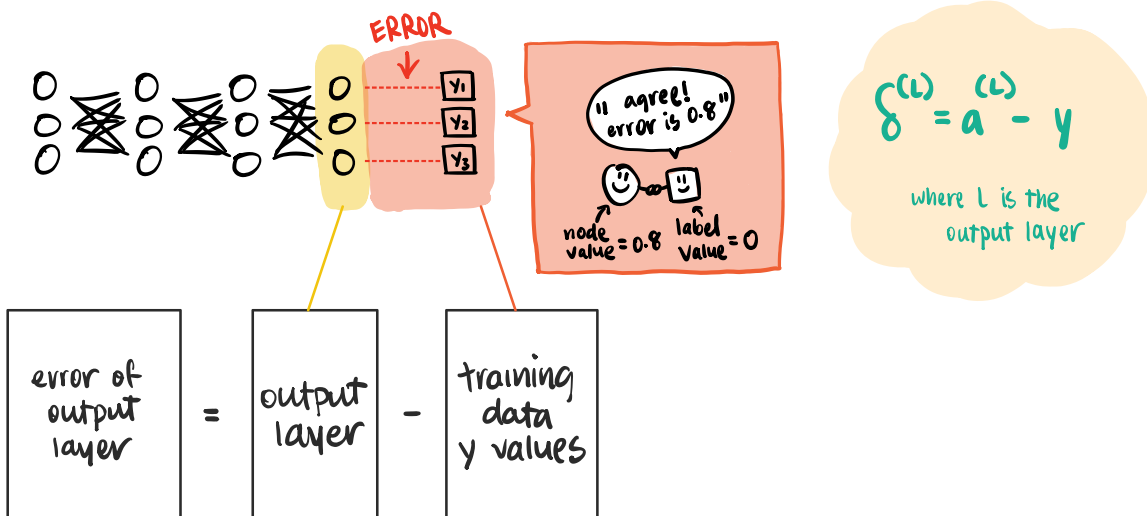
\*we dont create the proof of this partial derivative in the lectures

# Backpropation Algorithm

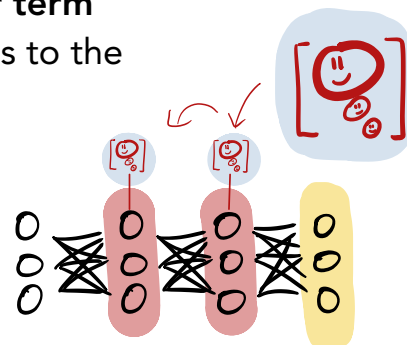
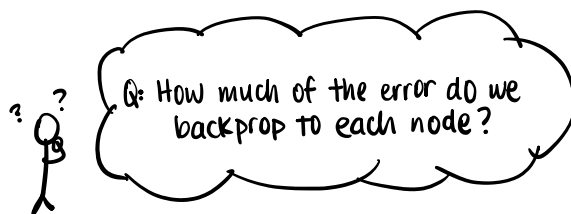
Through the feed forward mechanism our NN creates an output prediction layer



This output can be directly compared to the actual label values from the training set. The output layer **error terms** are straightforward



Q: How can we use the output layer **error term** ("delta") to propagate back error terms to the other layers?



## Calculating delta for layers

output layer delta = output layer - y labels

$$\delta^{(k)} = a^{(k)} - y$$

hidden layer delta is a function of  $\theta$ , delta of the next layer and derivative of sigmoid(z) and values

$$\delta^{(n)} = (\theta^{(n)})^T \delta^{(n+1)} .* g'(z^{(n)})$$

$$\text{where } g'(z^{(n)}) = a^{(n)} .* (1 - a^{(n)})$$

fully expanded:

$$\delta^{(n)} = (\theta^{(n)})^T \delta^{(n+1)} .* (a^{(n)} .* (1 - a^{(n)}))$$

These delta values are used in partial derivative to determine changes in theta "parameters"

single example  $\rightarrow$

$$\frac{\partial J(\theta)}{\partial \theta_{ij}^{(n)}} = a_j^{(n)} \delta_i^{(n+1)}$$

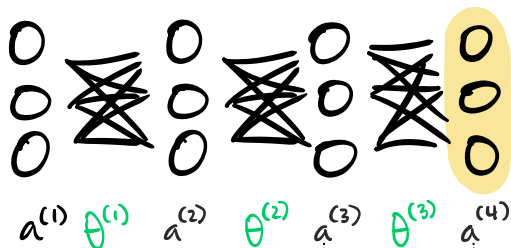
For our training set, the partial derivatives are averaged where  $m$  is the total number of training examples  $t$

multiple example  $\rightarrow$

$$\frac{\partial J(\theta)}{\partial \theta_{ij}^{(n)}} = \frac{1}{m} \sum_{t=1}^m a_j^{(t)(n)} \delta_i^{(t)(n+1)}$$

\* note: for simplicity we leave out the regularization term

Forward propagation calculates activation nodes and the output layer



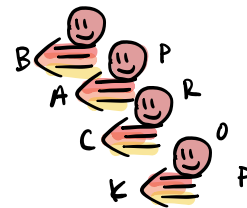
the output layer calculates the delta error term and backpropagate error terms for every layer

$$\delta^{(4)} = a^{(4)} - y$$

$$\delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} \cdot * g'(z^{(3)})$$

$$\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} \cdot * g'(z^{(2)})$$

$\nearrow$  theta used to compute next layer  
 $\uparrow$  delta of next layer  
 $\uparrow$  dot product  
 $\nwarrow$  sigmoid derivative =  $a^{(n)} \cdot * (1 - a^{(n)})$



Input layer error is not calculated



error terms are used to calculate the partial derivatives

$$\frac{\partial J(\theta)}{\partial \theta_{ij}^{(3)}} = a_j^{(3)} \delta_i^{(4)}$$

$$\frac{\partial J(\theta)}{\partial \theta_{ij}^{(2)}} = a_j^{(2)} \delta_i^{(3)}$$

$$\frac{\partial J(\theta)}{\partial \theta_{ij}^{(1)}} = a_j^{(1)} \delta_i^{(2)}$$

this multiplication creates a matrix that has dimensions for the matrix  $\theta^{(3)}$

\* this notation is for one training example. The actual training set will average over all training examples

complete notation	$\frac{\partial J(\theta)}{\partial \theta_{i,j}^{(l)}} = \frac{1}{m} \sum_{t=1}^m a_j^{(t)(l)} \delta_i^{(t)(l+1)}$
-------------------	--

"If this is difficult,  
**you are not alone**"

"I've actually used back propagation pretty successfully for many years and even today I still don't, sometimes, feel like I have a very good sense of just what it's doing or sort of intuition about what back propagation is doing"

- Andrew Ng

