

# Machine Learning

No. \_\_\_\_\_

Date \_\_\_\_\_

## Introduction:

Machine Learning — example: search engine, photo recognition, <sup>(web search, photo tagging, email anti-spam)</sup>  
↳ mimic algorithm to know how human <sup>learn</sup> spam filter in email.

(most recent stages of development)  
⇒ to learn state-of-the-art ML algorithm

⇒ to know how to get this stuff work on problem you care about

ML-grown out of field of AI

↳ ~~new~~ new capability for computers and touch a lot of things in science

Learning algorithm

Example: 1. Database Mining

— Large ~~data~~ datasets from growth of automation/web

e.g. — web <sup>click data</sup> ~~data~~ <sup>stream</sup>, medical record,  
(clickstream data)  
collecting data so to understand user/human better

2. Applications can't program by hand

Eg. — Autonomous helicopter, handwriting recognition, most of Natural Language Processing (NLP), Computer Vision  
we can't write the helicopter program, but write to let computer understand and learn by itself to fly

3. Self-Customising programs

Eg. — Amazon, Netflix product recommendations

4. Understanding human learning (brain, real AI)

What is Machine Learning?

2 definitions

- i) The field of study that gives computer the ability to learn without being explicitly programmed.
- ii) A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .

eg.: Playing checkers

$E$  = experience of playing many games of checkers

$T$  = task of playing checkers

$P$  = The probability that program will win next game.

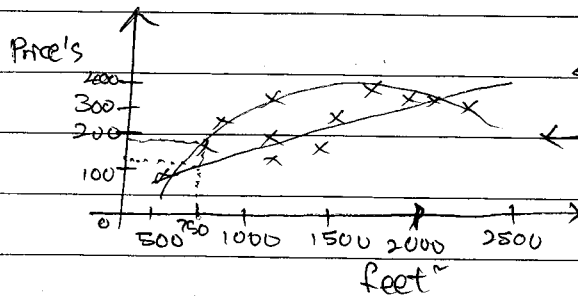
Machine learning problem be assigned to one of two broad classifications;

~~Sup~~ - Supervised learning and Unsupervised learning

# Supervised Learning

Example :

Housing price prediction :



can use a computer algorithm to draw straight line to check the price of 750 feet  
or even make a quadratic function to the data

So, to discuss whether use straight line or quadratic function

Supervised Learning can be said as - give right answer to computer.

↳ give a right data set how much houses are sell

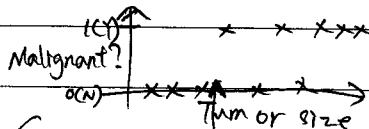
→ task of algorithm is to produce more right answer

in term  $\Rightarrow$  Regression = Predict continuous value output (price)

in this case

Example :

Breast cancer (malignant, benign)



if patient size is here, machine learning question is to estimate the probability / chance for tumor is malignant or benign.

Kind of classification problem

Discrete valued output (0 or 1)

$\Rightarrow$  trying to predict 0 or 1 output.

or there's other type of breast cancer, it will be

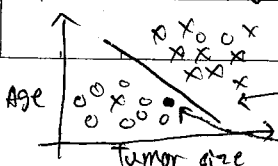
In this problem can have more than 2 values for the two possible value for output. 0, 1, 2, 3  
no first second third type

Another way plot data

Benign use "o" and malignant remain "x"

o o o x o x o x o x  
Tum or size (only 1 attribute)

if 2 attributes



Learning algorithm will draw a straight line to separate 2 classes of algorithm.

if someone falls in here, we will predict most likely is benign.

⇒ An algorithm, Support Vector Machine

Conclusion :

By regression, our goal is to predict a continuous valued output.

一个 class 来 predict  
在训练一个 class

## Supervised learning

regression      classification

we are instead  
trying to  
predict results  
in a discrete  
output. we are  
trying to map input  
variables into discrete categories.

- if we turn this into making our output about whether the house "sells for more or less than asking price" Then we are classifying houses based on price into two discrete categories.

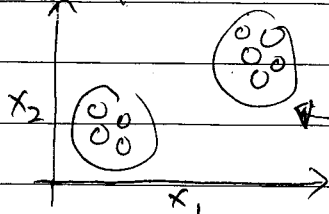
Classification - Given patient with a tumor, we have to predict the tumor is malignant or benign.

## Unsupervised Learning

We are given a data set, and not given any label on the data, and not told what to do and what each data point is.

We just have to find out the structure of dataset.

Example



all same

Given a data set, unsupervised learning algorithm might decide data lives in 2 diff. clusters.

will break these data into 2 separate clusters.

∴ clustering algorithm

example: Google news

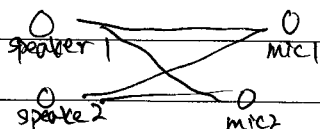
→ 一个标题是 news, online 有很多不同的 stories, google cluster 他们全部在一个大标题是

labelled data is like "spam" and "not spam" 垃圾邮件分类.

Used for other application:

- ⇒ Organize large computer clusters (figure out which machines tend to work together, put them together, so work better)
- social network analysts (identify the cohesive groups of friends)
- Market segmentation (group the customers)
- Astronomical data analysts (how galaxies are formed)

## Cocktail party problem



Both of speakers speak at the same time but microphone not receiving the same time, so algorithm can separate out the 2 voices of speakers by using Octave type programming environment

To approach problems with little or no idea what results should look like, we derive structure from dataset

— by clustering the data based on relationships among variables in data with unsupervised learning, no feedback based on prediction results.

example:

clustering: Take collection of 10 million diff. genes, group gene into groups that are somehow similar or related by different variables, such as lifespan, location, roles ... and so on.

Non-clustering: The "Cocktail party algorithm", allow you to find structure in a chaotic environment. (identify individual voices and music from a mesh of sounds at a cocktail party).

# Model Representation

## Example

Training set of housing prices

Size in feet <sup>2</sup> ( $x$ )	Price in 1000's ( $y$ )
2104	460
1234	232
⋮	⋮

$m=47$  example

$m$  = Number of training examples

$x$  = "input" variable / features

$y$  = "output" variable / "target" variable

$(x, y)$  = one training example

$(x^{(i)}, y^{(i)})$  —  $i$ th training example

$x^1 = 2104$        $y^1 = 460$

$x^2 = 1234$        $y^2 = 232$

$(x^{(i)}, y^{(i)})$ ;  $i = 1, \dots, m$  — is called a training set.

use  $X$  to denote the space of input values and  $Y$  to denote space of output values.

function  $h: X \rightarrow Y$ , to make  $h(x)$  is a good "predictor" for corresponding value of  $y$ .  $h$  can be called a hypothesis.

## Training Set

↓  
Learning  
algorithm

↓  
 $x \rightarrow [h] \rightarrow \text{predicted } y$   
(living area of house)      (predicted price of house)

When target variable we trying to predict is continuous, such as housing example, problem is a regression problem, when  $y$  take on only a small number of discrete value, (if giving living area, we predict it is a house or apartment), is a classification problem.

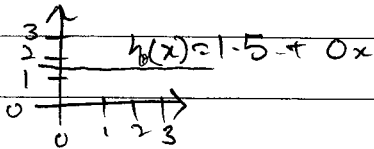
## Cost Function

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

$\theta_1$ 's: Parameters

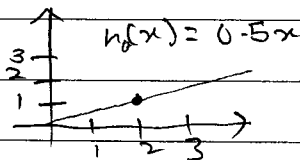
How to choose  $\theta_0, \theta_1$ 's?

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



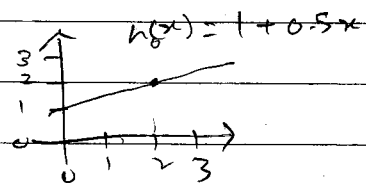
$$\theta_0 = 1.5$$

$$\theta_1 = 0$$



$$\theta_0 = 0$$

$$\theta_1 = 0.5$$



$$\theta_0 = 1$$

$$\theta_1 = 0.5$$

Idea: Choose  $\theta_0, \theta_1$  so that  $h_{\theta}(x)$  is close to  $y$  for our training examples ( $x, y$ )

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

$\uparrow$  # of training set  
 $\uparrow$  hypothesis of prediction of house price  
 $\uparrow$  actual house price

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

minimize  $J(\theta_0, \theta_1)$

cost function / Squared error function

- Measure the accuracy of our hypothesis function using cost function

- This takes an average difference of all results of hypothesis with input  $x$  and actual output  $y$ .

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

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

- The function is otherwise called "Square error function" or "Mean square error". Mean is halved ( $\frac{1}{2}$ ) as convenience for computation of gradient descent, as the derivative term of square function will cancel out  $\frac{1}{2}$  term.

# Cost Function - Intuition 1

No.

Date

Hypothesis:

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

Parameters:

$$\theta_0, \theta_1$$

Cost Function:

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

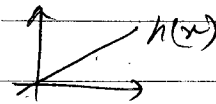
Goal = minimize  $J(\theta_0, \theta_1)$   
 $\theta_0, \theta_1$

simplified

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

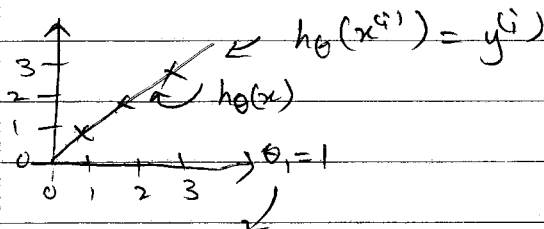
$$\theta_0 = 0$$

$\theta_1$



$h_{\theta}(x)$

for fixed  $\theta$ , this is a function of  $x$



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

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

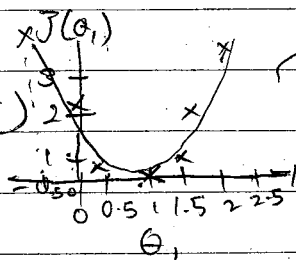
$$= \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0 \quad m=3, \text{ 所以为0个}$$

↓

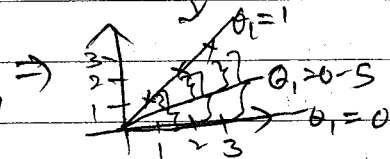
$$J(1) = 0$$

$J(\theta_1)$

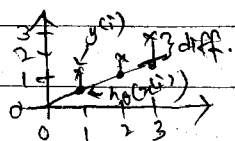
(function of parameter  $\theta_1$ )



Keep increasing when leaving further from the line



When graph go to  $\theta_1 = 0.5$



$\theta_1 = 1 = \text{global minimum}$

choose the value of  $\theta_1$  that can bestly is to minimize  $J(\theta_1)$

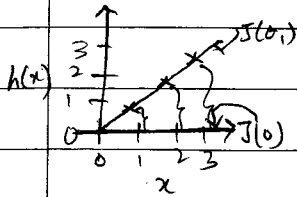
$$J(0.5) = \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

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



Example:

Training set  $m=3$ , hypothesis representation is  $h_{\theta}(x) = \theta_1 x$  with parameter  $\theta_1$ .  
 The cost function  $J(\theta_1)$  is  $J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$ . What is  $J(0)$ ?



$$J(0) = \frac{1}{2m} ((0-1)^2 + (0-2)^2 + (0-3)^2)$$

$$= \frac{14}{6}$$

Objective  $\Rightarrow$  get best possible line

$\Rightarrow$  will be such so that average squared vertical distance of scatter points from the line will be the least. Ideally, line should pass through all points of our training data set. In such a case, value of  $J(\theta_0, \theta_1)$  will be 0. The following example

Linear Regression with one variable Cost function intuition 1) (contour plot) 3-D

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

Parameters:  $\theta_0, \theta_1$

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

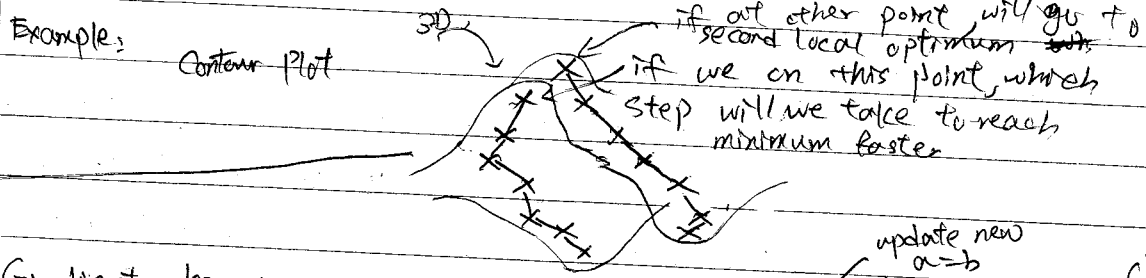
Goal: minimize  $J(\theta_0, \theta_1)$

Gradient descent  $\Rightarrow$  minimize cost function  $J$  by minimizing each  $\theta$  (minus derivative) to reach minimum value of  $J(\theta)$  using slope

Have some function  $J(\theta_0, \theta_1) \Rightarrow J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$   
 Want to  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1) \Rightarrow \min_{\theta_0, \dots, \theta_n} J(\theta_0, \dots, \theta_n)$   
 $\theta$  are to minimum value of  $J(\theta)$

Outline:

- Start with some  $\theta_0, \theta_1$  (initial guess)  $\Rightarrow$  (say  $\theta_0 = 0, \theta_1 = 0$ )
- Keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$   $\nearrow$  keep changing until we hopefully end up at a minimum



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  $j = 1$ )

Update (assign new variable)  $\nearrow$  learning rate (control how big a step we take downhill with creating descent)  $\nearrow$  derivative  $\nearrow$  to simultaneously update  $\theta_0$  and  $\theta_1$

$\alpha$  is step  $\nearrow$  increment

$j = 0, 1$  represents the feature index number

Assignment  $a := b$   
 $a := a + 1$   
 Truth assertion  $a = b$   
 $a = a + 1$  X

Correct: Simultaneous update

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

Incorrect:

temp0  $:= \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$   
 $\theta_0 := \text{temp0}$   
 temp1  $:= \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$  cannot because it will be updated first.  
 $\theta_1 := \text{temp1}$

The way we do is by taking derivative (the tangential line to a function) of our cost function. Slope of tangent is the derivative at the point and it will give us the direction to move towards.

Distance btwn each cross, "X" on the hill are determined by parameter  $\alpha$ .  
 Direction in which step is taken is determined by partial derivative of  $J(\theta_0, \theta_1)$ . Depending on ~~where~~ where one starts on graph, one could end up at different points.

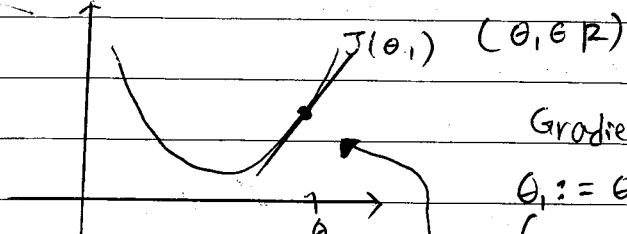
# Gradient Descent Intuition

$\frac{d}{d\theta_1} \Rightarrow$  Partial derivative

$\frac{d}{d\theta_1} \Rightarrow$  derivative

depending on number of parameters in function  $J$ .

in this is the same thing



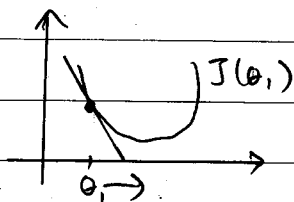
Gradient Descent will update  $\theta_1$

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

to find tangent line to function (slope  $\geq 0$ )

decrease to get closer to min-point

$$\theta_1 := \theta_1 - \alpha (\text{true num.})$$



$$\frac{d}{d\theta_1} J(\theta_1)$$

going to be -ve num.

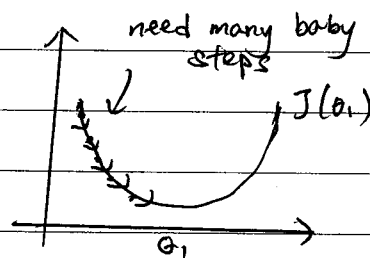
$$\theta_1 := \theta_1 - \alpha (-\text{ve num.})$$

$$\theta_1 := \theta_1 + \alpha (\text{num.})$$

increase to get closer to min-point

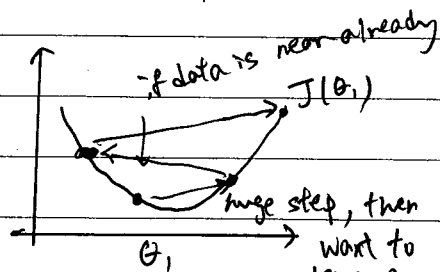
$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

if  $\alpha$  is too small, gradient descent can be slow to global minimum as need lot of steps.



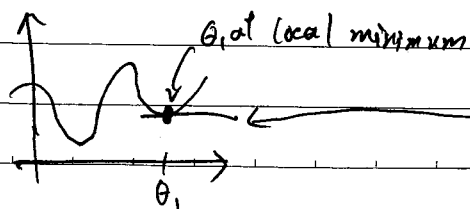
if  $\alpha$  is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.

won't converge, rather diverge



if data is near already huge step, then want to decrease, will become it keeps getting further away from the minimum

Suppose  $\theta_1$  is at local optimum of  $J(\theta_1)$ , what will one step of gradient descent  $\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$  do?



Answer: Will leave  $\theta_1$  unchanged

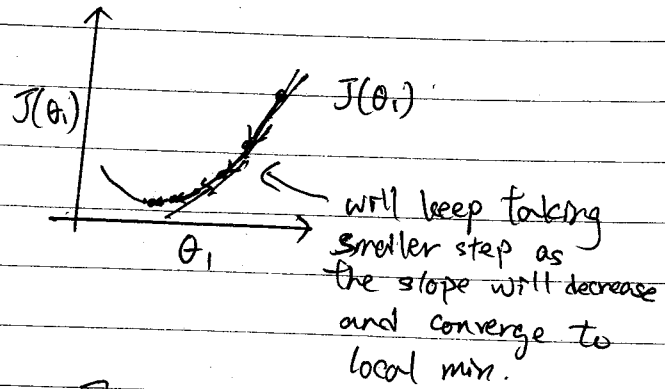
because slope of this line will become zero, thus derivative term is zero, gradient descent update  $\neq$  zero

And This is why gradient descent can converge to a local minimum, even with learning rate,  $\alpha$  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  $\alpha$  over time.



$\theta_1$  减一个大的 slope,  $\theta_1$  变小  
 $\theta_1$  再减小一点的 slope,  $\theta_1$  变得更小  
 $\theta_1$  减更小的 slope,  $\theta_1$  变得更更小

# Linear regression with 1 variable

## Gradient descent for linear regression

No. \_\_\_\_\_

Date \_\_\_\_\_

⇒ linear regression model

⇒ squared error cost function

Put together gradient descent ~~and~~ with cost function to give us an algorithm for linear regression or putting a straight line to our data

### Gradient descent algorithm

repeat until convergence {

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

(for  $j=1$  and  $j=0$ )

}

left term

### Linear Regression Model

linear hypothesis

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

squared error cost function

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

minimize

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

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

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

### Gradient descent algorithm

repeat until convergence {

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

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

}

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

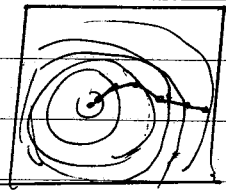
update  $\theta_0$  and  $\theta_1$  simultaneously

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

## "Batch" Gradient Descent

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

While gradient descent can be susceptible to local minima in general, the optimization problem we have posed here for linear regression has only 1 global and no other local, optima. Thus, gradient descent always converges (assuming  $\alpha$  is not too large) to global minimum. Indeed,  $J$  is a convex quadratic function. Here is an example of gradient descent as it is run to minimize quadratic function.



The ellipses are the contours of a quadratic function. Also shown is the trajectory taken by gradient descent, which was initialized at  $(48, 30)$ . " " in figure (joined by straight line) more the

successive values of  $\theta$  that gradient descent went through as it converged to its minimum.

# Matrices and Vectors

No. \_\_\_\_\_

Date \_\_\_\_\_

rectangular array of numbers

ex:  $\begin{bmatrix} 1402 & 191 \\ 1371 & 821 \\ 949 & 1437 \\ 147 & 1448 \end{bmatrix}$

rows

could be feature from learning problem/data

$\begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix}$

$\rightarrow 2 \times 3 \text{ matrix}$   
 $\mathbb{R}^{2 \times 3}$

Dimension of matrix: num of rows  $\times$  num of columns

$\rightarrow 3 \times 2 \text{ matrix}$   
 $\mathbb{R}^{3 \times 2}$

## Matrix Elements (entries of matrix)

$\rightarrow$  numbers inside the matrix

$$A = \begin{bmatrix} 1402 & 191 \\ 1371 & 821 \\ 949 & 1437 \\ 147 & 1448 \end{bmatrix}$$

$A_{ij}$  = "i, j entry" in the  $i^{\text{th}}$  row,  $j^{\text{th}}$  column.

$A_{11}$  refer to 1<sup>st</sup> row 1<sup>st</sup> column

$$A_{11} = 1402$$

$A_{43}$  = undefined (error)

$$A_{12} = 191$$

$$A_{32} = 1437$$

$$A_{41} = 147$$

Vector : An  $n \times 1$  matrix (matrix that has only 1 column)

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix} \xleftarrow{n=4} \text{4-dimensional vector } (\mathbb{R}^4)$$

$y_i = i^{\text{th}}$  element

$$y_1 = 460 \quad y_3 = 315$$

1-indexed vs 0-indexed.

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \quad y = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

### Linear Algebra

Addition and scalar multiplication & (multiply matrix)

Addition

$$\begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix}_{3 \times 2} + \begin{bmatrix} 4 & 0.5 \\ 2 & 5 \\ 0 & 1 \end{bmatrix}_{3 \times 2} = \begin{bmatrix} 5 & 0.5 \\ 4 & 10 \\ 3 & 2 \end{bmatrix}_{3 \times 2} \quad \text{add the } 3 \times 2 \text{ matrix of same dimension}$$

$$\begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix}_{3 \times 2} + \begin{bmatrix} 4 & 0.5 \\ 2 & 5 \end{bmatrix}_{2 \times 2} = \text{error}$$

real number  
Scalar Multiplication

$$3 \times \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 6 & 15 \\ 9 & 3 \end{bmatrix}$$

$$\begin{bmatrix} 4 & 0 \\ 6 & 3 \end{bmatrix} / 4 = \frac{1}{4} \begin{bmatrix} 4 & 0 \\ 6 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ \frac{3}{2} & \frac{3}{4} \end{bmatrix}$$



$$3 \times \begin{bmatrix} 1 \\ 4 \\ 2 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix} - \begin{bmatrix} 3 \\ 0 \\ 2 \end{bmatrix} \mid 3$$

$$= \begin{bmatrix} 3 \\ 12 \\ 6 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ \frac{2}{3} \end{bmatrix}$$

$$= \begin{bmatrix} 2 \\ 12 \\ \frac{31}{3} \end{bmatrix}$$

$$\begin{bmatrix} 4 \\ 6 \\ 7 \end{bmatrix} \mid 2 - 3 \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 2 \\ 3 \\ 7\frac{1}{2} \end{bmatrix} - \begin{bmatrix} 6 \\ 3 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} -4 \\ 0 \\ 7\frac{1}{2} \end{bmatrix}$$

3-dimensional vector

Matrix vector multiplication

$$\begin{bmatrix} 1 & 3 \\ 4 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 5 \end{bmatrix} = \begin{bmatrix} 16 \\ 4 \\ 7 \end{bmatrix}$$

(3) × 2      2 × 1      3 × 1

$$\begin{bmatrix} 1 & 2 & 1 & 5 \\ 0 & 3 & 0 & 5 \\ -1 & -2 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 14 \\ 13 \\ 7 \end{bmatrix}$$

$$1 \times 1 + 3 \times 5 = 16$$

$$4 \times 1 + 0 \times 5 = 4$$

$$2 \times 1 + 1 \times 5 = 7$$

House sizes:

→ 2104

→ 1416

→ 1534

→ 852

matrix

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

vector

$h_0(2104)$

$$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix} \times \begin{bmatrix} -40 \\ 0.25 \end{bmatrix} = \begin{bmatrix} 1 \times -40 + 0.25 \times 2104 \\ 1 \times -40 + 0.25 \times 1416 \\ 1 \times -40 + 0.25 \times 1534 \\ 1 \times -40 + 0.25 \times 852 \end{bmatrix}$$

(4) × 2      2 × 1      4 × 1

4 dimensional vector

Octave → prediction = DataMatrix \* parameters

for j=1 to 4

matrix-matrix multiplication

$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 0 & 1 \\ 5 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 10 \\ 9 & 16 \end{bmatrix}$$

$\begin{matrix} 2 \times 3 & 3 \times 2 & 2 \times 2 \end{matrix}$

$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 0 \\ 5 \end{bmatrix} = \begin{bmatrix} 9 \\ 14 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 3 \\ 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 10 \\ 14 \end{bmatrix}$$

House sizes: Have 3 competing hypotheses:

2104

1416

1534

852

$$1. h_0(x) = -40 + 0.25x$$

$$2. h_0(x) = 300 + 0.1x$$

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

Matrix

Matrix

$$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix} \times \begin{bmatrix} -40 & 200 & -150 \\ 0.25 & 0.1 & 0.4 \end{bmatrix} = \begin{bmatrix} 486 & 410 & 692 \\ 314 & 342 & 416 \\ 200 & 353 & 664 \\ 173 & 285 & 191 \end{bmatrix}$$

Matrix multiplication properties

$3 \times 5 = 5 \times 3 \rightarrow$  order not important ~~because~~ <sup>called</sup> "Commutative of multiplication of real numbers"

Let  $A$  and  $B$  be matrices. In general,

$A \times B \neq B \times A$  (not commutative)

$$\text{Eg. } \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 2 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 \\ 2 & 0 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 2 & 2 \end{bmatrix}$$

$A \times B$   
 $m \times n \quad n \times m$

$A \times B$  is  $m \times m$

$B \times A$  is  $n \times n$

$$3 \times (5 \times 2) = (3 \times 5) \times 2 \quad \text{"Associative"}$$

Matrix:

$$A \times B \times C = \begin{matrix} A \times (B \times C) \\ (A \times B) \times C \end{matrix} \rightarrow \text{same answer (Associative property)}$$

Identity Matrix

1 is identity

$$1 \times z = z \times 1 = z$$

only  $z$

Denoted  $I$  (or  $I_{n \times n}$ )

Examples of identity matrices:

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

$2 \times 2$

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

$3 \times 3$

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

$4 \times 4$

Informally:

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

For any matrix  $A$ ,

$$A \cdot I = I \cdot A = A$$

$\begin{matrix} \nearrow & \uparrow & \uparrow & \uparrow & \nwarrow \\ m \times n & n \times n & m \times m & m \times n & m \times n \end{matrix}$

Note:

 $AB \neq BA$  in general

$$AI = IA \quad \checkmark$$

Inverse and Transpose

 $I = \text{"identity"}$ 

$$3(3^{-1}) = 1$$

$$12 \times (12^{-1}) = 1$$

Not all numbers have inverse.  $\Rightarrow 0(0^{-1})$  undefined  
Matrix inverse.

If  $A$  is an  $m \times m$  matrix, and if it has an inverse,

$$A(A^{-1}) = A^{-1}A = I$$

Eg.  $\begin{bmatrix} 3 & 4 \\ 2 & 16 \end{bmatrix} \begin{bmatrix} 0.4 & -0.1 \\ -0.05 & 0.075 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I_{2 \times 2}$

Matrix transpose

Example:

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

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

imagine a 45° mirror,  
and all number reflect here

Let  $A$  be an  $m \times n$  matrix,  $B = A^T$ . $B$  is an  $n \times m$  matrix, and

$$B_{ij} = A_{ji}$$

$$B_{12} = A_{21} = 2$$

row  $\nearrow$   
column

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

$$\begin{bmatrix} 3 & -5 & 4 \end{bmatrix}$$