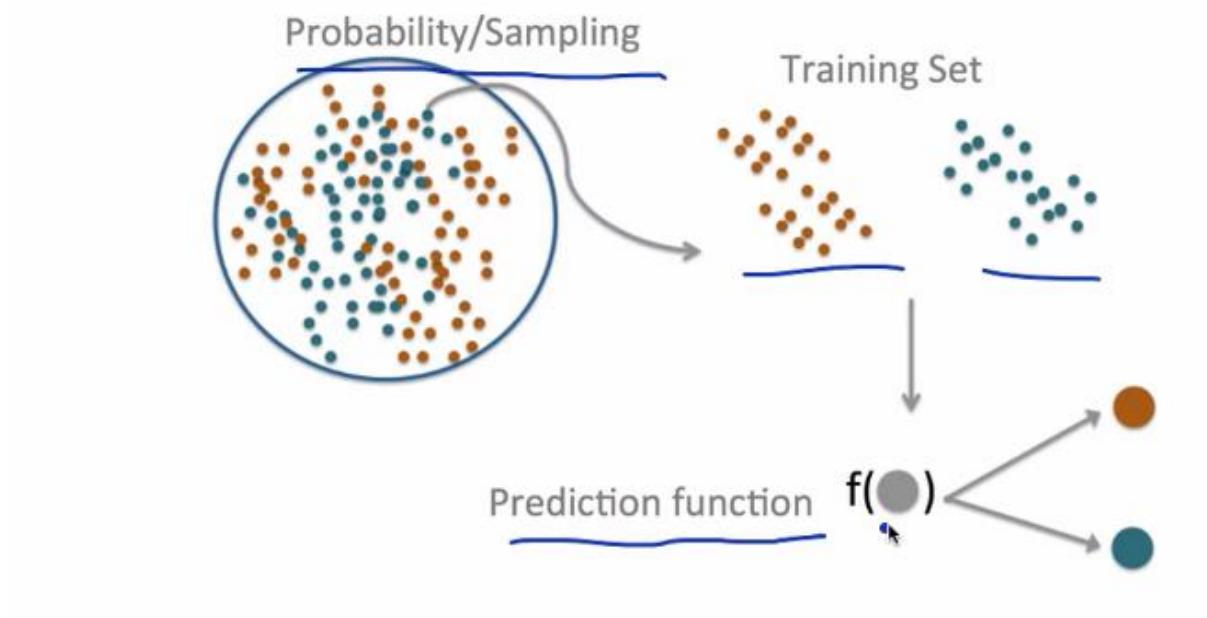


PREDICTION:

The central dogma of prediction



Suppose we have a dataset of red and blue balls. Now we take a sample of these red and blue balls each and train the sample and create a prediction function. Now, what this prediction function does is that it classifies a random ball whose color is not known into whether it is red or blue .

Now, we create a very basic prediction model , we want to predict spam or non – spam emails on the basis of frequency of the word ‘your’ in an email:

If frequency is > 0.5 , email is a spam and not otherwise.

```
prediction<-ifelse(spam$your>0.5,"spam","nonspam")
```

Now , we compare our prediction model to the actual prediction made in the data set:

```
xtabs(~prediction+spam$type)
      spam$type
prediction nonspam spam
  nonspam     2112   468
    spam       676 1345
```

We see that out of our total predicted non spam emails i.e. 2580 , 2112 are in actual non-spam and 468 are spam. Similarly for spam emails , out of the total 2021 emails , 676 are in actual non-spam and 1345 are spam.

Hence our prediction is correct $2112/2580$ i.e. 81.8% correct for non spam email and $1345/2021$ i.e. 66.5% correct for spam emails.

INSAMPLE AND OUTSAMPLE ERRORS:

1.In sample error is in the same data which we have used to train our predictor , also called resubstitution error.

In Sample error is always going to be a little bit optimistic , from what the error is that you would get from a new sample. And the reason why is, in your specific sample, sometimes your prediction algorithm will tune itself a little bit to the noise that you collected in that particular data set. And so when you get a new data set, there'll be different noise, and so the accuracy will go down a little bit.

Hence , In sample error< out of sample error always.

2.Out of sample error : The error we get when we perform our prediction on some other sample data set . We usually want to test our prediction using another sample data to see how it would actually work on the real dataset.

OVERFITTING: When we train our model ‘too perfectly’ i.e. it performs very accurately on the sample chosen to train the predictor, there can be some errors when it is predicting for some our sample or for the population itself , since it is built in a way to tackle the noise there in the sample used to train it.

Overfitting

- Data have two parts
 - Signal
 - Noise
- The goal of a predictor is to find signal
- You can always design a perfect in-sample predictor
- You capture both signal + noise when you do that
- Predictor won't perform as well on new samples

<http://en.wikipedia.org/wiki/Overfitting>

Rules of thumb for prediction study design

- If you have a large sample size
 - 60% training
 - 20% test
 - 20% validation
- If you have a medium sample size
 - 60% training
 - 40% testing
- If you have a small sample size
 - Do cross validation
 - Report caveat of small sample size

Some principles to remember

- Set the test/validation set aside and *don't look at it*
- In general randomly sample training and test
- Your data sets must reflect structure of the problem
 - If predictions evolve with time split train/test in time chunks (called [backtesting](#) in finance)
- All subsets should reflect as much diversity as possible
 - Random assignment does this
 - You can also try to balance by features - but this is tricky

TYPES OF ERROR:

Basic Terms:

In general, **Positive** = identified and **negative** = rejected. Therefore:

True positive = correctly identified

False positive = incorrectly identified

True negative = correctly rejected

False negative = incorrectly rejected

Explaining the above with an example:

Medical testing example:

True positive = Sick people correctly diagnosed as sick

False positive = Healthy people incorrectly identified as sick

True negative = Healthy people correctly identified as healthy

False negative = Sick people incorrectly identified as healthy.

Key Quantities:

Key quantities

		DISEASE	
		+	-
TEST	+	TP	FP
	-	FN	TN

Sensitivity

→ $\Pr(\text{positive test} \mid \text{disease})$

Specificity

→ $\Pr(\text{negative test} \mid \text{no disease})$

Positive Predictive Value

→ $\Pr(\text{disease} \mid \text{positive test})$

Negative Predictive Value

→ $\Pr(\text{no disease} \mid \text{negative test})$

Accuracy

→ $\Pr(\text{correct outcome})$

1. Sensitivity means diseased if a person is diseased , what is the probability that we get it right from our model?

2. Specificity means if a person is healthy , what is the probability we get it right from our model.

3. Positive predictive value means if our test shows that you are diseased, what is the probability that you are actually diseased or if we look at all the people we call diseased , what fraction of them are actually diseased.

4. Negative predictive value means that if our test/model shows that you are healthy, what is the probability that you are actually healthy.

You were healthy but we predicted you to be sick.

You're sick but we predicted you to be healthy.

5. Accuracy is the probability of the correct outcomes, in the table below , it is the true positives (TP) and the true negatives(TN) added up,

		DISEASE	
		+	-
TEST	+	TP	FP
	-	FN	TN

NOTE: In the Above table:

1.Say our prediction model predicted TP + FP people as diseased, out of that , from the lab results , we saw that TP were actually diseased whereas FP where healthy so TP was the correct prediction whereas FP was the wrong prediction.

2.Similarly Say our prediction model predicted FN + TN people as Healthy, out of that , from the lab results , we saw that TN were actually healthy whereas FN where sick so TN was the correct prediction whereas FN was the wrong prediction.

Key Quantities as Fractions :

→ people actually diseased (+) or healthy (-)

		DISEASE	
		+	-
TEST	+	TP	FP
	-	FN	TN

→ people our model shows as diseased (+) or healthy (-)

Sensitivity → $TP / (TP+FN)$

Specificity → $TN / (FP+TN)$

Positive Predictive Value → $TP / (TP+FP)$

Negative Predictive Value → $TN / (FN+TN)$

Accuracy → $(TP+TN) / (TP+FP+FN+TN)$

Question:

Screening tests

Assume that some disease has a 0.1% prevalence in the population. Assume we have a test kit for that disease that works with 99% sensitivity and 99% specificity. What is the probability of a person having the disease **given the test result is positive**, if we randomly select a subject from

- ▶ the general population?
- ▶ a high risk sub-population with 10% disease prevalence?

Ans: Let us assume that our population is 100000

① Now, as given in the question, only 0.1% of these 100000 have the disease, i.e. only 100 have disease.

② Thus 99900 are healthy.

③ sensitivity = 94.1%, specificity = 99.1%.
hence out of the 100 diseased people we predict $99.1 \times 100 = 99$ correctly.

e) And, out of the 99900 healthy people, we predict 99% i.e. $\frac{99}{100} \times 99900$. To be healthy i.e. 98901.

		DISEASE	
		+	-
TEST	+	99	999
	-	1	98901

		DISEASE	
		+	-
TEST	+	99	999
	-	1	98901

Sensitivity

$$\rightarrow 99 / (99+1) = 99\%$$

Specificity

$$\rightarrow 98901 / (999+98901) = 99\%$$

Positive Predictive Value

$$\rightarrow 99 / (99+999) \approx 9\%$$

Negative Predictive Value

$$\rightarrow 98901 / (1+98901) > 99.9\%$$

Accuracy

$$\rightarrow (99+98901) / 100000 = 99\%$$

In above picture we see that positive predicted value is very less, it is only 9.1%, hence there is only a 9.1% chance that a person actually has a disease, given that we predict him/her to be sick.

This can be attributed to the fact that only 0.1% of the people in our population have a disease, so the probability of a person having a disease decreases considerably.

If instead we assume that 10% of our population has a particular disease:

$$\begin{aligned}\text{No. of diseased people} &= 10\% \text{ of } 100000 \\ &= 10000\end{aligned}$$

$$\begin{aligned}\text{No. of healthy people} &= 100000 - 10000 \\ &= 90000\end{aligned}$$

		DISEASE	
		+	-
TEST	+		
	-	100	89100
		9900	900

positive Predicted value = $\frac{9900}{9900 + 900}$

$$= \frac{9900}{10800} = 91.6\%$$

hence, 91.6% of our predicted diseased patients are actually diseased.

		DISEASE	
		+	-
TEST	+		
	-	100	89100
		9900	900

- Sensitivity → $9900 / (9900+100) = 99\%$
- Specificity → $89100 / (900+89100) = 99\%$
- Positive Predictive Value → $9900 / (9900+900) \approx 92\%$
- Negative Predictive Value → $89100 / (100+89100) \approx 99.9\%$
- Accuracy → $(9900+89100) / 100000 = 99\%$

ERROR IN PREDICTED VS TRUE VALUE:

For continuous data

① Mean squared error (MSE):

$$\frac{1}{n} \sum_{i=1}^n (\underline{\text{Prediction}_i} - \underline{\text{Truth}_i})^2$$

② Root mean squared error (RMSE):

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (\text{Prediction}_i - \text{Truth}_i)^2}$$

Formula 1 given above is just like calculating average residuals in Regression.

In formula 2, we take the sq. root to make the units back to what they were.

Receiver Operating Characteristic curve OR The

ROC CURVE

In ROC curve, we plot $P(TP)$ vs $P(FP)$. $P(TP)$ tells us out of the actual total diseased people, how many we correctly predicted to be diseased.

Refer to this table:

		DISEASE	
		+	
TEST	+	TP	FP
	-	FN	TN

$P(TP) = \frac{\text{No. of people correctly predicted to be diseased}}{\text{Total no. of people actually diseased}}$

$$= \frac{TP}{TP+FN}$$

$P(FP)$ tells us out of all people who are actually healthy, how many did we predict to be sick.

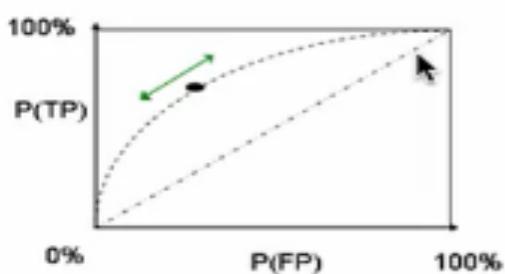
$P(FP) = \frac{\text{No. of people we wrongly predicted to be sick}}{\text{Total healthy people}}$

$$= \frac{FP}{FP+TN}$$

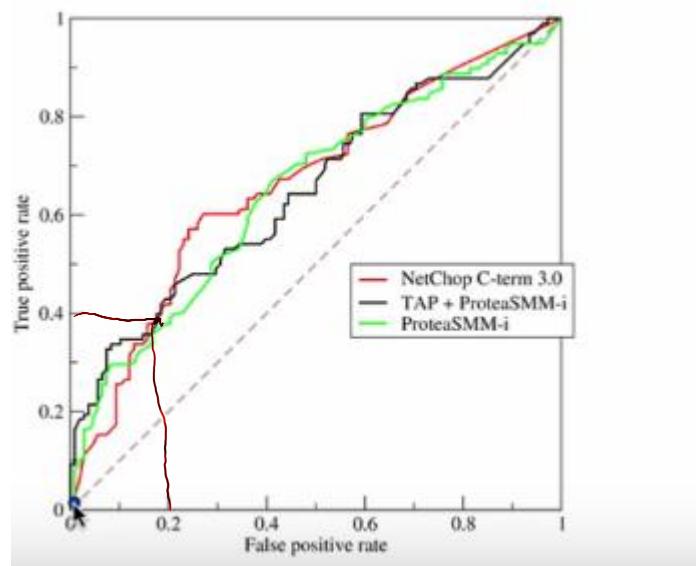
Note: $P(TP) = \text{sensitivity}$, where $P(FP) = 1 - \text{specificity}$

$$= 1 - \frac{TN}{FP+TN}$$

ROC CURVE:



An actual eg of ROC curve:



Here in the above curve , a false positive rate of 0.2 means a specificity of 0.8 whereas the point 0.4 => a sensitivity of 0.4 . Hence, this point has a high specificity but a low sensitivity.

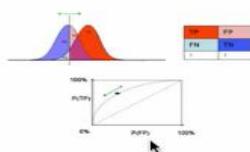
COMPARING TWO OR MORE ROC's:

What we want in our ROC is that it should have low False positive rate and high true positive rate both of which indicate high sensitivity and high specificity. Thus, we want our points to be on the upper left side of the curve. The point on the upper left corner has the highest specificity and the highest sensitivity.

We compare two or more ROC's by looking at the area under the curve. High area indicates that the curve will be more towards the upper left side and hence the points lying on it will have high sensitivity and high specificity.

GENEREALLY:

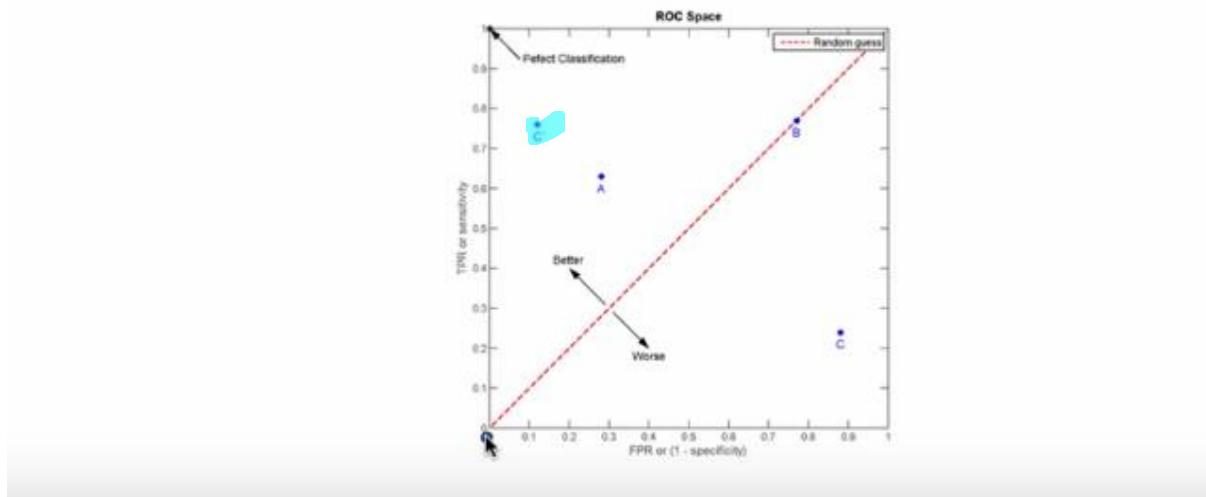
Area under the curve



- AUC = 0.5: random guessing
- AUC = 1: perfect classifier
- In general AUC of above 0.8 considered "good"

AUC = 0.5 is the dotted line and it indicates that is prediction is just random guessing, whereas AUC <0.5 indicates our prediction is worse than random guessing and the points on such a curve will be towards the lower right corner.

What is good?



In the above picture, the highlighted point C is the best predictor followed by point A . Point B is just a random guessing predictor whereas point C below the dotted line is the worst predictor.

CROSS VALIDATION:

We subset our sample data which we are going to use to train the prediction model into a training set and a testing set. We do this because we can't use another sample data as a test set while building the model since, this new sample dataset can then become a part of the training set and model can become acquainted with the noise there in this new dataset .

Hence we do all our testing while building the model on the training data set only by sub-setting it.

Cross-validation

Approach:

1. Use the training set
2. Split it into training/test sets
3. Build a model on the training set
4. Evaluate on the test set
5. Repeat and average the estimated errors

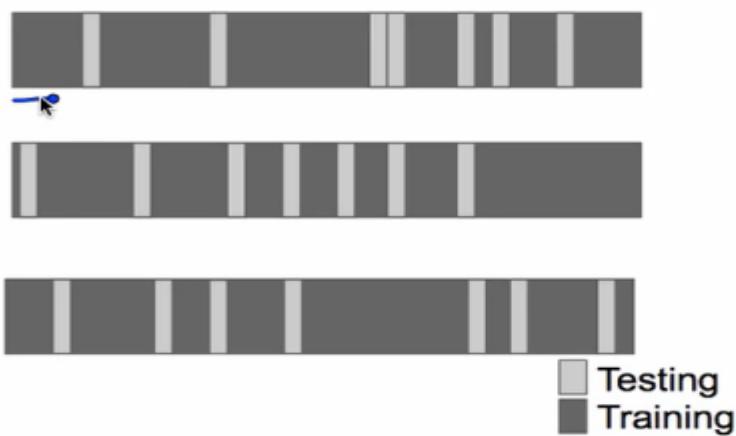
Used for:

1. Picking variables to include in a model
2. Picking the type of prediction function to use
3. Picking the parameters in the prediction function
4. Comparing different predictors

4/8

METHODS OF SUBSETTING:

1. RANDOM SUBSAMPLING:



#We randomly divide the data into testing and training datasets.

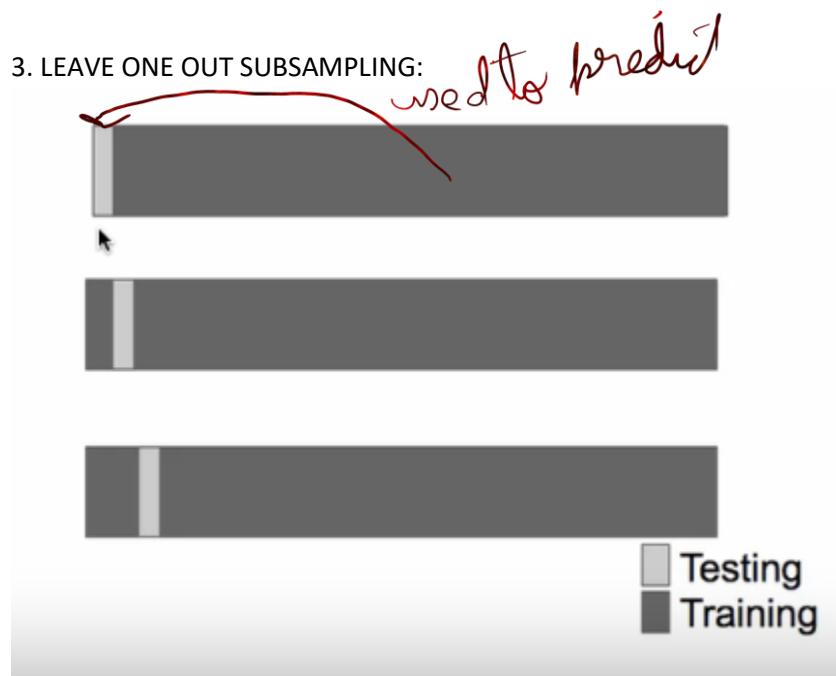
#We use the black coloured values to train our model and try to predict the grey coloured values to test the accuracy of our model.

2. K-FOLD SUBSAMPLING:



Above is an example of 3 fold testing. We divide our data into 1/3 testing , 2/3 training then into 1/3 training ,1/3 testing ,1/3 training and then into 2/3 training , 1/3 testing and do the training of black parts and testing on grey parts.

3. LEAVE ONE OUT SUBSAMPLING:



- a) We leave one sample out each time for testing and do the training using the other samples.
- b) In the first picture, we leave the first sample out for testing and use other for training.
- c) In the second picture, we leave the second sample out for testing and use other samples for training.
- d) We do this until we reach the end of the dataset.

SOME THINGS ABOUT CROSS VALIDATION TO KEEP IN MIND

Considerations

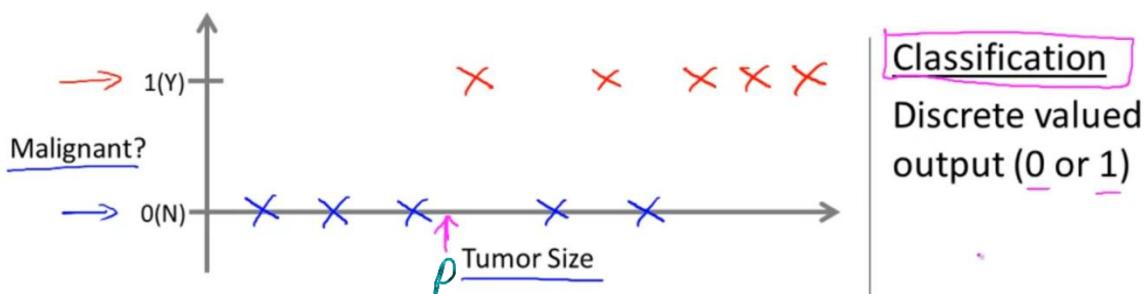
- For time series data data must be used in "chunks"
- For k-fold cross validation
 - Larger k = less bias, more variance
 - Smaller k = more bias, less variance
- Random sampling must be done *without replacement*
- Random sampling with replacement is the bootstrap
 - Underestimates of the error
 - Can be corrected, but it is complicated ([0.632 Bootstrap](#))
- If you cross-validate to pick predictors estimate you must estimate errors on independent data.

Supervised learning:

1. In supervised learning we give the programme certain amount of information and it trains itself using and tries to predict further things.

Eg:

Breast cancer (malignant, benign)



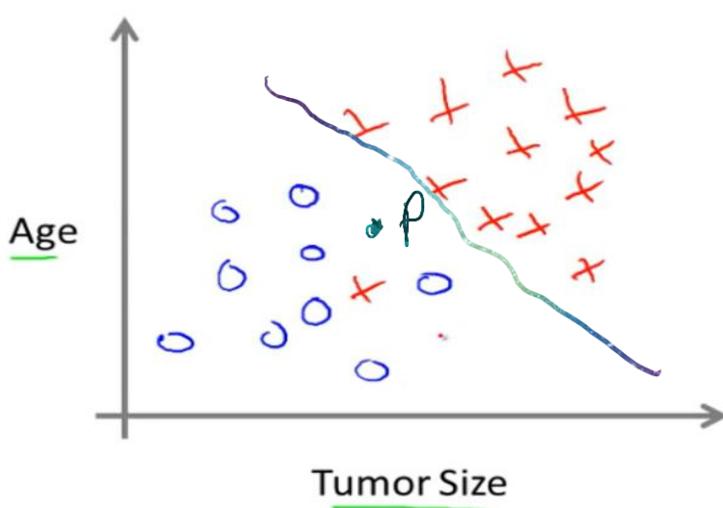
It turns out that in classification problems,

Andrew Ng

Here the x axis corresponds to certain tumour sizes and y axis tells that whether a tumour is malignant(1 on the y line)i.e. harmful or benign(0 on the y line) i.e. harmless.

We have inputted some pre-existing values pertaining to tumours -identified in the graph by blue and red crosses into the programme. The programme now tries to classify the point p as malignant or benign on the basis of the information we have inputted.

Eg2:



Here we try to find the relation between age and tumor size. The blue dots represent benign tumors and the red crosses represent malignant tumors.

If we ask the programme to classify point p as malignant or benign , it would first try and separate crosses and dots from each other and then try to classify p.

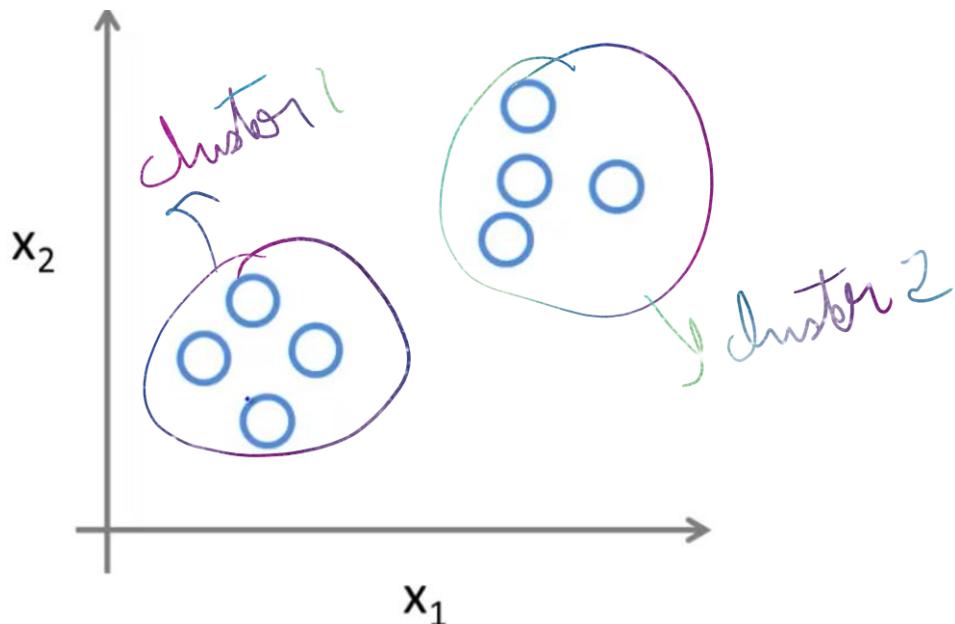
Classification:

It means grouping data into discrete values like 0 , 1 ,2 etc.

Eg: In eg 1 given above, benign tumors are represented by 0 and malignant tumors by 1.

UNSUPERVISED LEARNING:

In unsupervised learning we present the programme with just the data , which has not been grouped or classified. We ask the programme to classify it.

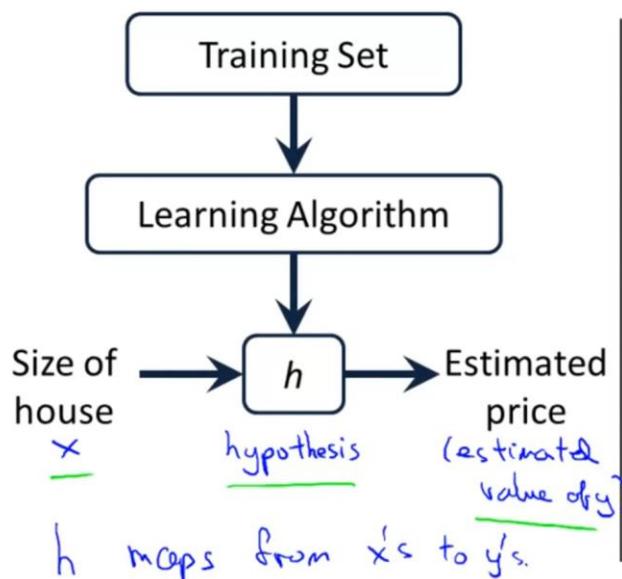


Here we ask the programme to classify the points plotted on the graph . The programme may classify the points as cluster 1 and cluster 2.

Another example may be that we may have two audio recordings with two overlapping voices. Suppose voice 1 is more significant in audio 1 and voice 2 in audio 2. Now we can ask the

programme to separate voice 1 and 2 from both the recordings . Now after the separation recording 1 has only voice 1 and recording 2 has only voice 2. This is an example of unsupervised learning.

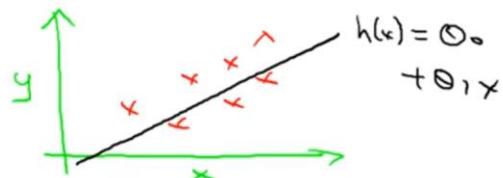
MODEL REPRESENTATION:



How do we represent h ?

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

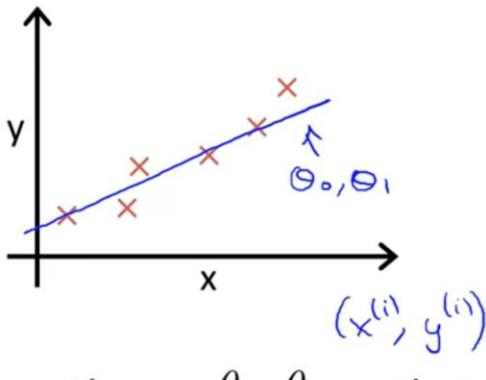
Shorthand: $h(x)$



Linear regression with one variable.
Univariate linear regression.

Andrew Ng

COST FUNCTION:



$$\text{minimize}_{\theta_0, \theta_1} \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

\uparrow

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

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Idea: Choose θ_0, θ_1 so that
 $h_\theta(x)$ is close to y for our
 training examples (x, y)

x, y

Minimize $J(\theta_0, \theta_1)$
 θ_0, θ_1 $\underbrace{\quad}_{\text{Cost function}}$
 Squared error function

Andrew Ng

Cost function is also known as the squared error function:

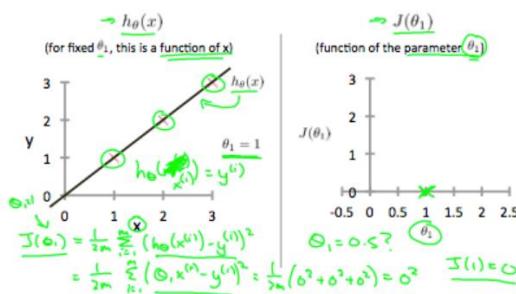
Let us consider our constant to be 0, so the equation becomes:

$$h_\theta(x) = \theta_0 x$$

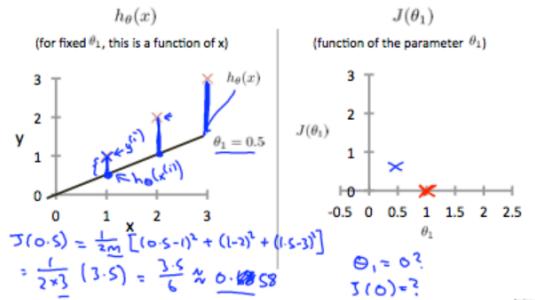
Cost Function - Intuition I

If we try to think of it in visual terms, our training data set is scattered on the x-y plane. We are trying to make a straight line (defined by $h_\theta(x)$) which passes through these scattered data points.

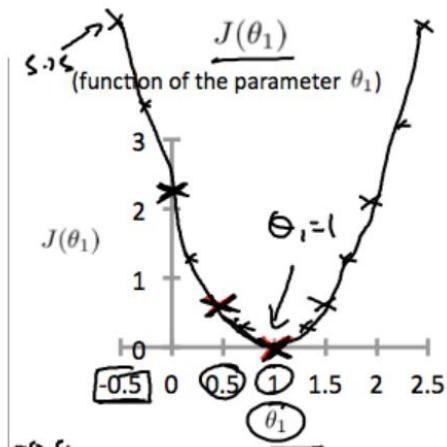
Our objective is to get the best possible line. The best possible line will be such so that the average squared vertical distances of the scattered points from the line will be the least. Ideally, the line should pass through all the points of our training data set. In such a case, the value of $J(\theta_0, \theta_1)$ will be 0. The following example shows the ideal situation where we have a cost function of 0.



When $\theta_1 = 1$, we get a slope of 1 which goes through every single data point in our model. Conversely, when $\theta_1 = 0.5$, we see the vertical distance from our fit to the data points increase.



This increases our cost function to 0.58. Plotting several other points yields to the following graph:



Thus as a goal, we should try to minimize the cost function. In this case, $\theta_1 = 1$ is our global minimum.

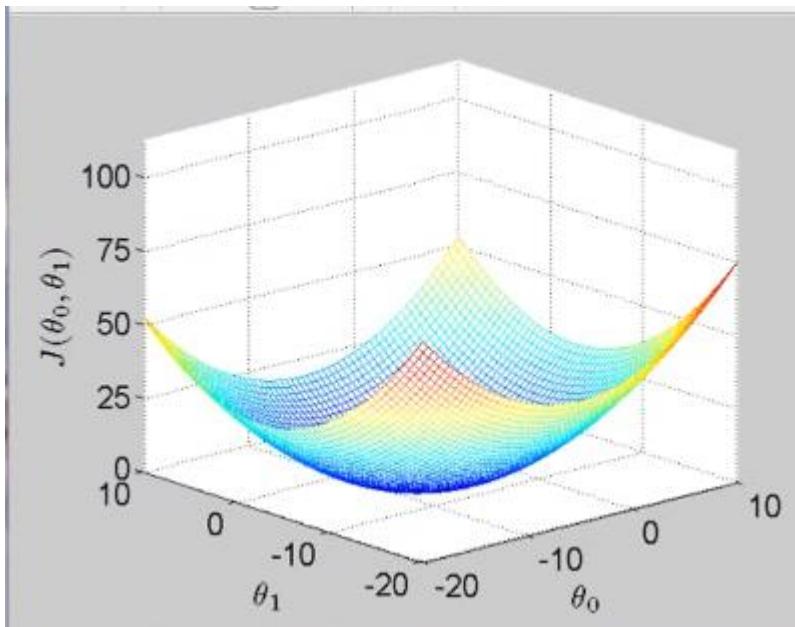
And hence $h_\theta(x) = 1 \cdot x$
 $h_\theta(x) = x$ is our final equation for which cost function is minimum.

Now let us assume our equation to be:

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

Our aim to find the value of θ_0 & θ_1 , for which the cost function is minimum.

Graph of cost function vs θ_0 & θ_1



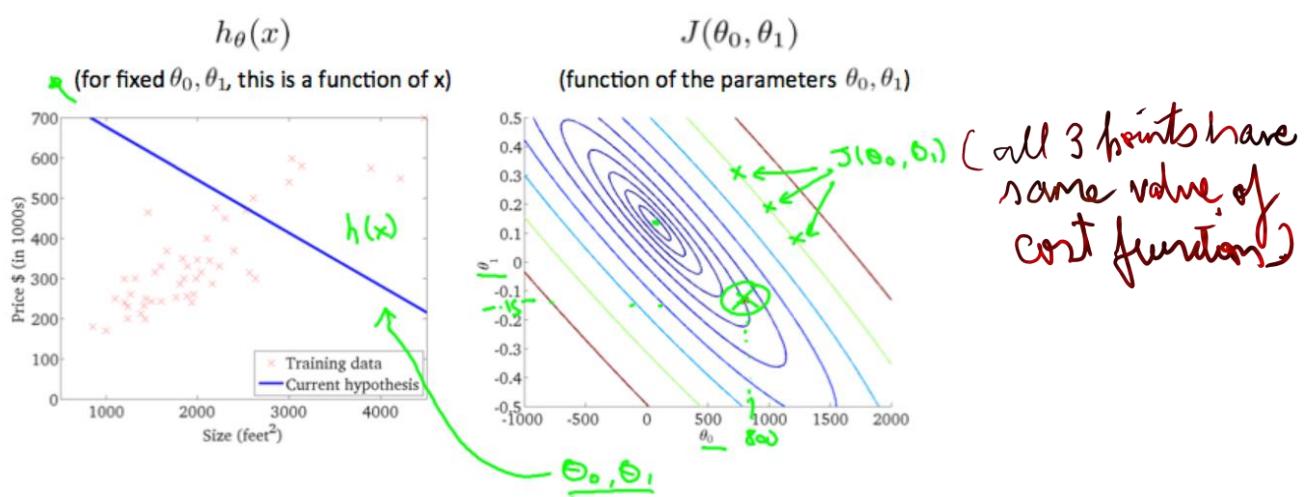
We infer from the above graph that a particular value of the cost function can be achieved through different combinations of different values of θ_0 & θ_1 .

For eg , the value 25 of the cost function may be achieved through various combinations of values of θ_0 & θ_1 .

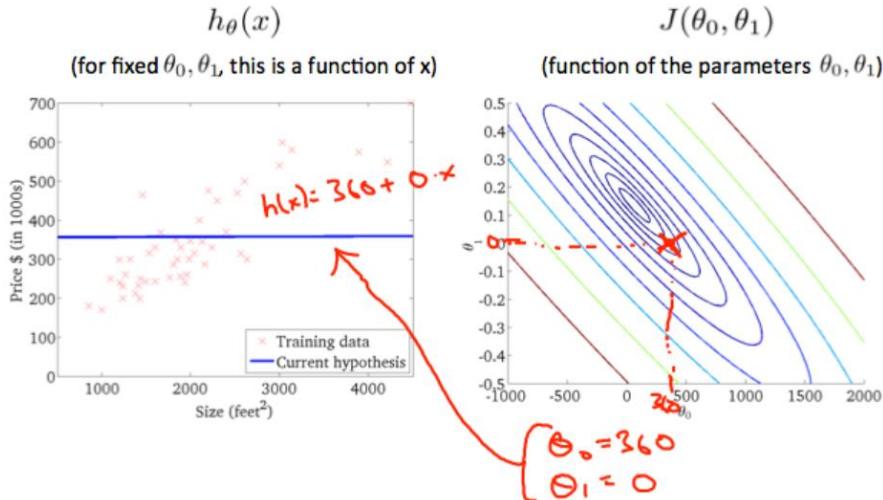
Above visualization in 2-d using contour plots:

Cost Function - Intuition II

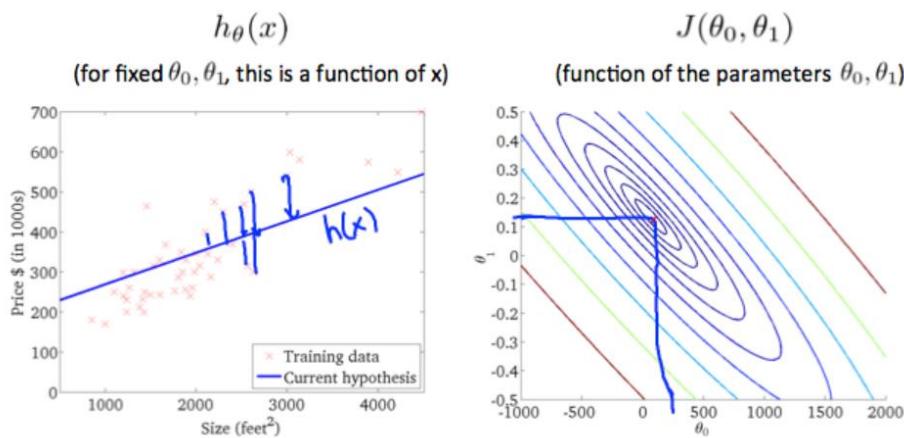
A contour plot is a graph that contains many contour lines. A contour line of a two variable function has a constant value at all points of the same line. An example of such a graph is the one to the right below.



Taking any color and going along the 'circle', one would expect to get the same value of the cost function. For example, the three green points found on the green line above have the same value for $J(\theta_0, \theta_1)$ and as a result, they are found along the same line. The circled x displays the value of the cost function for the graph on the left when $\theta_0 = 800$ and $\theta_1 = -0.15$. Taking another $h(x)$ and plotting its contour plot, one gets the following graphs:



When $\theta_0 = 360$ and $\theta_1 = 0$, the value of $J(\theta_0, \theta_1)$ in the contour plot gets closer to the center thus reducing the cost function error. Now giving our hypothesis function a slightly positive slope results in a better fit of the data.



The graph above minimizes the cost function as much as possible and consequently, the result of θ_1 and θ_0 tend to be around 0.12 and 250 respectively. Plotting those values on our graph to the right seems to put our point in the center of the inner most 'circle'.

GRADIENT DESCENT

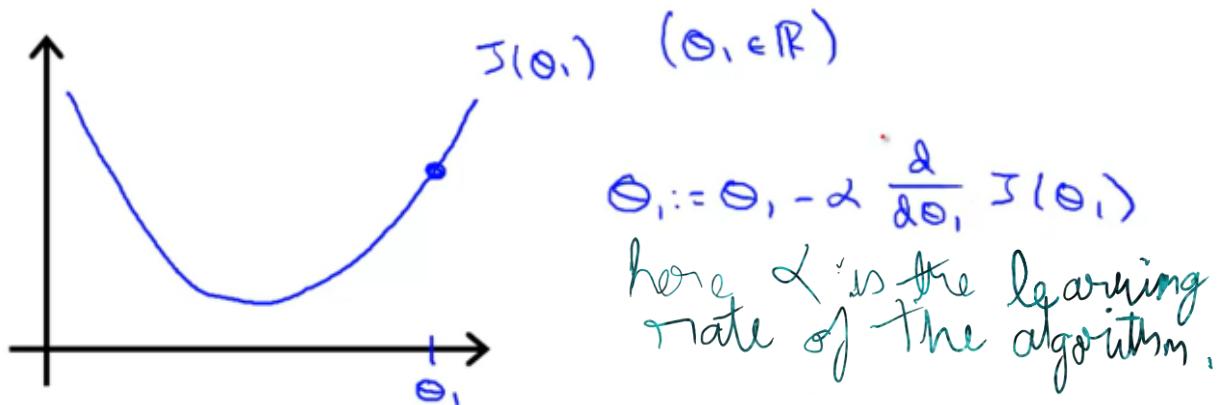
Gradient descent is a technique used to find minimum value of any function $J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

Gradient descent algorithm

```
repeat until convergence {
     $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$    (for  $j = 0$  and  $j = 1$ )
}
```

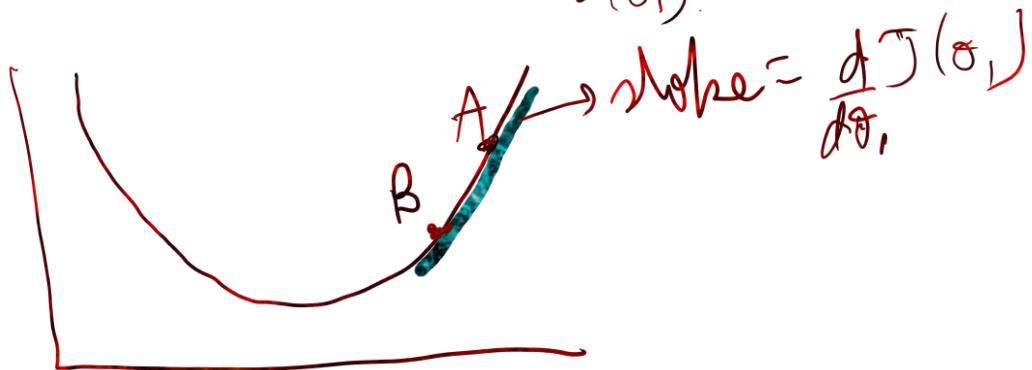
First we look at a very simple example where we have a function $J(\theta_1)$.

Now, we have to minimize the value of this function.



$J(\theta_1) := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$, α means we are updating value of θ_1 to $\theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$.

Updating theta1 means that we are changing value of theta 1 till we reach the minimum value of theta1 which in turn minimizes the function $J(\theta_1)$.



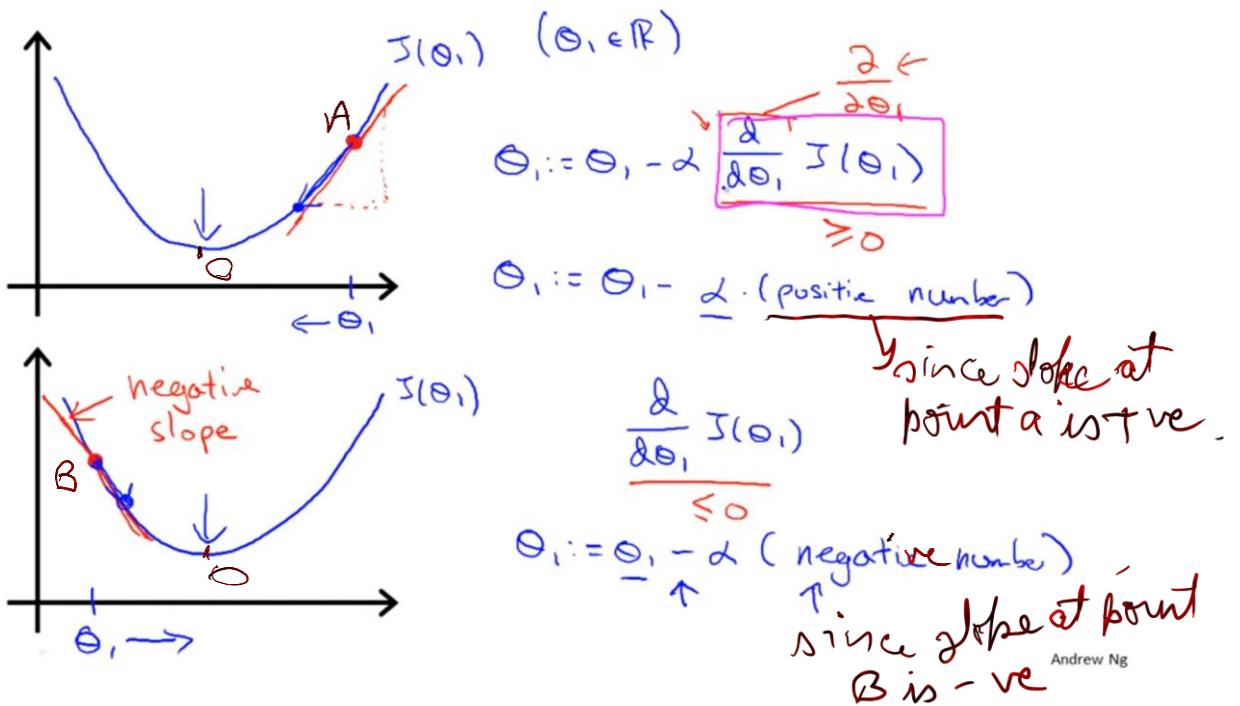
Let slope at point A = a. Hence, Slope = $\frac{d(J(\theta_{1A}))}{d\theta_1} = a$

As we move from point A to point B, θ_1 changes:

$$\theta_{1B} = \theta_{1A} - \alpha \frac{d}{d\theta_1} (J(\theta_{1A}))$$

here θ_{1B} represents change of we move from A to B.

NOTE: Higher the value of slope alpha, faster the change in theta1.



In graph 1, θ_1 value keeps on decreasing as we move from point A to point o (point of minima).

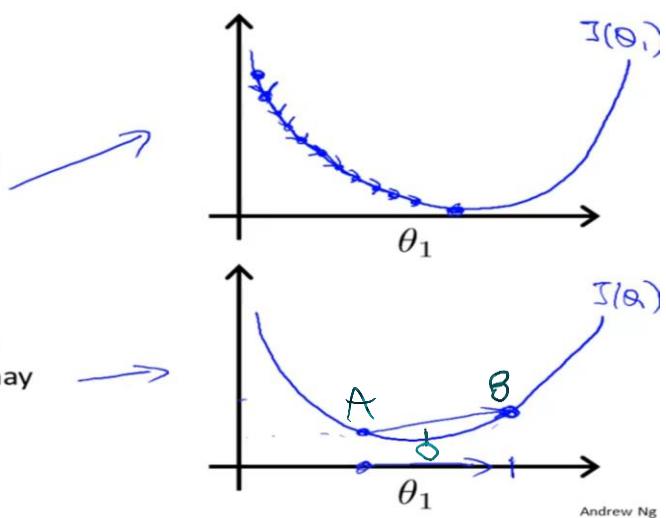
In graph 2, θ_1 value keeps on increasing as we move from point B to point o.

Some imp. points pertaining to alpha:

$$\theta_1 := \theta_1 - \alpha \frac{\partial J(\theta_1)}{\partial \theta_1}$$

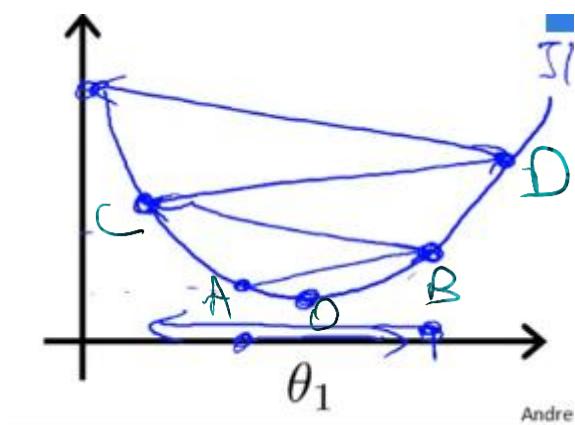
If α is too small, gradient descent can be slow.

If α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.

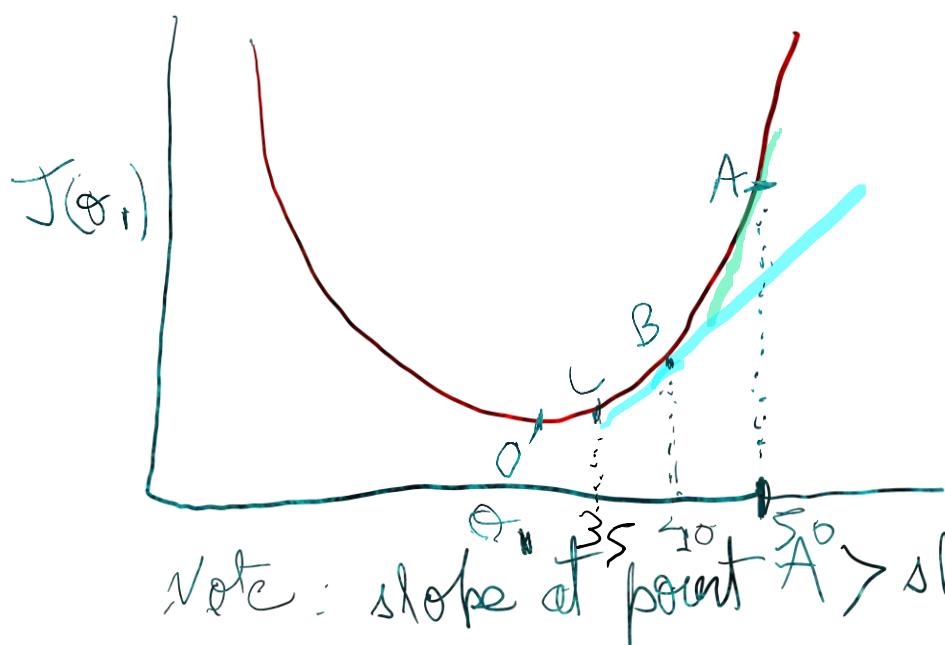


In graph 2, as we change try to go from point A to point O (point of minima), there is positive change in value of θ_1 , i.e. value of θ_1 increases yet we do not reach point O

and reach point B since value of alpha is so large , then we try to come back to point O but since alpha is so large , we reach a point C , and then point D and so on. Eventually we see we keep getting farther and farther from point O instead of getting closer to it.



Divergence of gradient descent.



Note: slope at point A > slope at point B.

As we move from point A to point B , value of θ_1 changes from 35 to 40 , i.e. change of 10 units.

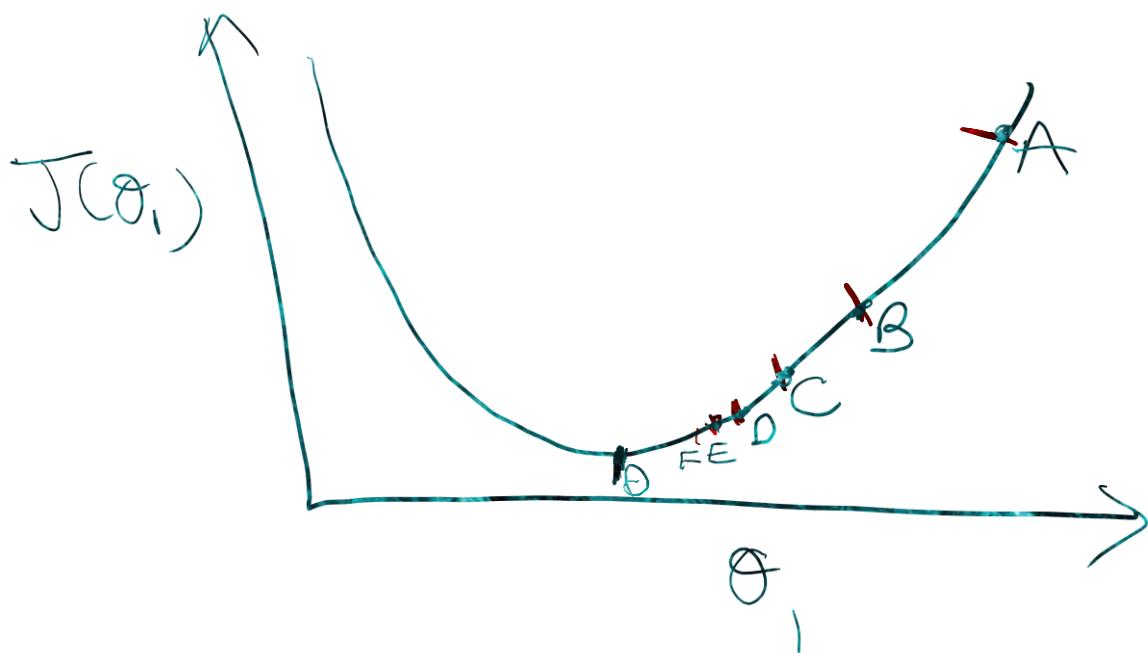
As we move from point B to point C , value of θ_1 changes from 40 to 35 , i.e. change of 5 units.

We see less change in value of θ , because the slope at point A is \leftarrow slope at point A.

$$\theta_1 = \theta_0 - \frac{d}{J(\theta_0)} (J(\theta_1))$$

Hence as we keep getting closer to point of minima θ , the steps we take from one point to another keep decreasing in magnitude.

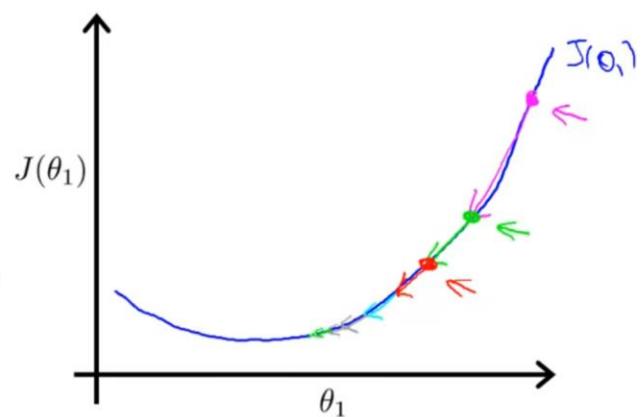
As we get very close to θ , the steps taken from one point to another become extremely small.



Gradient descent can converge to a local minimum, even with the learning rate α fixed.

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease α over time.



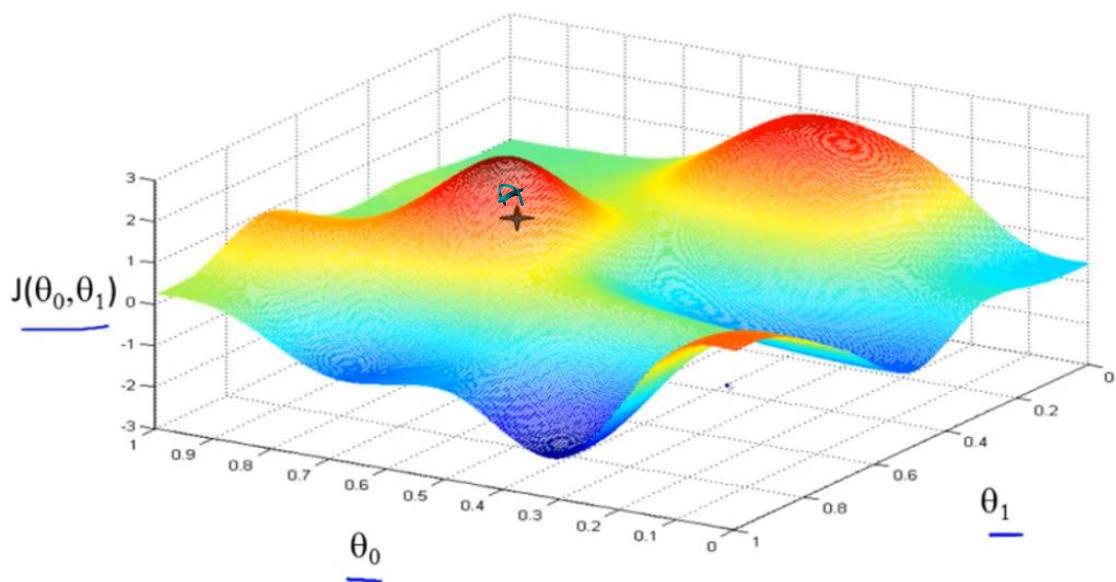
Andrew Ng

s

Now, We see another example in which we have the function $J(\theta_0, \theta_1)$

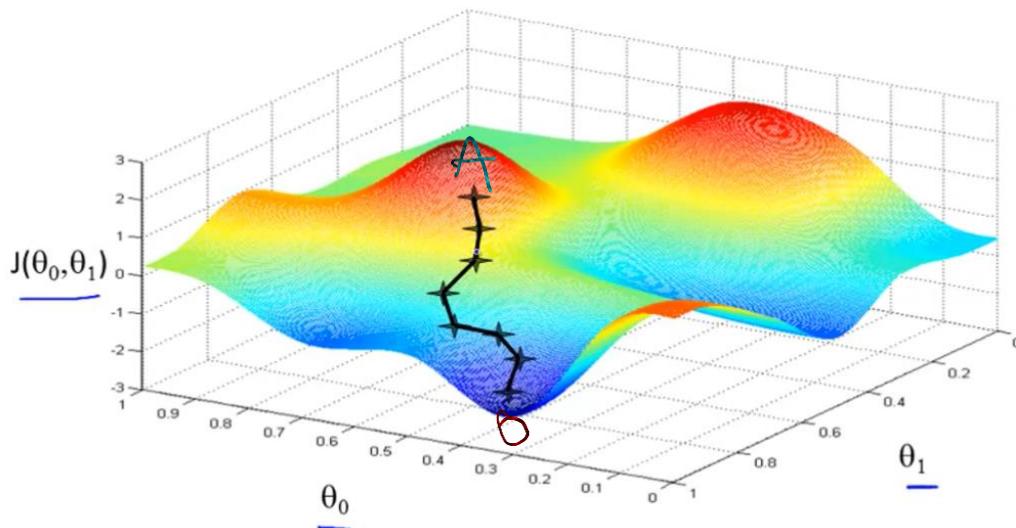
We want to reach point of minima for which the function has minimum value for the given θ_0, θ_1

Let us assume we are at point A which is represent by the cross on the graph of the function:



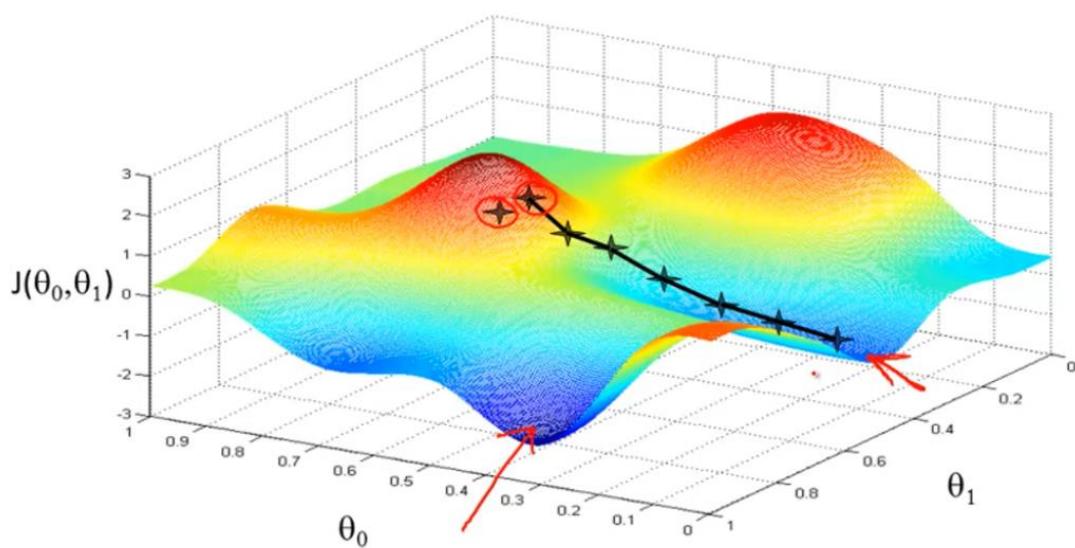
Now, we want to reach the point of local minima. Let's suppose that in this pursuit of reaching the point of minima, we start taking baby steps down the slope from point A .

We keep going until we reach a point O.

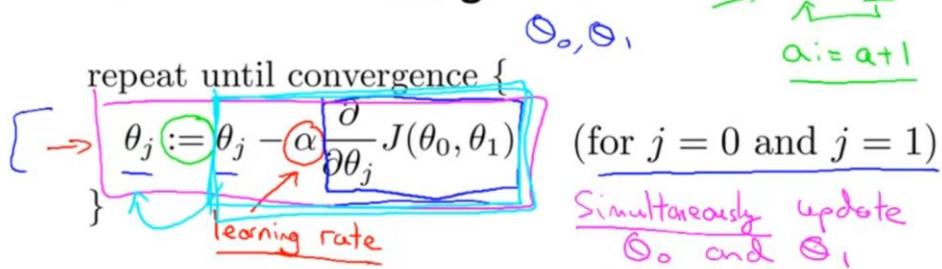


As we take each step starting from point A , value of θ_1 / θ_0 changes with each step till we reach the point of minima.

##If we would have moved a bit towards right from A and then started moving downwards , we would have reached a different local minima.



Gradient descent algorithm



Assignment
 $a := b$
 $a := a + 1$

Truth assertion
 $a = b$ ✓
 $a = a + 1$ ✗

Correct: Simultaneous update

→ $\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$
 → $\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$
 → $\theta_0 := \text{temp0}$
 → $\theta_1 := \text{temp1}$

Incorrect:

→ $\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$
 → $\theta_0 := \text{temp0}$
 → $\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$
 → $\theta_1 := \text{temp1}$

wrong because value of θ_0 has been updated in this equation
 Andrew Ng

Getting updated value of θ_0 & θ_1 at each step as we move from any point to the point of minima:

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{2}{2\theta_j} \cdot \frac{1}{2m} \sum_{i=1}^m (\underline{h_\theta(x^{(i)}) - y^{(i)}})^2$$

$$= \frac{2}{2\theta_j} \frac{1}{2m} \sum_{i=1}^m (\underline{\theta_0 + \theta_1 x^{(i)}} - y^{(i)})^2$$

$$\theta_0, j = 0 : \underline{\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)} = \frac{1}{m} \sum_{i=1}^m (\underline{h_\theta(x^{(i)}) - y^{(i)}})$$

$$\theta_1, j = 1 : \underline{\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)} = \frac{1}{m} \sum_{i=1}^m (\underline{h_\theta(x^{(i)}) - y^{(i)}}) \cdot \underline{x^{(i)}}$$

Explanation:

$$\theta_0 = \underline{\frac{1}{2\theta_0}} \sum_{i=1}^m (\underline{\theta_0 + \theta_1 x^{(i)}} - \underline{y^{(i)}})^2$$

$$\begin{aligned}
 &= \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2 \\
 &= \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) \times 1
 \end{aligned}$$

Similar goes for $\frac{\partial}{\partial \theta_1}$

Note: Above logic can also be extended to multivariate regression:

Gradient Descent

Previously ($h=1$):

Repeat {

$$\begin{aligned}
 \rightarrow \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\
 &\quad \boxed{\frac{\partial}{\partial \theta_0} J(\theta)}
 \end{aligned}$$

}

$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

(simultaneously update θ_0, θ_1)

New algorithm ($n \geq 1$):

Repeat {

$$\downarrow \frac{\partial}{\partial \theta_j} J(\theta)$$

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

(simultaneously update θ_j for $j = 0, \dots, n$)

}

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$\rightarrow \theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

...

Andrew Ng

Note: Superscript (i), is not usually not used to denote rows or columns but simply refers to the i^{th} training example. For eg: $x^{(2)}$, refers to training 2nd example.

Let us assume the equation has 2 variables i.e. $n=2$.

$$h(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 \text{ where } x_0 = 1$$

$$\textcircled{1} \frac{\partial}{\partial \theta_2} \left[\frac{1}{2m} \sum_{i=1}^m (\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2)^2 \right] \rightarrow \alpha$$

$$\text{let } j=2 \rightarrow \frac{1}{2m} \sum_{i=1}^m (\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2) \times (0 + 0 + x_2)$$

$$\begin{aligned}
 \textcircled{1} \text{ becomes: } \frac{\partial}{\partial \theta_2} \alpha &\Rightarrow \frac{1}{2m} \sum_{i=1}^m (\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2) x_2 \\
 &= \frac{1}{m} \sum_{i=1}^m (\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2) x_2
 \end{aligned}$$

QUESTION:

When there are n features, we define the cost function as

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

For linear regression, which of the following are also equivalent and correct definitions of $J(\theta)$?

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

Correct

$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left(\left(\sum_{j=0}^n \theta_j x_j^{(i)} \right) - y^{(i)} \right)^2$ (Inner sum starts at 0)

Correct

$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left(\left(\sum_{j=1}^n \theta_j x_j^{(i)} \right) - y^{(i)} \right)^2$ (Inner sum starts at 1)

Un-selected is correct

$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left(\left(\sum_{j=0}^n \theta_j x_j^{(i)} \right) - \left(\sum_{j=0}^n y_j^{(i)} \right) \right)^2$

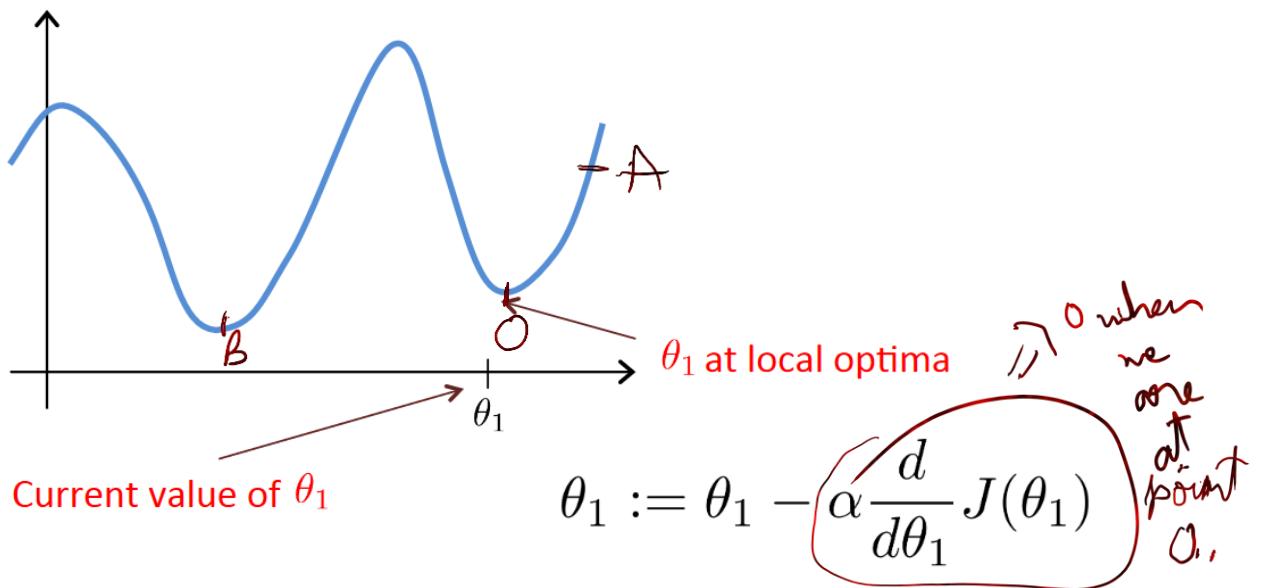
Un-selected is correct

Continue

so that theta here is a vector.

Andrew Ng

Theta always converges to the point of local minima due to gradient descent:



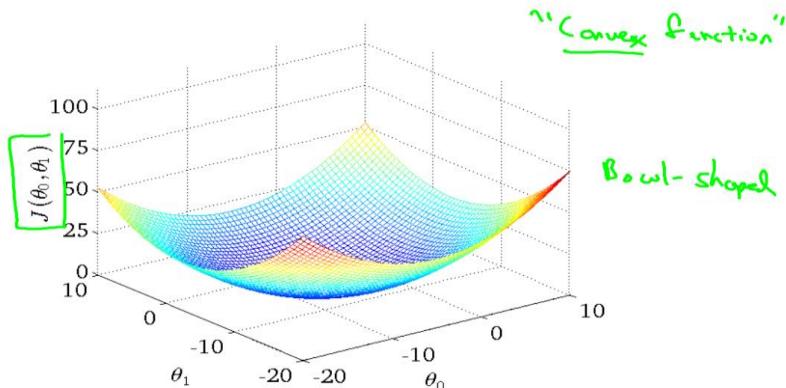
Let us assume that we start from point A and proceed towards point O. We know O is the local minima , while point B is the global minima. As we reach point O , there can be no further change in value of theta1 this is because slope at point O is = 0 , hence in accordance with the above given formula , change in theta 1 = 0. Since $\theta_1 := \theta_1 - \alpha \nabla J$

$$\theta_1 := \theta_1 - \alpha \nabla J$$

$$\theta_1 := \theta_1$$

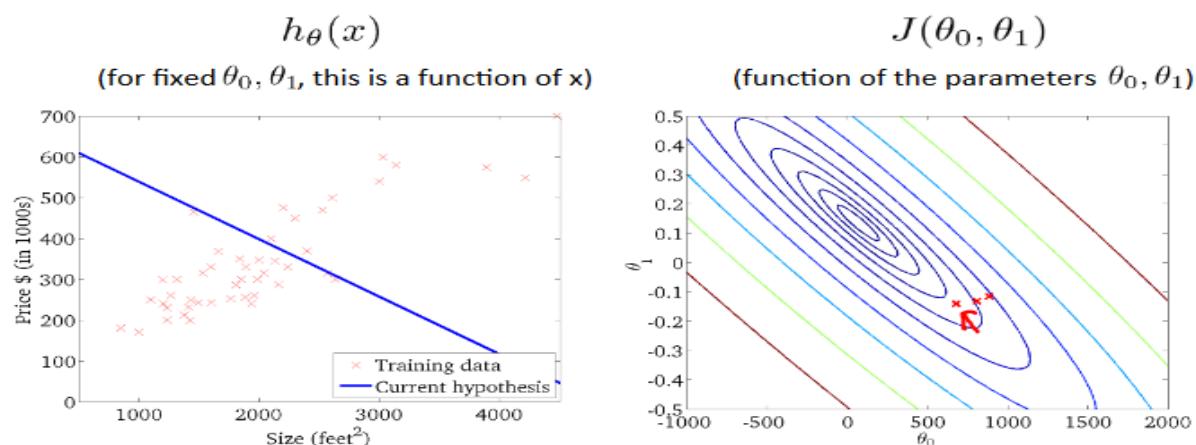
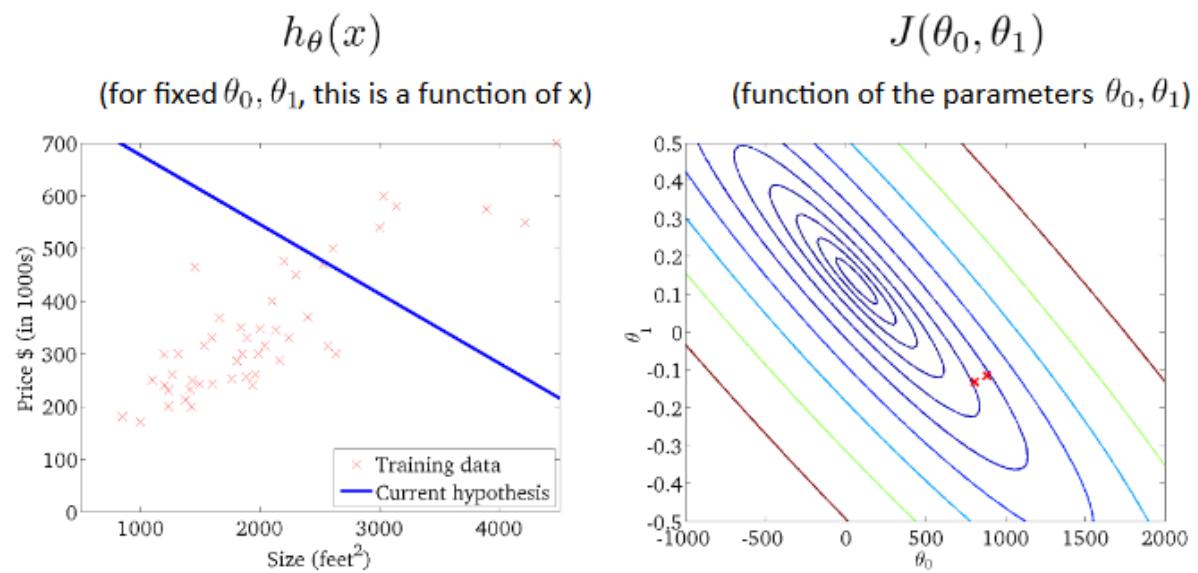
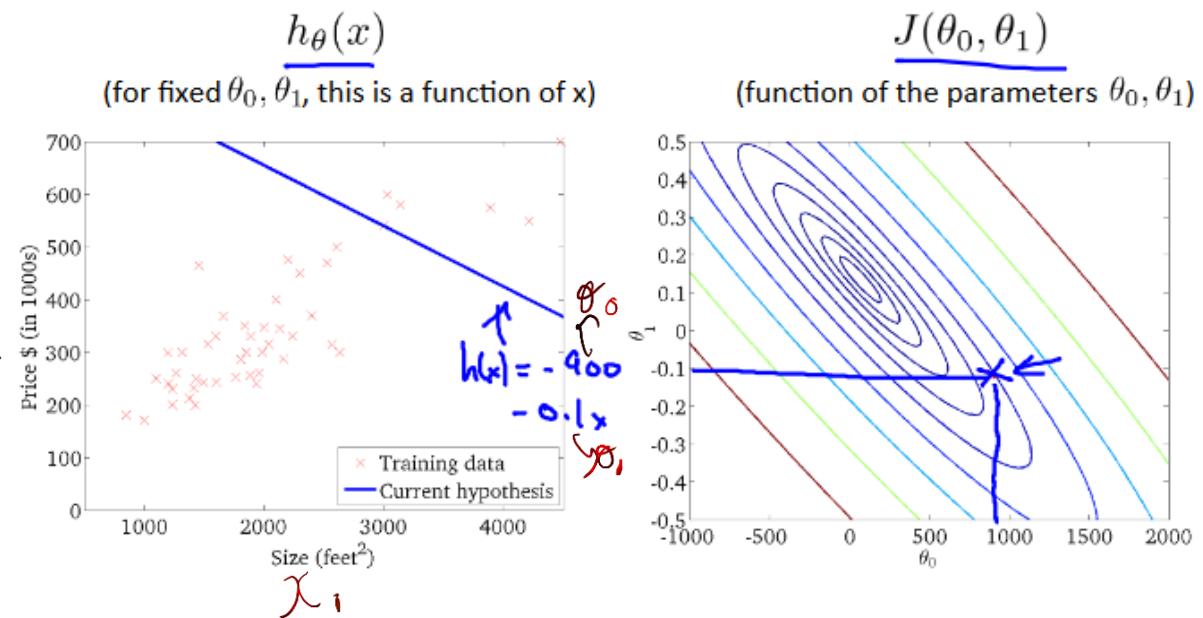
COST FUNCTION FOR LINEAR REGRESSION:

It turns out that that the cost function for linear regression is always going to be a bow shaped function like this:



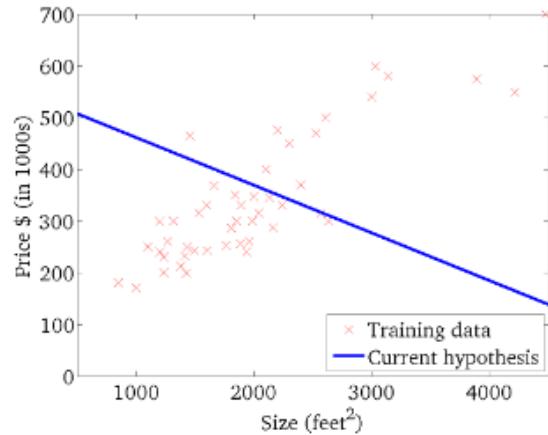
The technical term for this is that this is called a convex function.

AN EXAMPLE SHOWING GRADIENT DESCENT IN WORK:



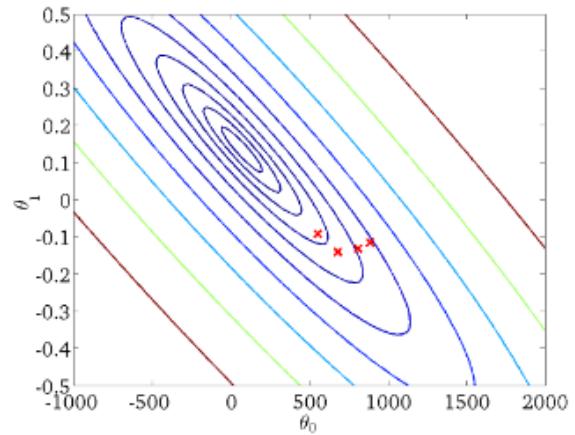
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



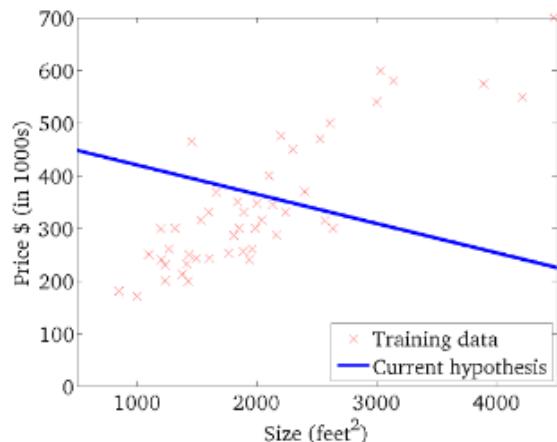
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



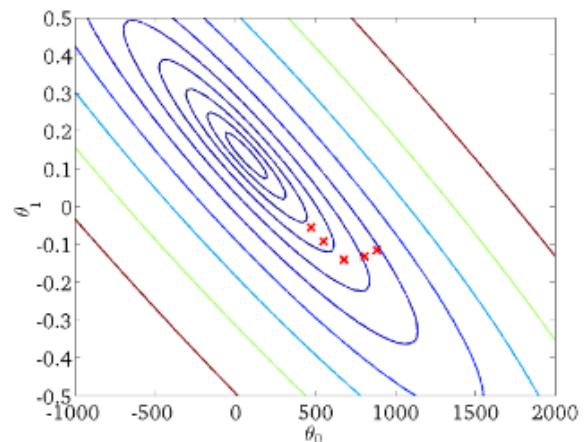
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



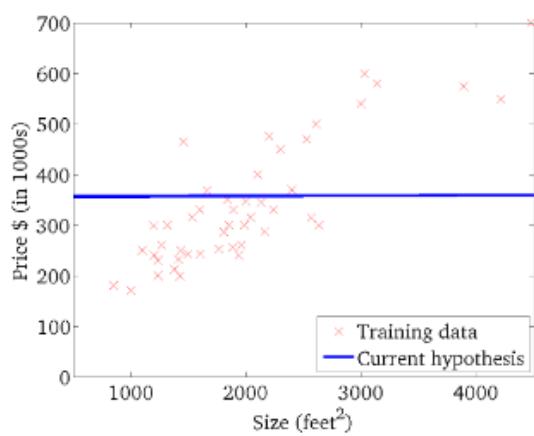
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



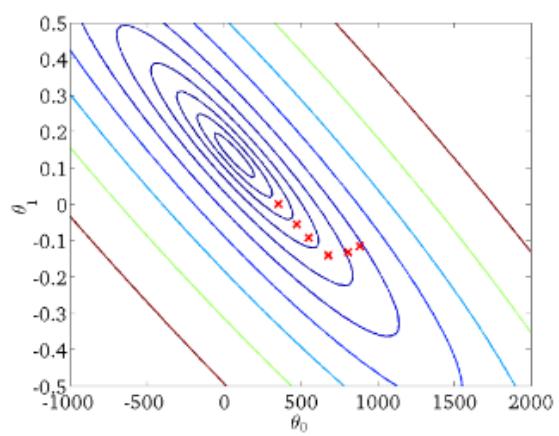
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)

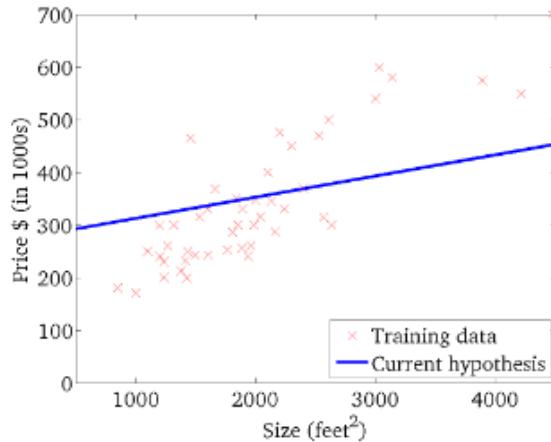


$$J(\theta_0, \theta_1)$$

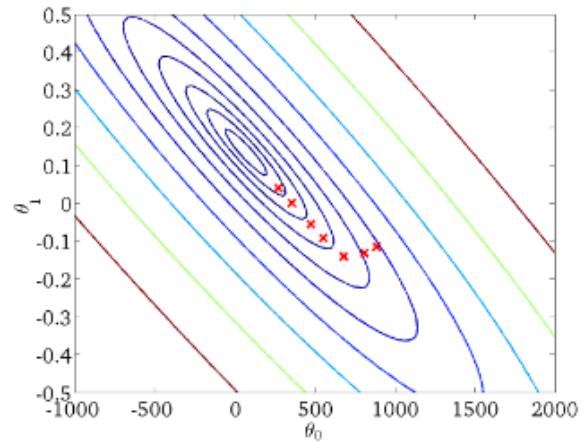
(function of the parameters θ_0, θ_1)



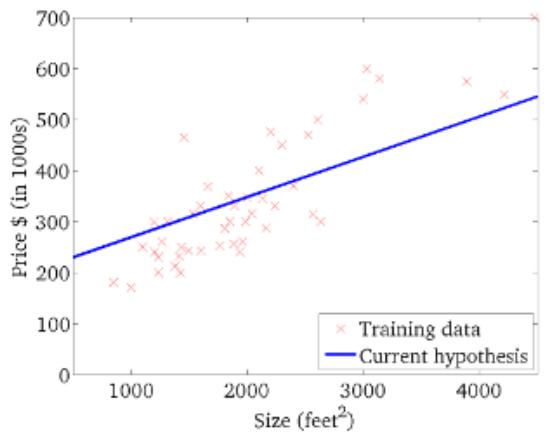
$h_\theta(x)$
(for fixed θ_0, θ_1 , this is a function of x)



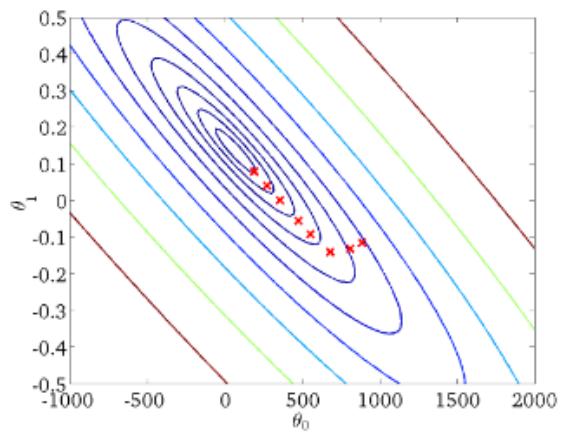
$J(\theta_0, \theta_1)$
(function of the parameters θ_0, θ_1)



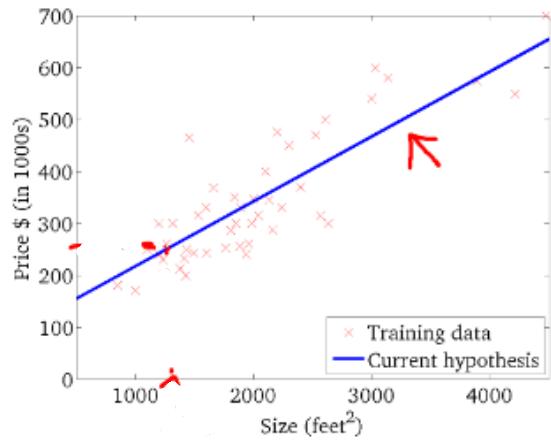
$h_\theta(x)$
(for fixed θ_0, θ_1 , this is a function of x)



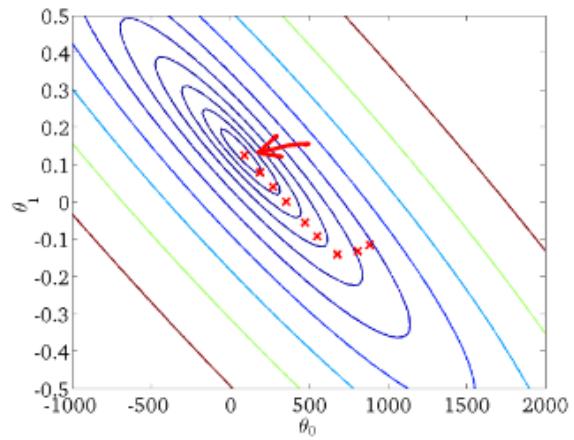
$J(\theta_0, \theta_1)$
(function of the parameters θ_0, θ_1)



$h_\theta(x)$
(for fixed θ_0, θ_1 , this is a function of x)



$J(\theta_0, \theta_1)$
(function of the parameters θ_0, θ_1)



In the above pictures as we keep getting closer to the perfect fit we start getting closer to the centre of the contour plot i.e., the point of minima.

FEATURE SCALING:

Why do we need feature scaling?

consider data (mtcars)

We make a basic model:

$$> lm(mtcars ~ mtcars$wt + mtcars$disp)$$

coefficients	Intercept	wt (θ_1)	disp (θ_2)
	34.96055	-3.35083	-0.01772

→ we see that θ_1 & θ_2 vary greatly in scale, this is due to the fact that wt & disp vary greatly in scale:

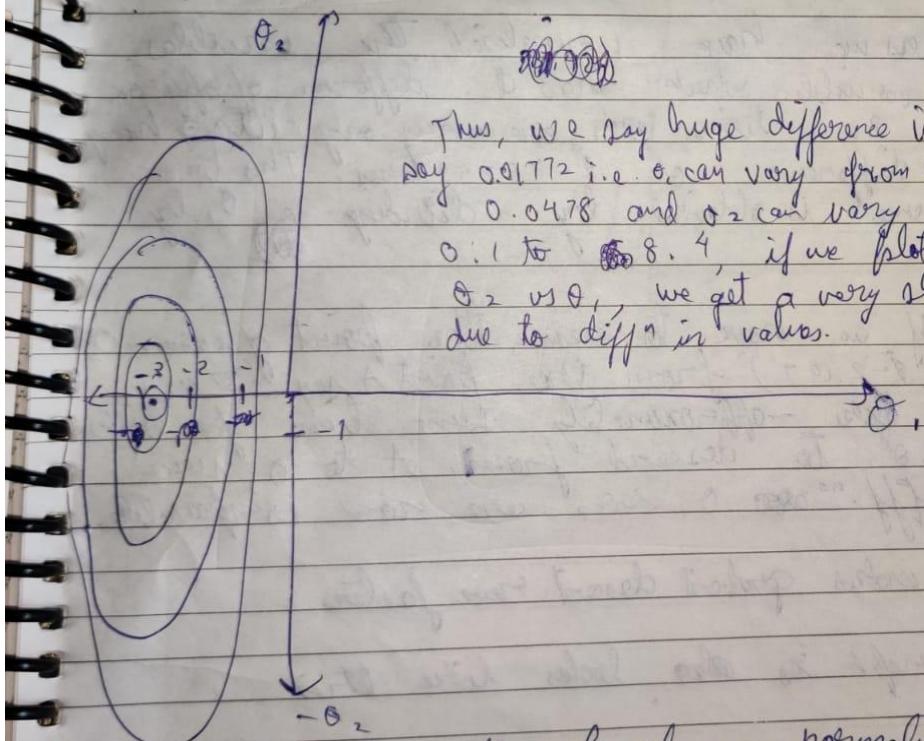
$$> range(mtcars$wt)$$

o/p (1) 1.513 5.424
$$> range(mtcars$disp)$$

o/p (1) 71.1 472.0

if we plot $J(\theta_1, \theta_2)$ vs (θ_1) vs (θ_2), we see the contour plot is highly skewed, this is due to the fact that θ_1 & θ_2 vary greatly in scale. (a difference of 300 times!) If we start from an extreme point say A on the contour plot & try to

reach point 0, i.e. the point of origin it would take a lot of time to reach this point since reaching $(-3.35, -0.01)$ from a point say $t = (37, 42)$ is very tough. Reaching 0.01 is especially tough. θ_2 will take a large time to descend from A to 0. hence gradient descent will never scale the The graph looks like this slow



Thus, we say huge difference in scale, say 0.01772 i.e. θ_1 can vary from 0.003 to 0.0478 and θ_2 can vary from 0.1 to 8.4, if we plot corresponding θ_2 vs θ_1 , we get a very skewed graph due to diff' in values.

If on the other hand we normalize the variables:
 $w \leftarrow (\text{mtcars} \$ wt - \text{mean}(\text{mtcars} \$ wt)) / \text{sd}(\text{mtcars} \$ wt)$
 $d \leftarrow (\text{mtcars} \$ disp - \text{mean}(\text{mtcars} \$ disp)) / \text{sd}(\text{mtcars} \$ disp)$

> $\ln(\text{mpg}) \approx w + d$

coefficients:

(Intercept)

20.091

$w(0_1) \& d(0_2)$

-3.279 -2.197

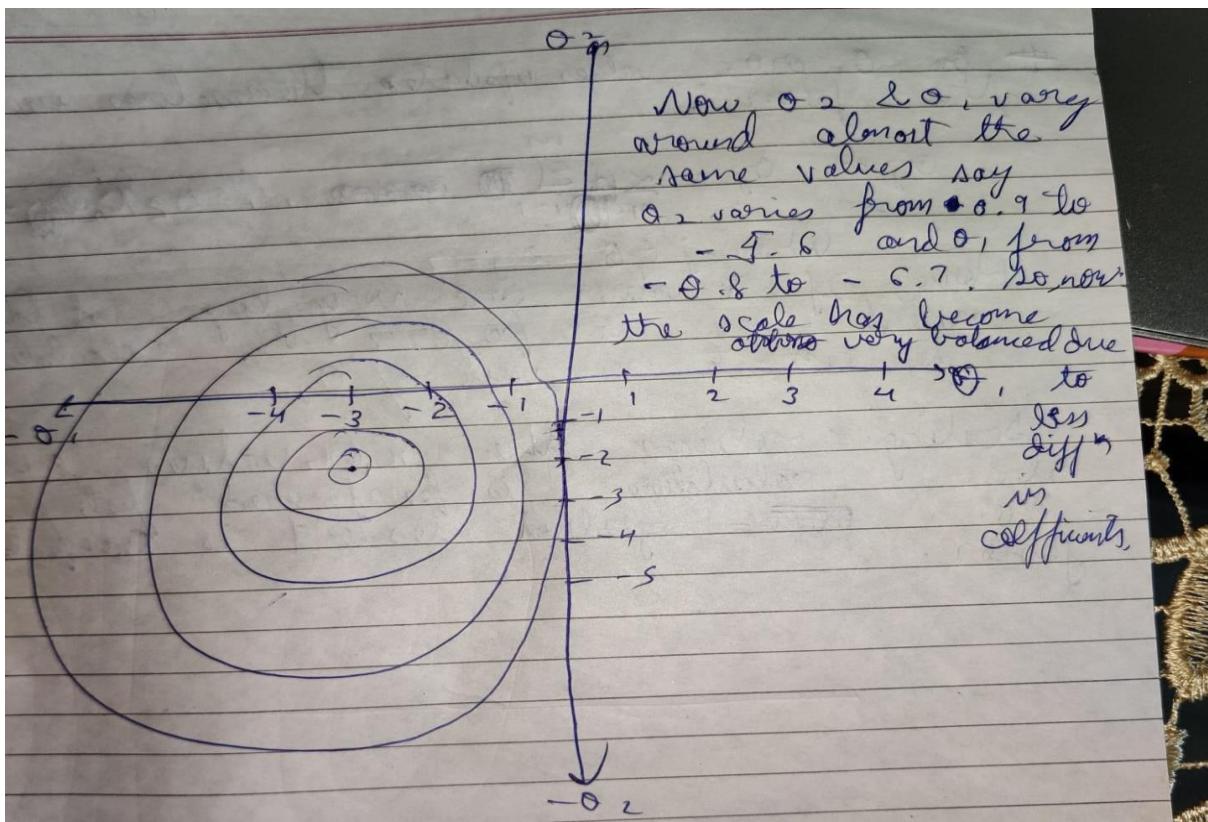
Now, as we have normalized the variables, the variables which had a difference of approximatively 300 times are now only having a difference of 1.4924 times. This times difference is obtained by dividing 0_1 by 0_2 .

Now, if we have to reach the point of minima 0 , (-3.279, -2.197) from the point A say (37.42), it would take approximately same time for both 0_1 & 0_2 , to descend from A to 0, since the coff. w & d are now comparable.

This makes gradient descent run faster.

The graph is also looks like this:

P.T.O



Page No. _____
Date : _____

the $\theta_0, \theta_1, \theta_2$ values update themselves using the formula:

$$\theta_j := \theta_j - \frac{\alpha}{m} \sum_{i=1}^m (y_i - \hat{y}_i) x_j^{(i)}$$

where $j \in \{1, 2\}$
 Hence, if x values are very large it takes more time to perform the calculation to compute new value of θ_j .

Scaling reduces the x_1, x_2 values & hence above calculation is performed faster & the ~~above~~ gradient descent is computed faster.

Feature Scaling

Get every feature into approximately a $-1 \leq x_i \leq 1$ range.

$$x_0 = 1$$

$$\begin{aligned}6 &\leq x_1 \leq 3 & \checkmark \\-2 &\leq x_2 \leq 0.5 & \checkmark \\-100 &\leq x_3 \leq 100 & \times \\-0.0001 &\leq x_4 \leq 0.0001 & \times\end{aligned}$$

$$-1 \leq x_i \leq 1$$

$$\begin{aligned}-3 &\rightarrow 3 & \checkmark \\-\frac{1}{3} &\rightarrow \frac{1}{3} & \checkmark\end{aligned}$$

These are also
correct scales

based on
requirement

Methods for scaling:

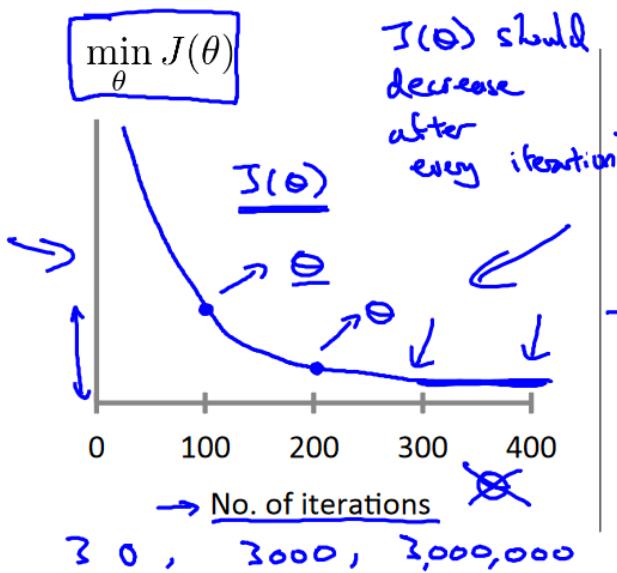
1. Divide the data by upper limit of range - lower limit of range

2. mean normalisation i.e.

$$\frac{\text{data} - \text{mean(data)}}{\text{std(data)}}$$

GRADIENT DESCENT AND ALPHA

Making sure gradient descent is working correctly.



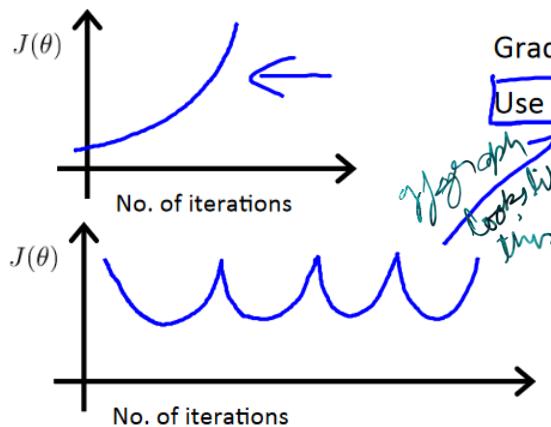
→ Example automatic convergence test:

→ Declare convergence if $J(\theta)$ decreases by less than 10^{-3} in one iteration.

$J(\theta)$ is not always easy to declare such a threshold value.

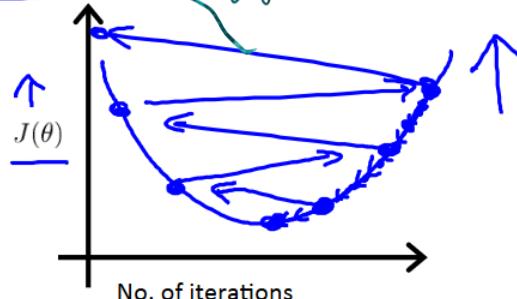
Andrew Ng

Making sure gradient descent is working correctly.



Gradient descent not working.
Use smaller α .

graph looks like this:
graph looks like this:

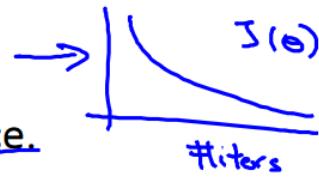


- For sufficiently small α , $J(\theta)$ should decrease on every iteration.
- But if α is too small, gradient descent can be slow to converge.

Andrew Ng

Summary:

- If α is too small: slow convergence.
- If α is too large: $J(\theta)$ may not decrease on every iteration; may not converge. (Slow converge also possible)



Slow converge
also possible

To choose α , try

$$\dots, \underbrace{0.001, 0.003}_{\xrightarrow{3x}}, \underbrace{0.01, 0.03}_{\approx 3x}, \underbrace{0.1, 0.3}_{\approx 3x}, 1, \dots$$

Andrew Ng

Above line means we can first assume alpha to be 0.001 and then increase alpha in multiples of 3 and see the no. of iterations.

Question:

Making graphs

Suppose a friend ran gradient descent three times, with $\alpha = 0.01$, $\alpha = 0.1$, and $\alpha = 1$, and got the following three plots (labeled A, B, and C):

Which plots corresponds to which values of α ?

A is $\alpha = 0.01$, B is $\alpha = 0.1$, C is $\alpha = 1$.
 A is $\alpha = 0.1$, B is $\alpha = 0.01$, C is $\alpha = 1$.

Correct
In graph C, the cost function is increasing, so the learning rate is set too high. Both graphs A and B converge to an optimum of the cost function, but graph B does so very slowly, so its learning rate is set too low. Graph A lies between the two.

A is $\alpha = 1$, B is $\alpha = 0.01$, C is $\alpha = 0.1$.
 A is $\alpha = 1$, B is $\alpha = 0.1$, C is $\alpha = 0.01$.

Continue

Andrew Ng

be slow to converge.

POLYNOMIAL REGRESSION:

Sometimes fitting a straight line through the data may not help, we may need to fit a line that corresponds to degree 2 , degree 3 , degree $\frac{1}{2}$,etc . equation .

Features and Polynomial Regression

We can improve our features and the form of our hypothesis function in a couple different ways.

We can **combine** multiple features into one. For example, we can combine x_1 and x_2 into a new feature x_3 by taking $x_1 \cdot x_2$.

Polynomial Regression

Our hypothesis function need not be linear (a straight line) if that does not fit the data well.

We can **change the behavior or curve** of our hypothesis function by making it a quadratic, cubic or square root function (or any other form).

For example, if our hypothesis function is $h_{\theta}(x) = \theta_0 + \theta_1 x_1$ then we can create additional features based on x_1 , to get the quadratic function $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$ or the cubic function $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$

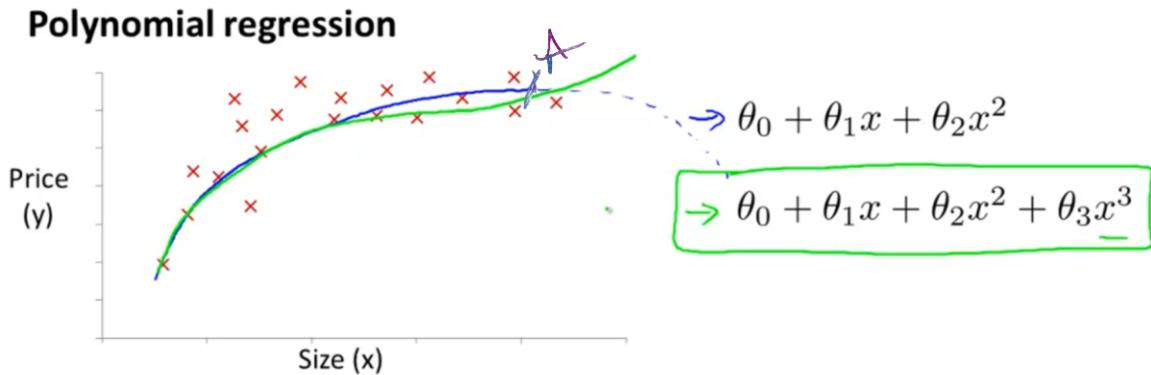
In the cubic version, we have created new features x_2 and x_3 where $x_2 = x_1^2$ and $x_3 = x_1^3$.

To make it a square root function, we could do: $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_1}$

One important thing to keep in mind is, if you choose your features this way then feature scaling becomes very important.

eg. if x_1 has range 1 - 1000 then range of x_1^2 becomes 1 - 1000000 and that of x_1^3 becomes 1 - 1000000000

As an example:

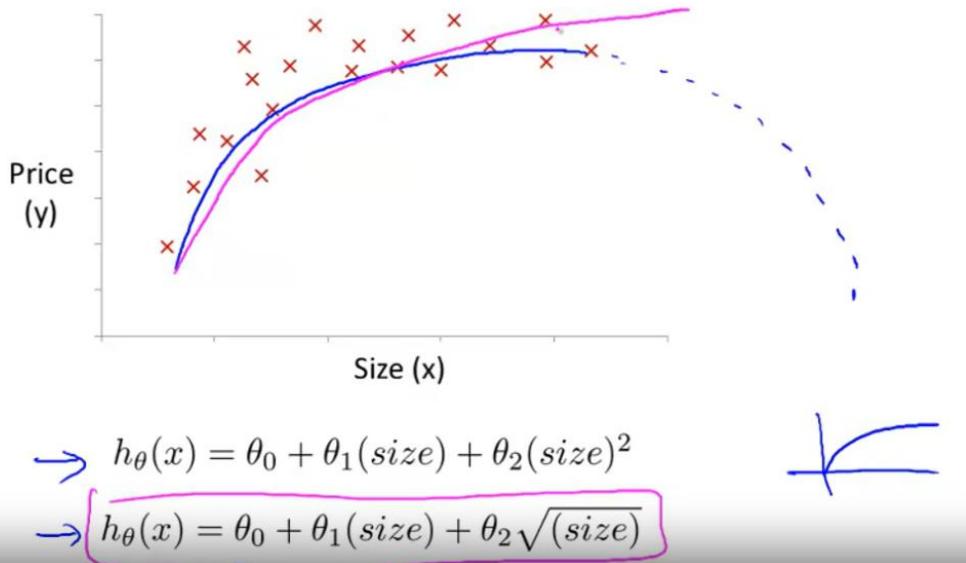


Andrew Ng

We first fit a degree 2 equation line through the data(represented by blue) , but this curve starts going down from point A and hence it is wrong.

We then fit a degree 3 equation line through the data (represented by green)

Choice of features



We can also fit a degree $\frac{1}{2}$ equation (represented by pink).

Normal Equation:

Gradient descent gives one way of minimizing J . Let's discuss a second way of doing so, this time performing the minimization explicitly and without resorting to an iterative algorithm. In the "Normal Equation" method, we will minimize J by explicitly taking its derivatives with respect to the θ_j 's, and setting them to zero. This allows us to find the optimum theta without iteration. The normal equation formula is given below:

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

Examples: $m = 4$.

x_0	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$ $y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$
 $m \times (n+1)$ m-dimensional vector
 $\theta = (X^T X)^{-1} X^T y$

There is **no need** to do feature scaling with the normal equation.

Note: X^T is a $(n+1) * m$ matrix, while X is a $m * (n+1)$ matrix. Thus $X^T * X$ is a $(n+1) * (n+1)$ matrix. . Thus $(X^T * X)^{-1}$ is also a $(n+1) * (n+1)$ matrix. $(X^T * X)^{-1} * X^T$ is a vector of size $(n+1) * m$. y is a vector of size m or

a $m \times 1$ matrix and hence $(X^T X)^{-1} * X^T y$ is a vector of size $(n+1) \times 1$, which also happens to be the dimensions of theta. Where $\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_n)$.

From the above formula we obtain values of all respective thetas $(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$ for which the cost function is minimum.

The following is a comparison of gradient descent and the normal equation:

Gradient Descent	Normal Equation
Need to choose alpha	No need to choose alpha
Needs many iterations	No need to iterate
$O(kn^2)$	$O(n^3)$, need to calculate inverse of $X^T X$
Works well when n is large	Slow if n is very large

With the normal equation, computing the inversion has complexity $O(n^3)$. So if we have a very large number of features, the normal equation will be slow. In practice, when n exceeds 10,000 it might be a good time to go from a normal solution to an iterative process.

The normal equation gets slow with many samples due to the reason that it takes time for performing computation of matrix related operations for large values of n.

Normal Equation Noninvertibility

When implementing the normal equation in octave we want to use the 'pinv' function rather than 'inv.' The 'pinv' function will give you a value of θ even if $X^T X$ is not invertible.

If $X^T X$ is **noninvertible**, the common causes might be having :

- Redundant features, where two features are very closely related (i.e. they are linearly dependent)
- Too many features (e.g. $m \leq n$). In this case, delete some features or use "regularization" (to be explained in a later lesson).

Solutions to the above problems include deleting a feature that is linearly dependent with another or deleting one or more features when there are too many features.

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

- **Redundant features (linearly dependent).**

E.g. $x_1 = \text{size in feet}^2$ $l_m = 3.28 \text{ feet}$
 $x_2 = \text{size in m}^2$
 $x_1 = (3.28)^2 x_2$ $\rightarrow m = 10$
 $\rightarrow n = 100$

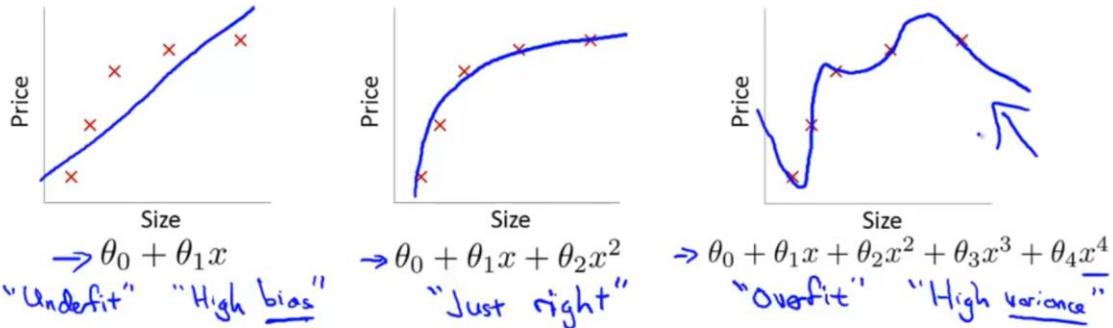
- **Too many features (e.g. $m \leq n$).**

- **Delete some features, or use regularization.**

\downarrow late

Overfitting:

Example: Linear regression (housing prices)



Overfitting: If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \approx 0$), but fail to generalize to new examples (predict prices on new examples).

Andrew Ng

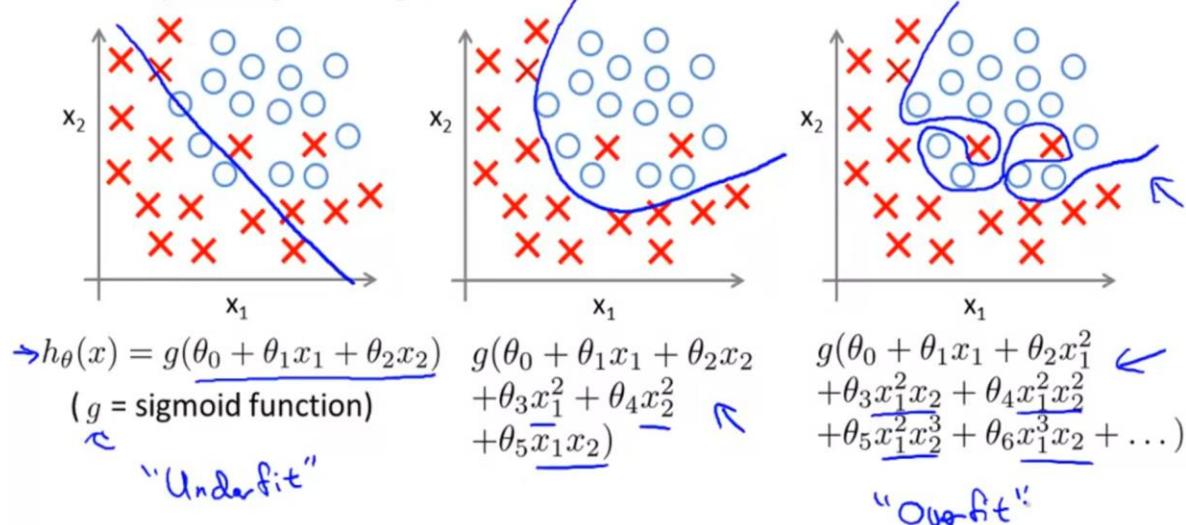
Consider the problem of predicting y from $x \in \mathbb{R}$. The leftmost figure below shows the result of fitting a $y = \theta_0 + \theta_1 x$ to a dataset. We see that the data doesn't really lie on straight line, and so the fit is not very good.

Instead, if we had added an extra feature x^2 , and fit $y = \theta_0 + \theta_1 x + \theta_2 x^2$, then we obtain a slightly better fit to the data (See middle figure). Naively, it might seem that the more features we add, the better. However, there is also a danger in adding too many features: The rightmost figure is the result of fitting a 5th order polynomial $y = \sum_{j=0}^5 \theta_j x^j$. We see that even though the fitted curve passes through the data perfectly, we would not expect this to be a very good predictor of, say, housing prices (y) for different living areas (x). Without formally defining what these terms mean, we'll say the figure on the left shows an instance of **underfitting**—in which the data clearly shows structure not captured by the model—and the figure on the right is an example of **overfitting**.

Underfitting, or high bias, is when the form of our hypothesis function h maps poorly to the trend of the data. It is usually caused by a function that is too simple or uses too few features. At the other extreme, overfitting, or high variance, is caused by a hypothesis function that fits the available data but does not generalize well to predict new data. It is usually caused by a complicated function that creates a lot of unnecessary curves and angles unrelated to the data.

Above terminology is applied to both linear and logistic regression

Example: Logistic regression



and, being unlikely to generalize well to new examples.

Andrew Ng

There are two main options to address the issue of overfitting:

1) Reduce the number of features:

- Manually select which features to keep.
- Use a model selection algorithm (studied later in the course).

2) Regularization

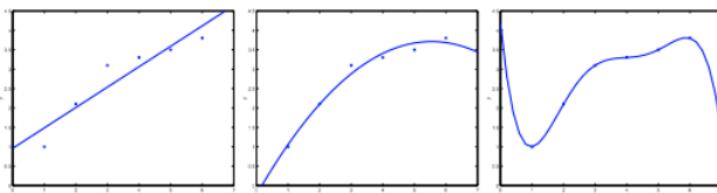
- Keep all the features, but reduce the magnitude of parameters θ_j .
- Regularization works well when we have a lot of slightly useful features.

Regularisation:

When our model is underfitting the data i.e. not able to predict the trend there in the training data properly, we may think of adding more parameters in it to increase the accuracy of prediction. Although this may seem to solve the problem , we have to realise that this can lead to overfitting that is , our model can adjust itself according to the loopholes there in the training data and hence predicts the training data very accurately , but can't do the same on new data. In the example given as we introduce x^3 and x^4 term in our equation $y= \theta_0 + \theta_1 *x + \theta_2 x^2$, though the predictions become more accurate , the problem of overfitting also arises.

The Problem of Overfitting

Consider the problem of predicting y from $x \in \mathbb{R}$. The leftmost figure below shows the result of fitting a $y = \theta_0 + \theta_1 x$ to a dataset. We see that the data doesn't really lie on straight line, and so the fit is not very good.



Instead, if we had added an extra feature x^2 , and fit $y = \theta_0 + \theta_1 x + \theta_2 x^2$, then we obtain a slightly better fit to the data (See middle figure). Naively, it might seem that the more features we add, the better. However, there is also a danger in adding too many features: The rightmost figure is the result of fitting a 5th order polynomial $y = \sum_{j=0}^5 \theta_j x^j$. We see that even though the fitted curve passes through the data perfectly, we would not expect this to be a very good predictor of, say, housing prices (y) for different living areas (x). Without formally defining what these terms mean, we'll say the figure on the left shows an instance of **underfitting**—in which the data clearly shows structure not captured by the model—and the figure on the right is an example of **overfitting**.

Underfitting, or high bias, is when the form of our hypothesis function h maps poorly to the trend of the data. It is usually caused by a function that is too simple or uses too few features. At the other extreme, overfitting, or high variance, is caused by a hypothesis function that fits the available data but does not generalize well to predict new data. It is usually caused by a complicated function that creates a lot of unnecessary curves and angles unrelated to the data.

So what we do is that we introduce a regularisation term in our cost function equation. What this term does is that ‘keeps a track’ on the new features we add in the equation for better fitting of the model. The way it keeps a track is by increasing the value of the cost function if the theta values corresponding to the features/regressors is large. As we want our cost function to be less, the values of theta corresponding to regressors are kept balanced. This also prevents assignment of unnecessary weight to any of the regressors.

##Some screenshots from an article on regularisation:

Find the article here: <https://towardsdatascience.com/regularization-in-machine-learning-76441ddcf99a>

Regularization

This is a form of regression, that constrains/ regularizes or shrinks the coefficient estimates towards zero. In other words, *this technique discourages learning a more complex or flexible model, so as to avoid the risk of overfitting.*

A simple relation for linear regression looks like this. Here Y represents the learned relation and β represents the coefficient estimates for different variables or predictors(X).

$$Y \approx \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

The fitting procedure involves a loss function, known as residual sum of squares or RSS. The coefficients are chosen, such that they minimize this loss function.

$$\text{RSS} = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 .$$

Now, this will adjust the coefficients based on your training data. *If there is noise in the training data, then the estimated coefficients won't*

generalize well to the future data. This is where regularization comes in and shrinks or regularizes these learned estimates towards zero.

Ridge Regression

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 = \text{RSS} + \lambda \sum_{j=1}^p \beta_j^2$$

Above image shows ridge regression, where the **RSS is modified by adding the shrinkage quantity**. Now, the coefficients are estimated by minimizing this function. Here, **λ is the tuning parameter that decides how much we want to penalize the flexibility of our model**. The increase in flexibility of a model is represented by increase in its coefficients, and if we want to minimize the above function, then these coefficients need to be small. This is how the Ridge regression technique prevents coefficients from rising too high. Also, notice that we shrink the estimated association of each variable with the response, except the intercept β_0 , This intercept is a measure of the mean value of the response when $x_{i1} = x_{i2} = \dots = x_{ip} = 0$.

When $\lambda = 0$, the penalty term has no effect, and the estimates produced by ridge regression will be equal to least squares. However, **as $\lambda \rightarrow \infty$, the impact of the shrinkage penalty grows, and the ridge regression coefficient estimates will approach zero**. As can be seen, selecting a good value of λ is critical. Cross validation comes in handy for this purpose. The coefficient estimates produced by this method are **also known as the L2 norm**.

The coefficients that are produced by the standard least squares method are scale equivariant, i.e. if we multiply each input by c then the corresponding coefficients are scaled by a factor of $1/c$. Therefore, regardless of how the predictor is scaled, the multiplication of predictor and coefficient ($X_j \beta_j$) remains the same. **However, this is not the case with ridge regression, and therefore, we need to standardize the predictors or bring the predictors to the same scale before performing ridge regression.** The formula used to do this is given below.

$$\tilde{x}_{ij} = \frac{x_{ij}}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}},$$

NOTES FROM COURSERA:

If we have overfitting from our hypothesis function, we can reduce the weight that some of the terms in our function carry by increasing their cost.

Say we wanted to make the following function more quadratic:

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

We'll want to eliminate the influence of $\theta_3 x^3$ and $\theta_4 x^4$. Without actually getting rid of these features or changing the form of our hypothesis, we can instead modify our **cost function**:

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \cdot \theta_3^2 + 1000 \cdot \theta_4^2$$

We've added two extra terms at the end to inflate the cost of θ_3 and θ_4 . Now, in order for the cost function to get close to zero, we will have to reduce the values of θ_3 and θ_4 to near zero. This will in turn greatly reduce the values of $\theta_3 x^3$ and $\theta_4 x^4$ in our hypothesis function. As a result, we see that the new hypothesis (depicted by the pink curve) looks like a quadratic function but fits the data better due to the extra small terms $\theta_3 x^3$ and $\theta_4 x^4$.

If we have overfitting from our hypothesis function, we can reduce the weight that some of the terms in our function carry by increasing their cost.

Say we wanted to make the following function more quadratic:

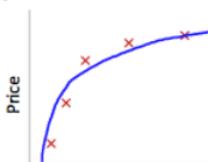
$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

We'll want to eliminate the influence of $\theta_3 x^3$ and $\theta_4 x^4$. Without actually getting rid of these features or changing the form of our hypothesis, we can instead modify our **cost function**:

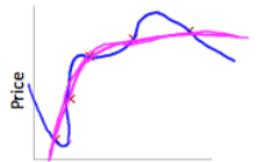
$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \cdot \theta_3^2 + 1000 \cdot \theta_4^2$$

We've added two extra terms at the end to inflate the cost of θ_3 and θ_4 . Now, in order for the cost function to get close to zero, we will have to reduce the values of θ_3 and θ_4 to near zero. This will in turn greatly reduce the values of $\theta_3 x^3$ and $\theta_4 x^4$ in our hypothesis function. As a result, we see that the new hypothesis (depicted by the pink curve) looks like a quadratic function but fits the data better due to the extra small terms $\theta_3 x^3$ and $\theta_4 x^4$.

Intuition



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



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

Suppose we penalize and make θ_3, θ_4 really small.

$$\rightarrow \min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \underline{\theta_3^2} + 1000 \underline{\theta_4^2}$$

$\underline{\theta_3 \approx 0} \quad \underline{\theta_4 \approx 0}$

We could also regularize all of our theta parameters in a single summation as:

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

The λ , or lambda, is the **regularization parameter**. It determines how much the costs of our theta parameters are inflated.

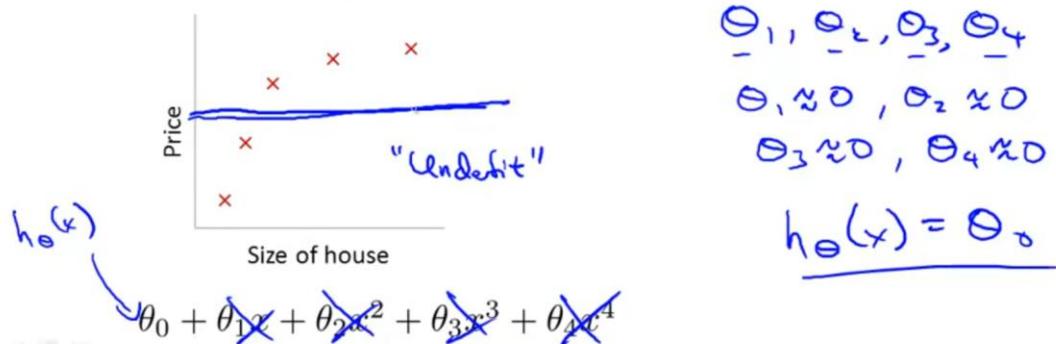
Using the above cost function with the extra summation, we can smooth the output of our hypothesis function to reduce overfitting. If lambda is chosen to be too large, it may smooth out the function too much and cause underfitting. Hence, what would happen if $\lambda = 0$ or is too small?

When lambda is too large we obtain a trend line parallel to the x axis. This is because to minimize the cost function $\theta_1, \theta_2, \dots, \theta_n$ will assume values that tend to 0 and hence our model will become $y = \theta_0 + 0$, which explains the horizontal line parallel to x axis.

In regularized linear regression, we choose θ to minimize

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

What if λ is set to an extremely large value (perhaps for too large for our problem, say $\lambda = 10^{10}$)?



Regularisation and linear regression:

How the formula changes on adding the regularisation term:

Gradient descent

$$\frac{\partial J(\theta)}{\partial \theta_0} \quad \theta_0, \theta_1, \dots, \theta_n$$

Repeat {

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\rightarrow \theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad (j = 1, 2, 3, \dots, n)$$

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

$1 - \alpha \frac{\lambda}{m} < 1$ $\theta_j \times 0.99$

Andrew Ng

Gradient Descent

We will modify our gradient descent function to separate out θ_0 from the rest of the parameters because we do not want to penalize θ_0 .

```
Repeat {
     $\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$ 
     $\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \quad j \in \{1, 2, \dots, n\}$ 
}
```

The term $\frac{\lambda}{m} \theta_j$ performs our regularization. With some manipulation our update rule can also be represented as:

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

The first term in the above equation, $1 - \alpha \frac{\lambda}{m}$ will always be less than 1. Intuitively you can see it as reducing the value of θ_j by some amount on every update. Notice that the second term is now exactly the same as it was before.

How the normal equation formula changes:

Normal equation

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \quad \text{m} \times (n+1)$$

$$\rightarrow \min_{\theta} J(\theta) \quad \frac{\partial}{\partial \theta_j} J(\theta) \stackrel{\text{set } 0}{=} 0 \quad \text{implies}$$

$$\rightarrow \theta = (X^T X + \lambda \underbrace{\begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}}_{(n+1) \times (n+1)})^{-1} X^T y$$

E.g. n=2 $\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (n+1) \times (n+1)$

Andrew Ng

Note: If $m \leq n$, then the matrix $(X^T X)$ becomes non-invertible. The regularisation term λ *the matrix shown in the figure below make this $(X^T X)$ invertible even when $m \leq n$. So adding this regularisation has one more advantage other than regularisation itself.

Non-invertibility (optional/advanced).

Suppose $m \leq n$, \leftarrow
 (#examples) (#features)

$$\theta = \underbrace{(X^T X)^{-1} X^T y}_{\text{non-invertible / singular}} \quad \text{pinv} \quad \frac{\text{inv}}{\lambda}$$

If $\lambda > 0$,

$$\theta = \left(X^T X + \lambda \begin{bmatrix} 0 & 1 & & & \\ 1 & 0 & 1 & & \\ & 1 & 0 & \ddots & \\ & & \ddots & \ddots & 1 \end{bmatrix} \right)^{-1} X^T y$$

invertible.



Coursera notes:

Normal Equation

Now let's approach regularization using the alternate method of the non-iterative normal equation.

To add in regularization, the equation is the same as our original, except that we add another term inside the parentheses:

$$\theta = X^T X + \lambda \cdot L^{-1} X^T y$$
$$0$$
$$1$$

where $L =$

$$\begin{matrix} & & 1 \\ & \ddots & \\ 1 & & \end{matrix}$$

L is a matrix with 0 at the top left and 1's down the diagonal, with 0's everywhere else. It should have dimension $(n+1) \times (n+1)$. Intuitively, this is the identity matrix (though we are not including x_0), multiplied with a single real number λ .

Recall that if $m < n$, then $X^T X$ is non-invertible. However, when we add the term $\lambda \cdot L$, then $X^T X + \lambda \cdot L$ becomes invertible.

Question:

Gradi
Repe
→
Suppose you are doing gradient descent on a training set of $m > 0$ examples, using a fairly small learning rate $\alpha > 0$ and some regularization parameter $\lambda > 0$. Consider the update rule:
$$\theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}.$$

Which of the following statements about the term $\left(1 - \alpha \frac{\lambda}{m}\right)$ must be true?

$1 - \alpha \frac{\lambda}{m} > 1$
 $1 - \alpha \frac{\lambda}{m} = 1$
 $1 - \alpha \frac{\lambda}{m} < 1$
 None of the above.

Correct

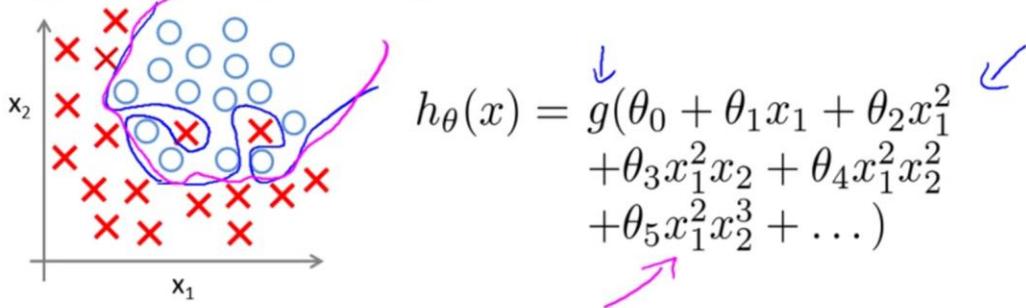
Continue

Andrew Ng

Regularisation and logistic regression:

Even when we fit high order terms like $x_1^2 * x_2^3$ in our equation , we don't overfit the model since due to regularisation , the coefficients like θ_5 corresponding to these terms decrease their magnitude and hence adjust their effect. The blue line shows an instance of overfitting and the line in pink is the regularized equation.

Regularized logistic regression.



Cost function:

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

| $\boxed{\theta_1, \theta_2, \dots, \theta_n}$

Andrew Ng

Changes in the formula :

Gradient descent

Repeat {

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\rightarrow \theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \leftarrow$$

$(j = \boxed{1, 2, 3, \dots, n})$
 $\theta_1, \dots, \theta_n$

}

$$\frac{\partial}{\partial \theta_j} J(\theta)$$

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

a previous slide that does use regularization.

Andrew Ng

The fminunc function and regularisation in octave:

Advanced optimization

```

fminunc (e costFunction)
θ = [θ₀ θ₁ ... θₙ] ← theta(1) ←
function [jVal, gradient] = costFunction(theta) ← theta(2) ←
    θ₀ ← theta(n+1) ←
    jVal = [ code to compute J(θ)]; ←
    → J(θ) =  $-\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))$  ←
    → gradient(1) = [ code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$ ]; ←
     $\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$  ←
    → gradient(2) = [ code to compute  $\frac{\partial}{\partial \theta_1} J(\theta)$ ]; ←
     $\left( \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)} \right) + \frac{\lambda}{m} \theta_1$  ←
    → gradient(3) = [ code to compute  $\frac{\partial}{\partial \theta_2} J(\theta)$ ]; ←
    : ←
     $\left( \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)} \right) + \frac{\lambda}{m} \theta_2$  ←
    gradient(n+1) = [ code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$ ];

```

Andrew Ng

Question:

Advanced Optimization

When using regularized logistic regression, which of these is the best way to monitor whether gradient descent is working correctly?

- Plot $[-\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))]$ as a function of the number of iterations and make sure it's decreasing.
- Plot $[-\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) - \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2]$ as a function of the number of iterations and make sure it's decreasing.
- Plot $[-\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2]$ as a function of the number of iterations and make sure it's decreasing.

Correct

- Plot $\sum_{j=1}^n \theta_j^2$ as a function of the number of iterations and make sure it's decreasing.

Continue

Andrew Ng



Quiz 1

Graded Quiz • 10 min

Due Mar 23, 12:29 PM IST

Quiz 1

TOTAL POINTS 5

1. Which of the following are components in building a machine learning algorithm?

1 point

- Statistical inference
- Machine learning
- Artificial intelligence
- Training and test sets



Collecting data to answer the question.





Quiz 1

Graded Quiz • 10 min

Due Mar 23, 12:29 PM IST



Connecting data to answer the question.



2. Suppose we build a prediction algorithm on a data set and it is 100% accurate on that data set. Why might the algorithm not work well if we collect a new data set?

1 point

- We have used neural networks which has notoriously bad performance.
- We may be using bad variables that don't explain the outcome.v
- We have too few predictors to get good out of sample accuracy.
- Our algorithm may be overfitting the training data, predicting both the signal and the noise.



3. What are typical sizes for the training and test sets?

1 point



 Quiz 1 Due Mar 23, 12:29 PM IST
Graded Quiz • 10 min

3. What are typical sizes for the training and test sets?

1 / 1 point

- 60% in the training set, 40% in the testing set.
 - 90% training set, 10% test set
 - 20% training set, 80% test set.
 - 100% training set, 0% test set.

✓ Correct

4. What are some common error rates for predicting binary variables (i.e. variables with two possible values like yes/no, true/false, normal/abnormal, clicked/didn't click)? Check the correct

1 / 1 point



 Quiz 1
Graded Quiz • 10 min

✓ Correct

4. What are some common error rates for predicting binary variables (i.e. variables with two possible values like yes/no, disease/normal, clicked/didn't click)? Check the correct answer(s).

1 / 1 point

- Median absolute deviation
 - R²
 - Root mean squared error
 - Correlation
 - Predictive value of a positive



✓ Correct



Quiz 1

Graded Quiz • 10 min

Due Mar 23, 12:29 PM IST

Correct

5. Suppose that we have created a machine learning algorithm that predicts whether a link will be clicked with 99% sensitivity and 99% specificity. The rate the link is clicked is 1/1000 of visits to a website. If we predict the link will be clicked on a specific visit, what is the probability it will actually be clicked?

1 / 1 point

- 99.9%
- 50%
- 99%
- 9%



Correct