

①

Supervised Learning:

- We are given a set of input and trying to predict what our output output should look like, having the idea that there is a relationship between the input and the output.
- Supervised learning problem can be split into "regression" and "classification" problems.
 - Regression; In regression problem, we are trying to predict output with a continuous output, meaning that we are trying to map input variable to some continuous function..
 - Classification; In classification problem, we are instead trying to predict results in a discrete output. In other words, we're trying to map input variable into discrete categories.

Ex 1.

Given data about the size of house on the real estate market, try to predict their price. Price is a function of size is a continuous output, so this is a regression problem.

We could turn this example into a classification problem by instead making our output about whether the house "sells for more or less than the asking price". Here we are classifying the house based on price into two discrete categories.

Ex 2.

- (a) Regression - Given a picture of male/female, we have to predict his/her age on the basis of picture.
 - (b) Classification - Given a picture of male/female, we've to predict whether he/she is at high school, college, graduation.
- Another Example for classification - Banker has to decide whether or not to give a loan to someone on the basis of her credit history.

Unsupervised Learning:

Unsupervised learning, on the other hand, allows us to approach problems with little or no idea what our results should look like. We can derive structure from data when we don't necessarily know the effect of the variables.

We can derive this structure by deriving the data are based on relationships among the variables in the data.

With unsupervised learning there is no feedback based on the prediction result, i.e., there is no teacher to reward you.

Ex:

- Clustering: Take a collection of 1000 essays written on the US Economy, and find a way to automatically group these essays into a small number that are somewhat similar or related by different variables, such as word frequency, sentence length, page count, and so on.

Non-clustering: The

- Non-clustering: The "cocktail party algorithm", which can find structure in messy data (such as the identification of individual voices and music from a room with ^{→ noise} of sounds at a cocktail party).

(5)

[ML: Linear Regression With One Variable]

* Model Refreshation:

Recall that in regression problems, we are taking input variables and trying to fit the output onto a continuous expected result function.

Linear regression with one variable is also known as "Univariate linear regression".

Univariate linear regression is used when you want to predict a single output value y from a single input value x . We're doing Supervised learning here, so that means we already have an idea about what the input/output curve and effect should be.

* The Hypothesis Function

Our hypothesis function has the general form:

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

Note that this is like the equation of a straight line. We give $h_{\theta}(x)$ values for θ_0 and θ_1 to get our estimated output \hat{y} . In other words, we're trying to create a function called h_{θ} that is trying to map our input data (the x 's) to our output data (the y 's).

Example:

Suppose we've the following set of training data:

Input x	Output y
0	4
1	1
2	1
3	8

Now we can make some random guess about our hypothesis function: $\theta_0 = 2$ and $\theta_1 = 2$. The hypothesis function becomes $h_{\theta}(x) = 2 + 2x$.

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

So for input of 1 to our hypothesis, \hat{y} will be 4 . This is off by 3 . Note that we will be trying out various values of θ_0 and θ_1 to try to find values which provide the best possible "fit" or the most representative "straight line" through the data points mapped on the X - y plane.

* Cost Function:

We can measure the accuracy of our hypothesis function by using a Cost function. This takes an average (actually a scalar division of an average) of all the results of the hypothesis with inputs from x_i 's compared to the actual output y_i 's.

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

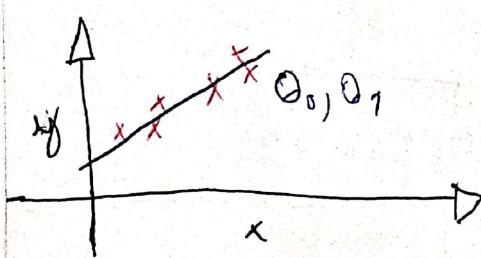
To break it apart, it is $\frac{1}{2} \bar{x}$ where \bar{x} is the mean of the square of $h_{\theta}(x_i) - y_i$, or the difference between the predicted value and the actual value.

This function is often times called the "Squred error function", or "Non squared error". The mean is halved ($\frac{1}{2m}$) as a compensation for the computation of the gradient descent, or the derivative term of the square function will cancel out the $\frac{1}{2}$ term.

Now we are able to correctly measure the accuracy of our prediction function against the correct results. What is that we can predict these results we don't know.

If we try to think of it in visual terms, our training dataset is scattered on the x - y plane. We are trying to find a straight line (defined by $h_\theta(x)$) which passes through this scattered set of data. Our objective is to get the best possible line. The best possible line will be such so that the sum of average squared vertical distances of the scattered points from the line will be the least. In the best case, the line should pass through all points of our training dataset. In such a case the value of $J(\theta_0, \theta_1)$ will be 0.

Video: Cost function: 4:54



$$\text{Training Example} \rightarrow \min_{\theta_0, \theta_1} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

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

Minimizing Training Example Expression

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

minimizing $J(\theta_0, \theta_1)$

θ_0, θ_1

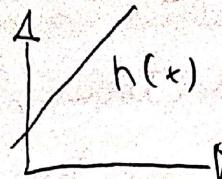
Cost function

Squared error cost function.

Videos: (a) function: $\theta_0 + \theta_1 x$

Input variable:

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



Parameter:

$$\theta_0, \theta_1$$

(a) function:

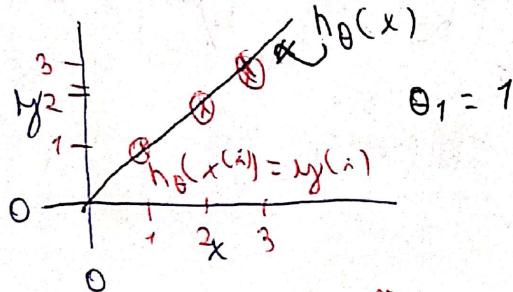
$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \text{ minimizing } J(\theta)$$

Goal: Minimizing $J(\theta_0, \theta_1)$

$$\theta_0, \theta_1$$

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

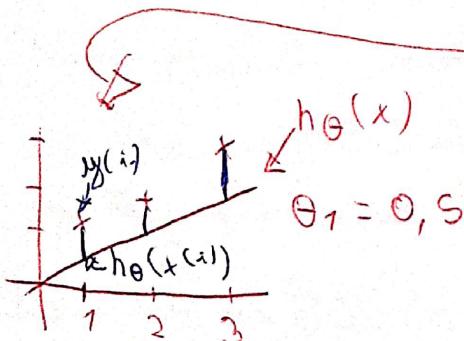
(for fixed θ_1 , this is a function of x)



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

$$= \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2 = \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0^2$$

$$J(1) = 0$$



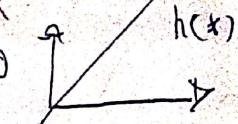
$$J(0.5) = \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

$$= \frac{1}{2 \times 3} (3.5) = \frac{3.5}{6} \approx 0.58$$

Simplified:

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

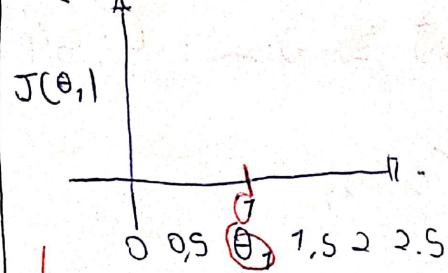
$$\theta_0 = 0$$



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

$$J(\theta_1)$$

(function of the parameter θ_1)



$$\theta_1 = 0.5?$$

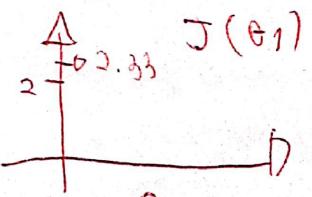
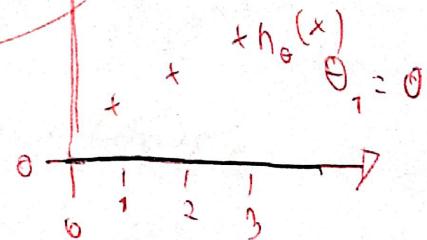
Werte:

$$J(0) ? \underline{m=3}$$

$$(0-1)^2 + (0-2)^2 +$$

$$(0-3)^2 = 14$$

$$\frac{1}{2 \cdot 3} = \frac{1}{6} \cdot 14 = \frac{14}{6}$$

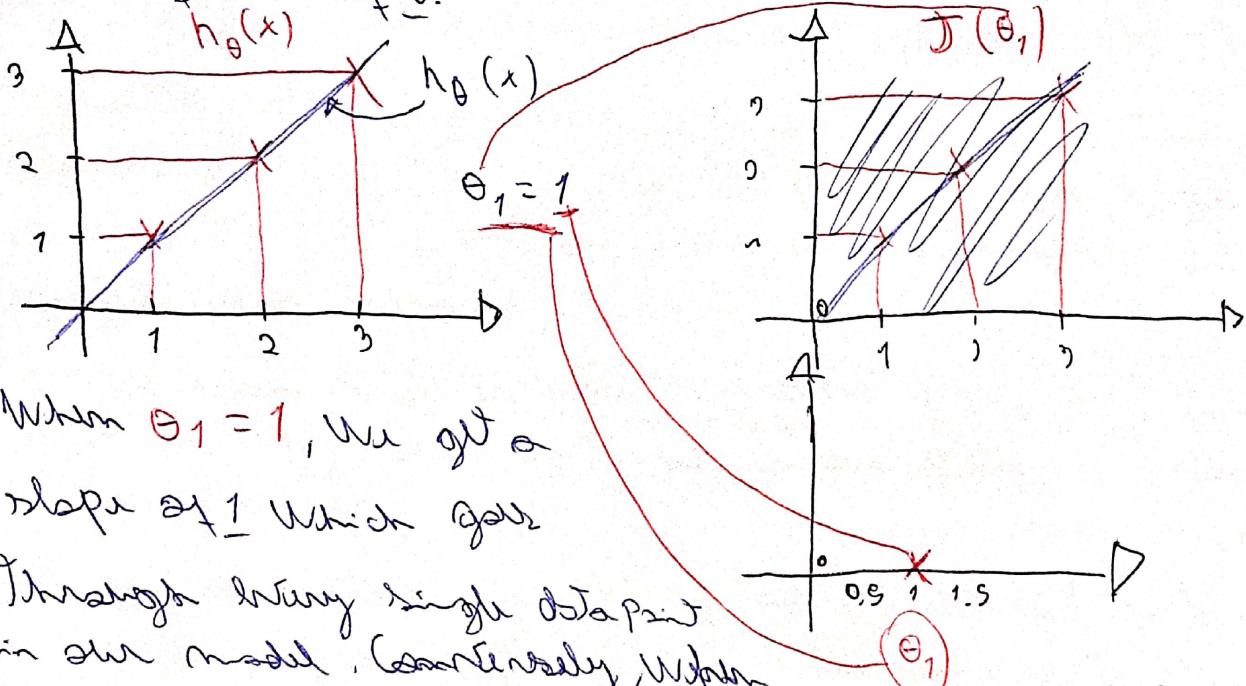


!! Vericles (or) Since ... I.1

9:40

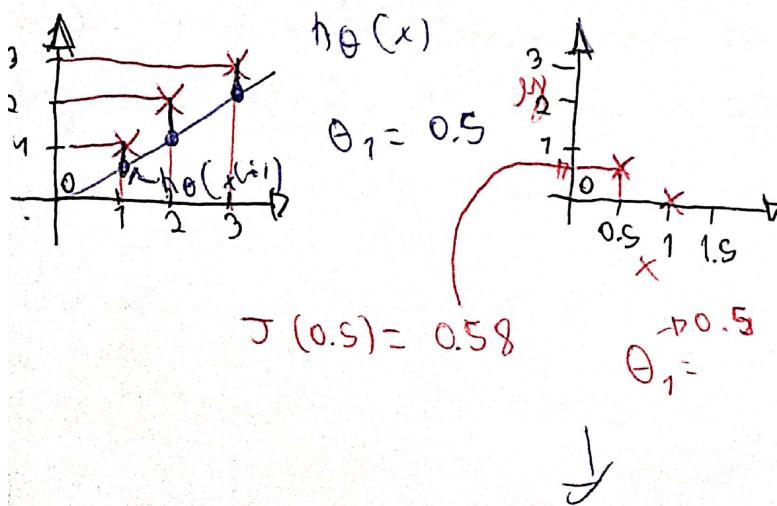
(4) Linear - Solution I

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 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 graphs show the ideal situation where we've a cost function of 0.

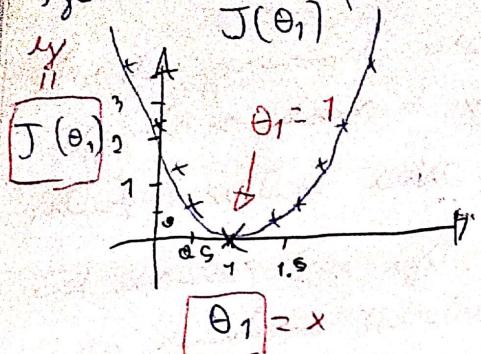


When $\theta_1 = 1$, we get a slope of 1 which goes through every single data point in our model. Conversely, when $\theta_1 = 0.5$, we see the vertical distance from our fit to the data points increase.

$$J(1) = 0$$



This is our cost function -
yields the following graph:



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

- (b) Function - Function II:

ML: Gradient Descent

So we have our hypothesis function and we have a way of measuring how well it fits into the data. Now we need to estimate the parameters in hypothesis function. That's where gradient descent comes in.

Imagine that we graph our hypothesis function based on its θ_0 and θ_1 (actually we are graphing the cost function as a function of the parameter estimator). This can be kind of confusing; we are moving up to a higher level of abstraction. We are not graphing x and y itself, but \rightarrow the parameter range of our hypothesis function and the cost resulting from selecting particular θ_0 & θ_1 parameters.

We put θ_0 on the x axis and θ_1 on the y axis, with the cost function on the vertical z axis. The points on our graph will be the result of the cost function using our hypothesis with those specific theta parameters.

We will know that we have succeeded when our cost function is at the very bottom of the pit in our graph, i.e. when its value is the minimum.

The way we do this is by taking the derivative (the longitudinal line is a function) of our cost function. The slope of the tangent is the derivative of that point and it will give us a direction to move to. We must step down the cost function in the direction with the ~~steepest~~ ^{negative} descent, and the size of each step is determined by the parameter α , which is called the learning rate.

The gradient descent algorithm is:

Repeat until convergence:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

Where

$J = 0, j$ represents the feature index number.

Intuitively, this could be thought of as:

Repeat until convergence:

$$\theta_j := \theta_j - \alpha [\text{Slope at } j\text{th feature dimension}] \quad [\quad]$$

- Gradient Descent for Linear Regression:

When specifically applied to the case of Linear Regression, a more form of gradient descent equation can be derived. We can substitute our cost function and our actual hypothesis function and modify the equation to (The derivation of the formulae are out of the scope of this course, but a really good one can be found here):

Repeat until convergence:

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m ((\hat{y}_i - y_i) x_i)$$

}

Where m is the size of the training set, θ_0 a constant that will be changing simultaneously with θ_1 , and x_i, y_i are values of the given training set (data).

Note that we have separated out the bias cost for θ_0 into separate equations for θ_0 and θ_1 ; and that for θ_1 we are multiplying x_i at the end due to the dot product.

The point of all this is that if we start with a guess for our hypothesis and then repeatedly apply these gradient descent equations, our hypothesis will become more and more accurate.

- Gradient Descent:

(6)

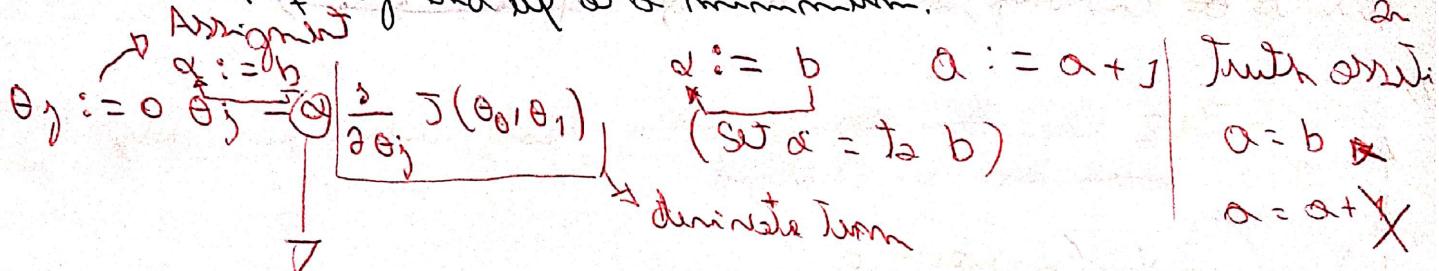
Hausse Kostenfunktion $J(\theta_0, \theta_1)$

Want min $J(\theta_0, \theta_1)$
 θ_0, θ_1

⇒ Outline:

- Start With same θ_0, θ_1 (say $\theta_0 = 0, \theta_1 = 0$)

- Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$
Until we hopefully end up at a minimum.



Simultaneously Learning Rate
Update θ_0 and θ_1 !

$$\text{Temp}_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{Temp}_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

~~$\text{Temp}_0 := \text{Temp}_0$~~

~~$\theta_1 := \text{Temp}_1$~~

True:

$$\theta_0 = 1$$

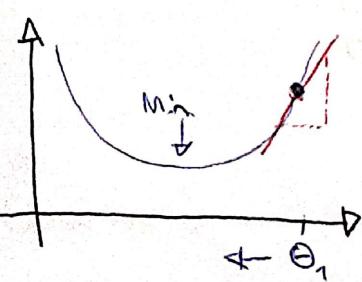
$$\theta_1 = 2$$

$$\theta_j := \theta_j + \sqrt{\theta_{j+1}} \quad (\text{for } J=0)$$

and $\gamma = 1$

Gradient Descent & Iteration:

$$\min_{\theta_1} J(\theta_1) \quad \theta_1 \in \mathbb{R}$$



$$J(\theta_1) \quad (\theta_1 \in \mathbb{R})$$

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

$$\theta_1 := \theta_1 - \underline{\alpha} \cdot (\text{Periode number})$$

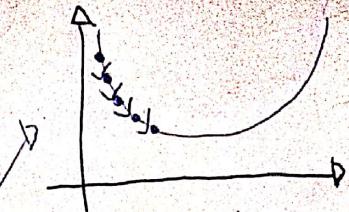


$$\frac{\alpha}{\partial \theta_1} J(\theta_1)$$

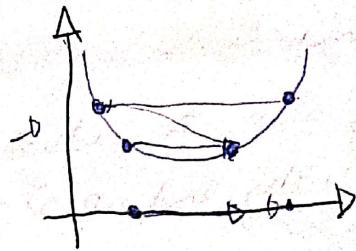
$$\theta_1 := \theta_1 - \alpha \cdot (\text{Nagende Number})$$

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

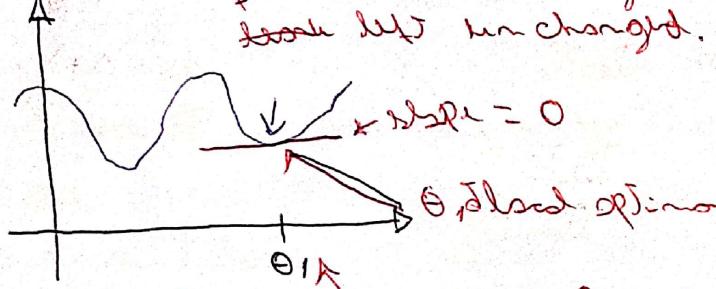
If α is too small, gradient descent can be slow.



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



If I don't stay at the best minimum my θ_1 will be stuck like this.



Current value of θ_1

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

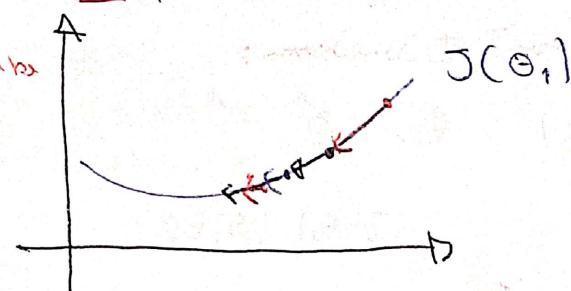
$$\theta_1^* := \theta_1 - \alpha \cdot 0$$

$$\theta_1 := \theta_1$$

Gradient descent can converge to a best minimum when with the learning rate α fixed.

$$\theta_1 := \theta_1 - \alpha \boxed{\frac{\partial}{\partial \theta_1} J(\theta_1)} \quad \begin{array}{l} \text{fixing} \\ \text{time with} \\ \alpha \rightarrow 0 \end{array}$$

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



Gradient Descent for Linear Regression:

Gradient descent algorithm of Linear Regression Model

Repeat until convergence:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

$$(3m \geq j \text{ and } j = 0)$$

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

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

$$\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

}

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

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

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

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

Gradient descent algorithm

Repeat until convergence:

$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

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

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

Up date
 θ_0 and θ_1 ,
 simultaneously

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

"Convex Function"

Bowl-Shaped

"Batch" Gradient Descent

"Batch" Each step of gradient descent uses all the training examples.

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

(Elliptical)

Then = Invert

Then = Is gen

? Gradient Descent can converge even if α is kept fixed. (But it could be too long, or else it may fail to converge.)

2. For the specific choice of loss function $J(\theta_0, \theta_1)$ used in Linear Regression
There're no local optima (other than the global optimum)

⑧ [ML: Linear Algebra Review]

- Matrices and Vectors

Matrices are 2-dimensional arrays:

$$[a \ b \ c \ d \ e \ f \ g \ h \ i \ j \ k \ l]$$

The above matrix has four rows and three columns, so it is a 4×3 matrix.

A vector is a matrix with one column and many rows:

$$[w \ x \ y \ z]$$

Scalars are a subset of matrices. The scalar vector is a 4×1 matrix

* Notation and Terms:

- A_{ij} refers to the element in the i th row and j th column of matrix A .
- A vector with ' m ' values is referred to as an ' m -dimensional vector.
- v_i refers to the element in the i th row of the vector.
- In general, all our vectors and matrices will be 1-indexed. Note this for some programming languages. The arrays are 0-indexed.
- Matrices are usually denoted by uppercase letters while vectors are lowercase.
- "Scalar" means that an object is a single value, not a vector = matrix.
- \mathbb{R} refers to the set of scalar real numbers
- \mathbb{R}^n refers to the set of n -dimensional vectors of real numbers.

- Addition and Scalar Multiplication
Addition and subtraction are element-wise, we can simply add or subtract each corresponding element:
$$[a \ b \ c \ d] + [m \ x \ y \ z] = [a+m \ b+x \ c+y \ d+z]$$
To add or subtract two matrices, their dimensions must be the same.
In scalar multiplication, we simply multiply every element by the scalar value:
$$[a \ b \ c \ d] * x = [a*x \ b*x \ c*x \ d*x]$$
- Matrix-Vector Multiplication
We map the columns of the vector onto each row of the matrix, multiplying each element and summing the result.
$$[a \ b \ c \ d] * [x \ y] = [a*x + b*x + c*y + d*y]$$
The result is a vector. The vector must be the second term of the multiplication. The number of columns of the matrix must equal the number of rows of the vector.
An $m \times n$ matrix multiplied by an $n \times 1$ vector results in an $m \times 1$ vector.

(9)

Matrix - Matrix Multiplication

We multiply two matrices by breaking it into several scalar multiplications and concatenating the result.

$$[a \ b \ c \ d \ e] * [u \ v \ w \ x \ y \ z] = [a*u + b*v \ a*x + b*z \ c*u + d*v \ c*x + d*z \ e*u + f*v \ e*x + f*z]$$

An $m \times m$ matrix multiplied by an $n \times 0$ matrix results in an $m \times 0$ matrix. In the above example, a 3×2 matrix times a 2×2 matrix resulted in a 3×2 matrix.

To multiply two matrices, the number of columns of the first matrix equal the number of rows of the second matrix.

- Matrix Multiplication Properties

- Not commutative. $A * B \neq B * A$
- Associative. $(A * B) * C = A * (B * C)$

The identity matrix, when multiplied by any matrix of the same dimensions, results in the original matrix. It's like multiplying numbers by 1. The identity matrix simply has 1's on the diagonal (upper left to bottom right diagonal) and 0's elsewhere.

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

When multiplying the identity matrix after some matrix ($A * I$), the square identity matrix should match the other matrix's columns. When multiplying the identity matrix before some other matrix ($I * A$), the square identity matrix should match the other matrix's rows.

- Inverse and Transpose

The inverse of a matrix A is denoted A^{-1} . Multiplying by the inverse results in the identity matrix.

A non square matrix does not have an inverse matrix. We can compute inverses of matrices in octave with the inv(A) function [5] and in matlab with the inv(A) function. Matrices that don't have an inverse are singular or degenerate.

The Transposition of a matrix is like rotating the matrix 90° in clockwise direction and then reflecting it. We can compute Transposition of matrices in matlab with the Transpose(A) function on A.

$$A = \begin{bmatrix} a & b & c & d & e \end{bmatrix}$$

$$A^T = \begin{bmatrix} a & c & b & d & e \end{bmatrix}$$

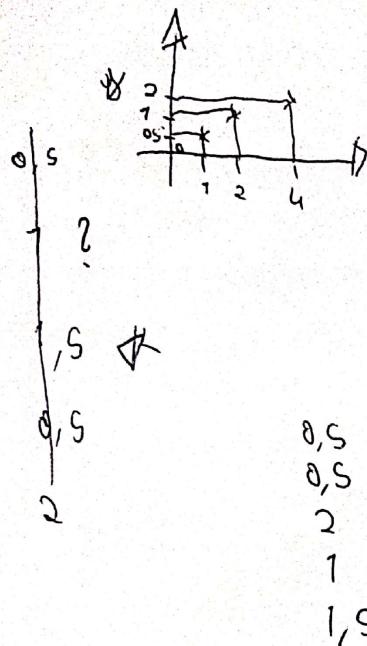
In other words:

$$A_{ij} = A_{ji}^T$$

[Test: Linear Regression With one Variable]

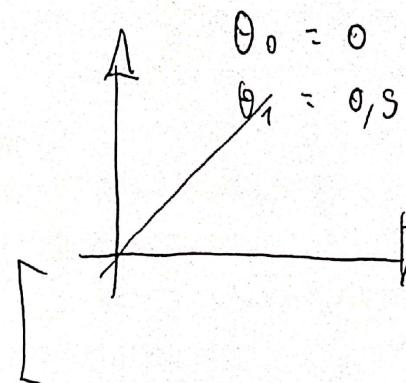
y_{pred}

$$h_{\theta}(x) = -1 + 2 \cdot 6 = 11$$



A

$$\frac{1}{2} \sum$$



$$3) \quad \theta_0 = -2, \theta_1 = 0.5$$

$$2) \quad J(\theta) = \frac{1}{2m} \left[\frac{(1-1)^2}{9} + \frac{(2-1)^2}{1} + \frac{(3-3)^2}{8} + \frac{(4-1)^2}{9} \right]$$

Hausse size:

$$\left\{ \begin{array}{l} 2104 \\ 1416 \\ 1534 \\ 852 \end{array} \right.$$

Module 3 Computing hypothesis:

① $h_0(x) = -40 + 0.25x$

② $h_0(x) = 200 + 0.1x$

3. $h_0(x) = -150 + 0.4x$

Matrix

$$\begin{bmatrix} 1 & 2104 \\ 1 & : \\ 1 & \\ 1 & \end{bmatrix}$$

Matrix Transfer

Matrix

$$\begin{bmatrix} -40 \\ 200 \\ 0.25 \\ 0.7 \end{bmatrix} \quad \begin{bmatrix} -150 \\ 0.4 \end{bmatrix}$$

$$= \begin{bmatrix} 486 \\ 314 \\ 344 \\ 113 \end{bmatrix} \quad \begin{bmatrix} 410 \\ 942 \\ 953 \\ 285 \end{bmatrix} \quad \begin{bmatrix} 892 \\ 418 \\ 464 \\ 191 \end{bmatrix}$$

Prediction

at 1st h_0

Prediction

at 2nd h_0

Ex:

$$A = \begin{bmatrix} 1 & 2 & 0 \\ 3 & 5 & 9 \end{bmatrix}$$

$$A^T = \begin{bmatrix} 1 & 3 \\ 2 & 5 \\ 0 & 9 \end{bmatrix}$$

Note:

-Matrix Vector Mult. 1.44:

Hand Rights:

2704
1416

1534
852

Matrix

$$\begin{bmatrix} 1 & 2704 \\ 1 & 1416 \\ 1 & \end{bmatrix}$$

4x2

$$h_0(x) = -40 + 0.25x$$

2x1

Viter

$$\begin{bmatrix} -40 \\ 0.25 \end{bmatrix}$$

4x1

$$\begin{bmatrix} h_0(1416) \\ h_0(2104) \end{bmatrix}$$

$$\begin{bmatrix} -40 \times 1 + 0.25 \times 2104 \\ -40 \times 1 + 0.25 \times 1416 \end{bmatrix}$$

G7

Code:

Prediction = DataMatrix * Parameters

=

Mult. a matrix by a vector of 1 is equal to sum all the number of the row in the matrix.

// $A \times B = C$

$$\begin{bmatrix} & \\ & \end{bmatrix}$$

x

$$\begin{bmatrix} & \\ & \end{bmatrix}$$

=

$$\begin{bmatrix} & \\ & \end{bmatrix}$$

$\boxed{n} \times \boxed{n}$ matrix
(n rows,
 n columns)

$\boxed{n} \times \boxed{0}$ matrix
(n rows,
 0 columns)

$n \times 0$
matrix