

Machine Learning Class Note

Reference : ① Prof. Andrew Ng's

learning class @ Coursera
Stanford University

② Wikipedia

③ CS231 class @ Stanford University

Topic 1 : Introduction

* What is Machine Learning?

[Wikipedia] Machine learning is the scientific study of algorithms and statistical models that computer systems use to effectively perform a specific task without using explicit instructions, relying on patterns and inference instead.

* can be ~~split~~ divided into :

ML	{	supervised learning
		unsupervised learning
		weakly supervised learning
		...

supervised learning: a machine learning task of learning a function that maps an input to an output based on example input-output pairs. ("right answer" given)

unsupervised learning: a machine learning task without a teacher. ("right answer" not given)

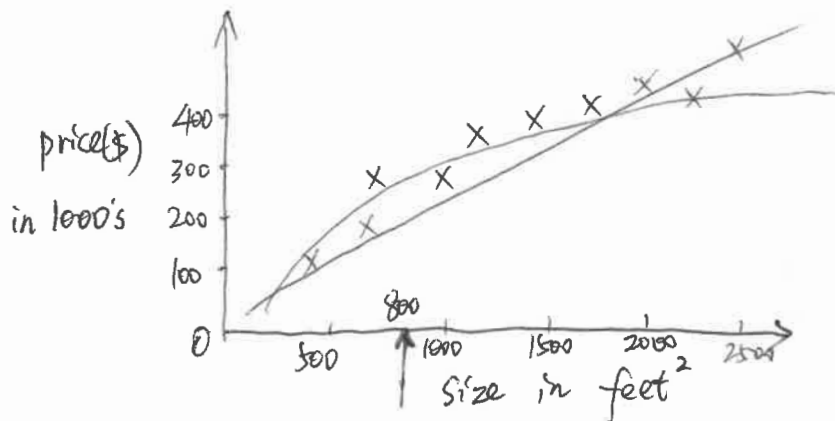
Given a data, find the structure of data without "right answer".

weakly supervised learning: similar to supervised learning, but with weak supervision, such as lower-quality labels.

("right answer" weakly given)

Example of Supervised Learning:

① Housing price prediction



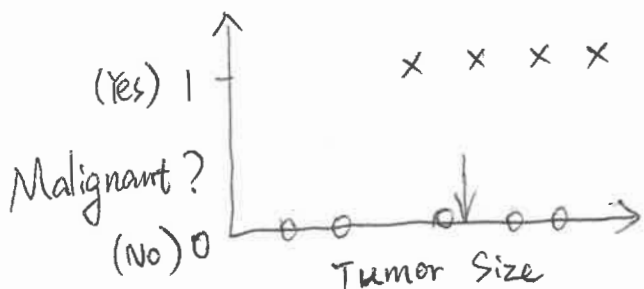
Plot Data in a city.

What is the estimate price for a house with 800 feet² in that city?

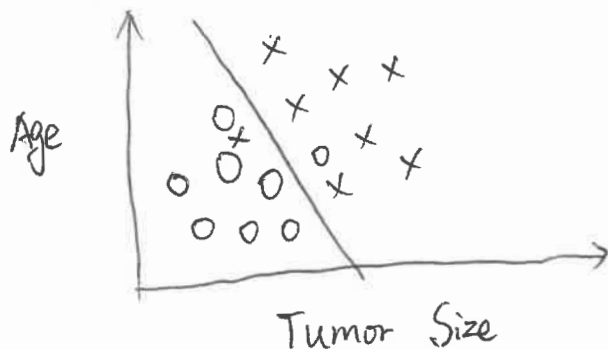
Supervised learning: "right answer" given

Regression: Predict continuous valued output (price)

② Breast Cancer Classification (malignant, benign)



One feature/attribute
multiple features/attributes



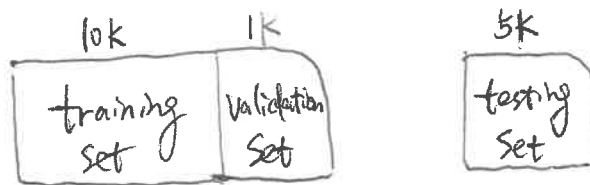
Supervised learning: "right answer" given

Classification: Discrete valued output (0 or 1)

Or: 0, 1, 2, 3, ...
benign cancer type 1 cancer type 2 cancer type 3

Supervised learning:

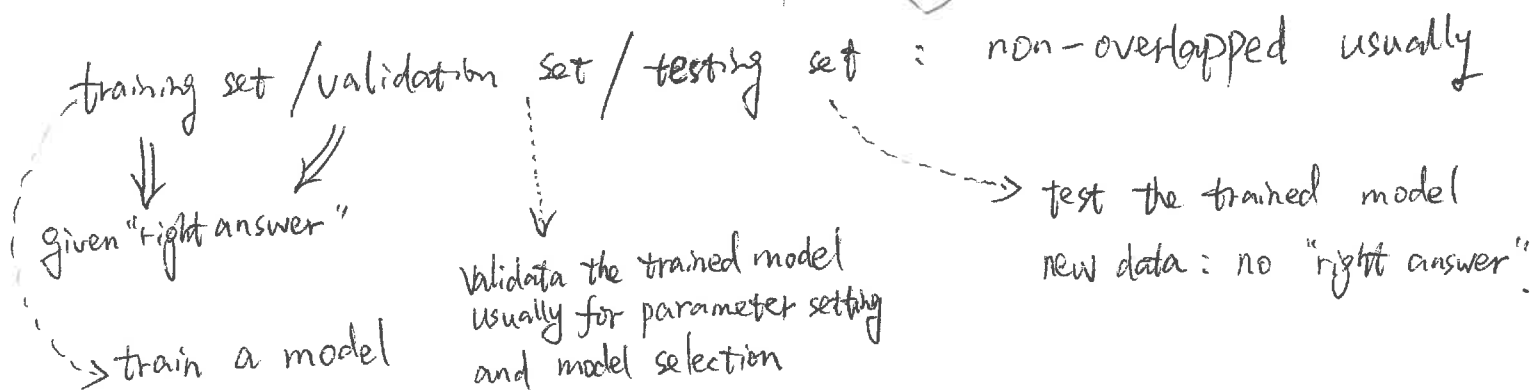
Data is divided into:



↓
 $m \times n$ matrix $\begin{cases} m: \text{sample \#} \\ n: \text{feature dimension/attribute} \end{cases}$

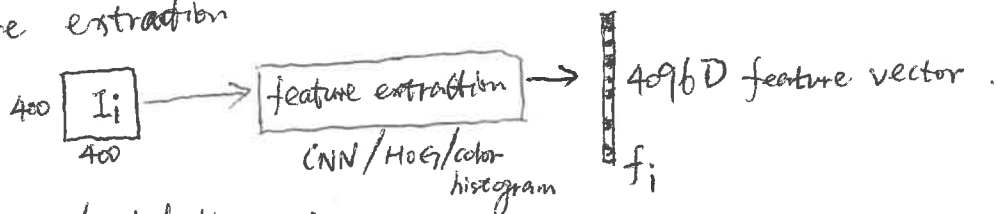
example data:

	Age	Tumor Size	Label/gt
Patient 1	20	5	1
Patient 2	25	2.5	0
⋮	⋮	⋮	⋮
Patient M	45	5	1



Take **K-NN** (K-nearest neighbor) based Image Classification as example:

① feature extraction



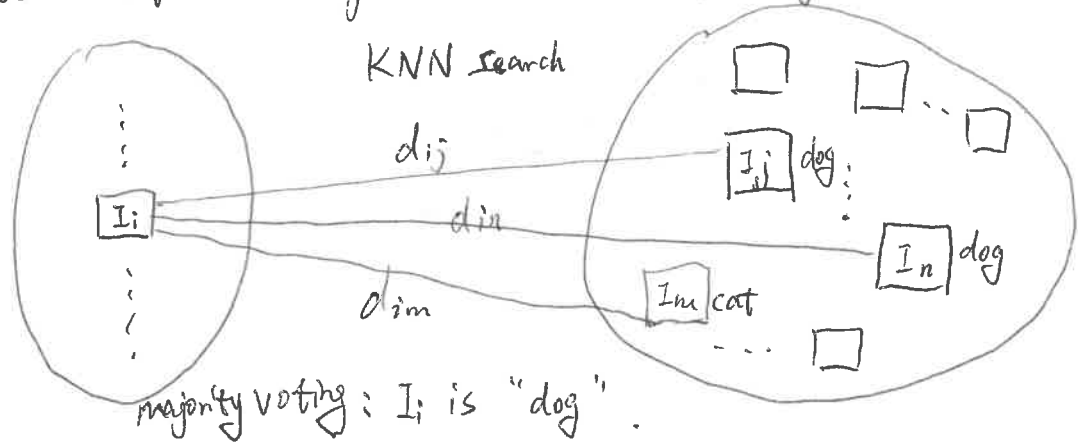
② training set / validation set

10,000 images 1,000 images : each image can be represented with one feature vector
 I_i f_i
each image has a classification label by human
 I_i dog, cat, monkey, ...
 1000 classes

③ naive training (Method 1)

Validation set: 1000 images

training set: 10,000 images



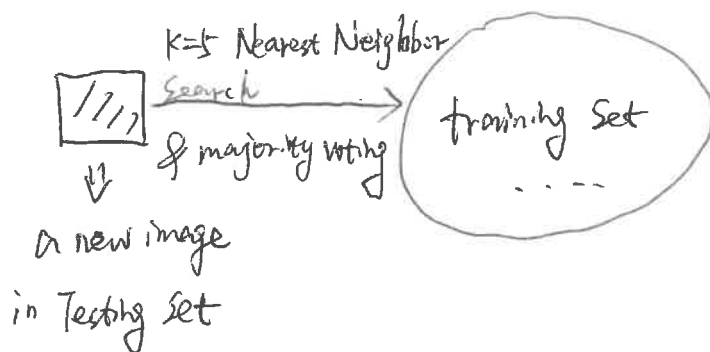
- * distance of Image i and Image j can be the l_2 distance of f_i and f_j feature vectors.
- * Given the label of each image on Validation Set, we can know the majority voting result is correct or not \Rightarrow an accuracy can be obtained
- * Using Validation Set, Search on Training Set. Suppose that:
 - $k=1 \Rightarrow \text{accuracy} = 50\%$
 - $k=2 \Rightarrow \text{accuracy} = 55\%$
 - $k=3 \Rightarrow \text{accuracy} = 58\%$
 - $k=4 \Rightarrow \text{accuracy} = 60\%$
 - $k=5 \Rightarrow \text{accuracy} = 62\%$
 - $k=6 \Rightarrow \text{accuracy} = 55\%$
 - \vdots

★ best validation accuracy obtained when $k=5$

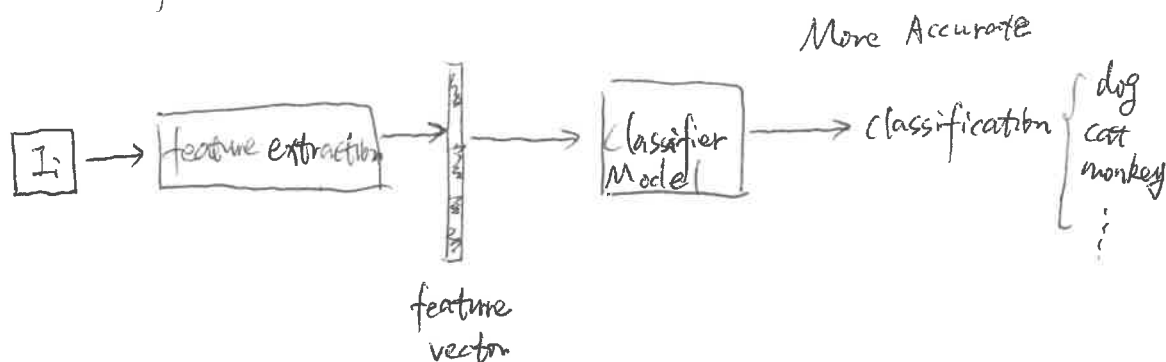


$k=5$ after training

④ Use the trained model on Testing Set (new 5000 images).



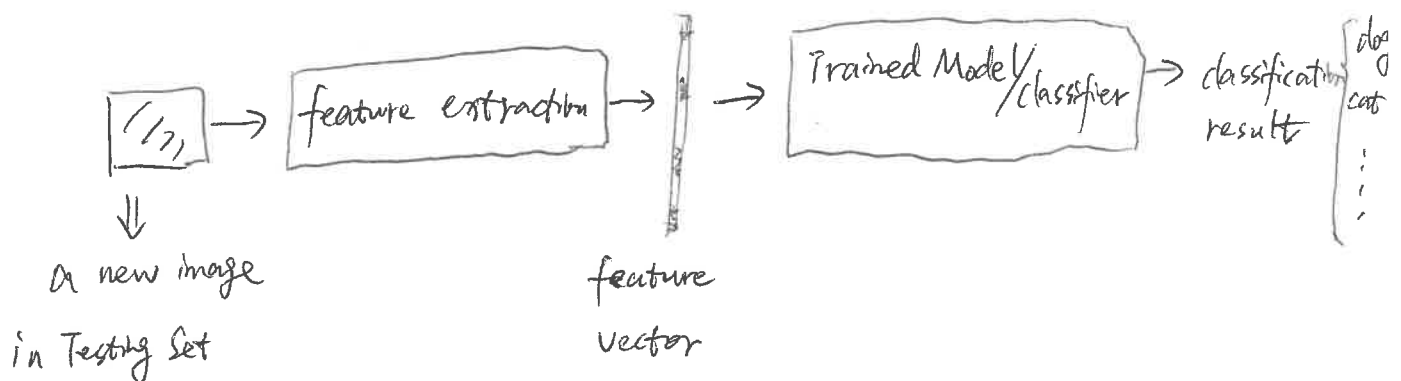
3.2 train a classifier, like SVM, NN, CNN, etc. (Method 2)



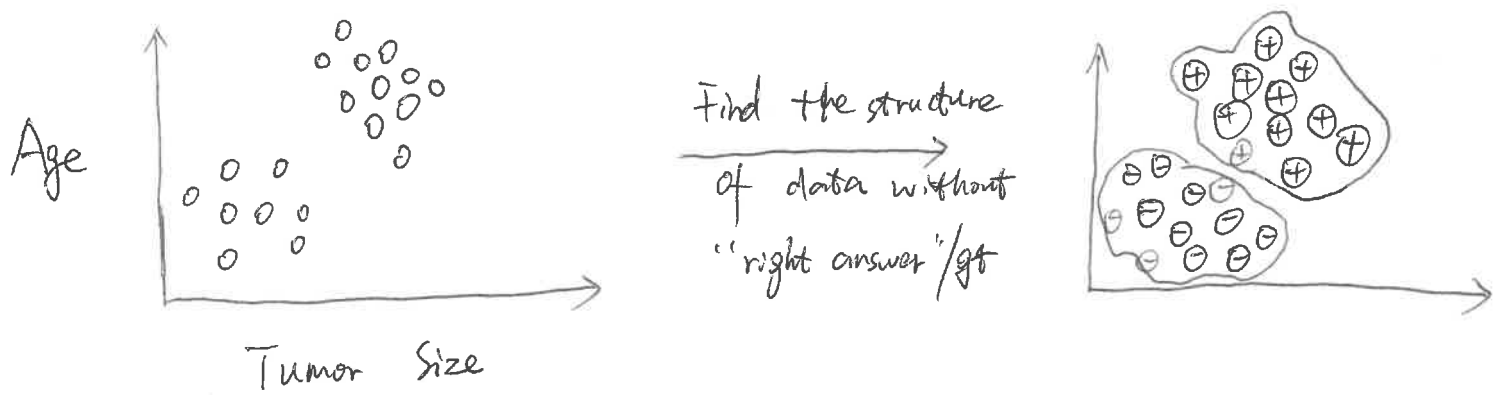
Training = adjust the classifier to make the prediction (output) by the classifier can be more close to the manually labeled result/gt.
 on Training \Leftarrow Repeat it for many iterations (for example: 20,000 times).
 Set 10,000 images

Validation: tune some hyperparameters to select the best classifier that has best performance on validation set.
 on Validation \Leftarrow hyperparameters: like learning rate, model structure, iteration number, etc.
 Set 1,000 images

4.2 use the trained model (classifier) on Testing Set (new 5000 images)



Example of unsupervised learning.

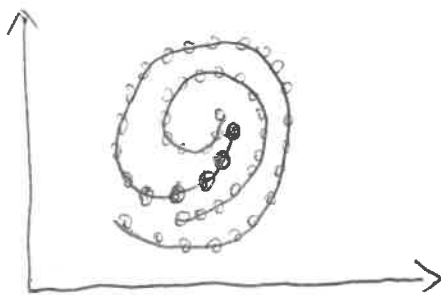


① clustering (kmeans, ...)

as above: cluster the data into k groups based on similarity

like Google News: Sports, Technology, Business, ...

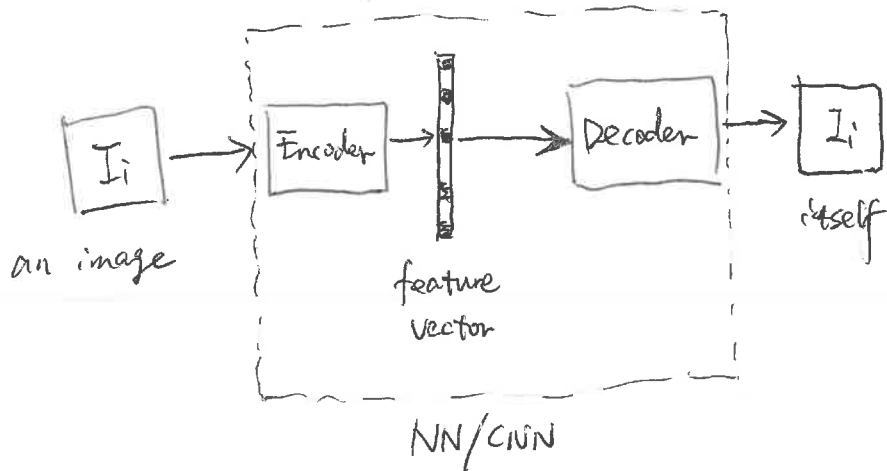
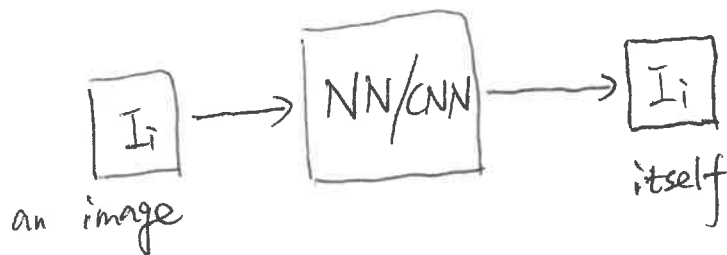
for example:
L2 distance



kmeans cannot solve this clustering problem.
Need more advanced algorithm.

② Feature Extraction by Encoder & Decoder.

Feature: color histogram, histogram
HOG, RGB color, ...
learned feature, ...



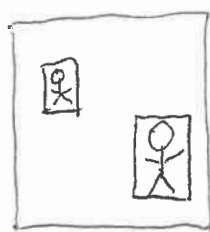
Train a NN/CNN on
Training Set (10,000 images)
for many iterations
(say, 20,000 times)

Example of Weakly Supervised Learning.

Similar as supervised learning, but with lower-quality labels.
("right answer" weakly given)

Supervised learning for Human detection:

On Training Set (10,000) images, for each image:

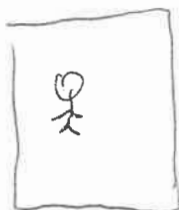


Strong supervision

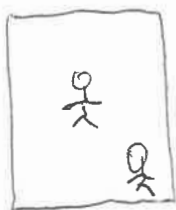


Weakly supervised learning for Human detection:

On Training Set (10,000) images, for each image:



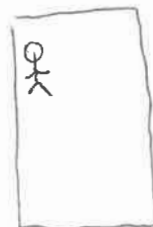
Yes/1



Yes/1



No/0



Yes/1

Weak supervision: No bounding-box annotation
but Yes/No annotation.

* Overfitting.

On supervised learning, if the trained model performs very great on training and validation set, but performs poor on testing set. It is called 'overfitting'.

training accuracy: 99%

validation accuracy: 98%

testing accuracy: 40%

Machine Learning Class Note

03/2019

Reference : ① Prof. Andrew Ng's Machine Learning class @ Coursera Stanford

② Wikipedia

③ CS231 class @ Stanford

Topic 2: Linear Regression with one variable

Training set of
housing prices

size in feet ² (x)	Price in 1000's (y)
2104	460
1416	232
1534	315
852	178
...	...

$\Rightarrow m=47$

Notation :

m : Number of training examples

x 's: "input" variable / features

y 's: "output" variable / "target" variable

$$x^{(1)} = 2104$$

$$x^{(2)} = 1416$$

$$y^{(1)} = 460$$

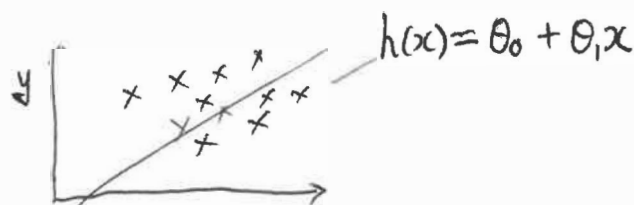
$$y^{(2)} = 232$$

(x, y) : one training example

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

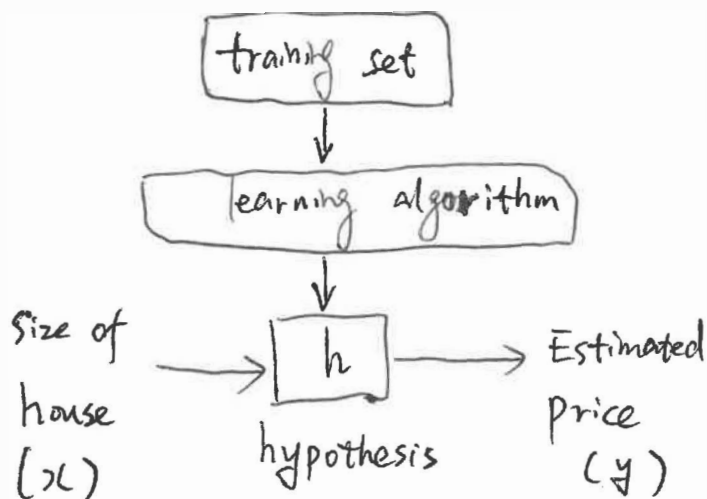
How to represent h ?

Say $h_0(x) = \theta_0 + \theta_1 x$
shorthand: $h(x)$



h maps from x 's to y 's.

Linear regression with one variable

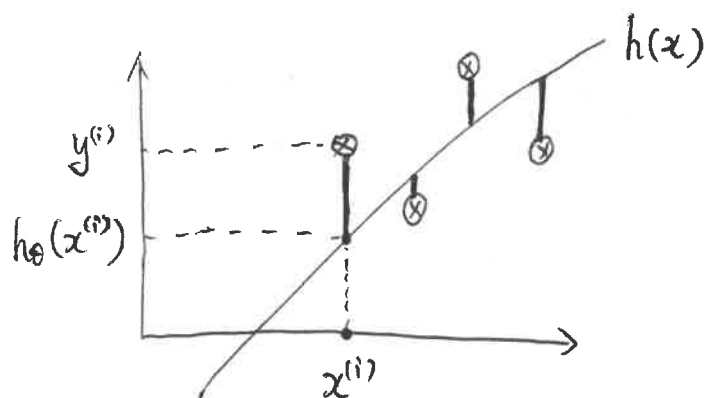


Hypothesis: $h_0(x) = \theta_0 + \theta_1 x$

Parameters: θ_0, θ_1

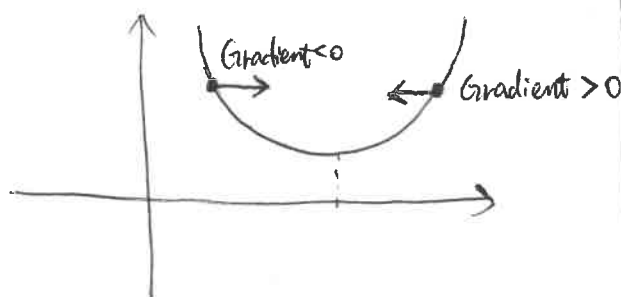
cost function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$
misalignment

Goal: Minimize $J(\theta_0, \theta_1)$ \Rightarrow Find a linear line (θ_0, θ_1) to well fit the training samples so as to minimize the cost function / misalignment.



Optimization: minimize the cost function

Gradient descent:



If the cost function is convex, moving to $-1 \times$ Gradient direction will ~~find~~ find smaller cost function. Repeat this until finding the global minimum.

Gradient descent Algorithm

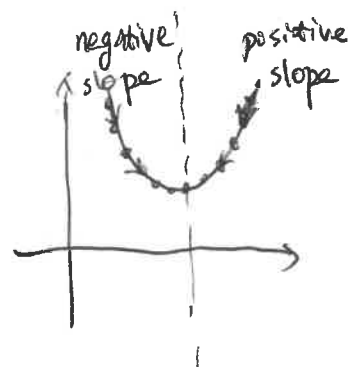
Repeat until convergence {

$$\theta_j = \theta_j - \alpha \cdot \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \text{for } j=0 \text{ and } j=1$$

}

\rightarrow learning rate

\rightarrow derivative



Correct:

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

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

$$\theta_0 = \text{temp0}$$

$$\theta_1 = \text{temp1}$$

Incorrect:

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

$$\theta_0 = \text{temp0}$$

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

$$\theta_1 = \text{temp1}$$

How to compute the gradient?

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

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

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

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

$$\frac{\partial}{\partial x} x^2 = 2x \quad \frac{\partial}{\partial x} (x+3)^2 = 2(x+3) \cdot 1$$

~~scribble~~ $\frac{\partial}{\partial x} (kx+b)^2 = 2(kx+b) \cdot \frac{\partial}{\partial x} (kx+b) = 2(kx+b) \cdot k$

$$\frac{\partial}{\partial b} (kx+b)^2 = 2(kx+b) \cdot 1$$

partial gradient

Topic 3: Linear Regression with multiple variables.

size (feet ²) x_1	Number of bedrooms x_2	number of floors x_3	Age (years) x_4	Price (in \$1000) y
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮

Notation: n : number of features

$x^{(i)}$: input (features) of i -th training example.

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

Hypothesis:

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Let us define $x_0 = 1 \Rightarrow x_0^{(i)} = 1$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

column vector

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

column vector

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



$$[\theta_0, \theta_1, \dots, \theta_n] \cdot \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Hypothesis:

$$h_\theta(x) = \theta^T \cdot x = \theta_0 \cdot x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

parameters: $\theta: \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$ $(n+1)$ -dimensional vector

cost function:

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

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

New Algorithm

Repeat {

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

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

}

Feature scaling / normalization.

$$0 \leq x_1 \leq 3 \quad \checkmark$$

$$-2 \leq x_2 \leq 0.5 \quad \checkmark$$

$$-100 \leq x_3 \leq 100 \quad \times$$

Mean normalization:

$$x_j = \frac{x_j - \mu_j}{s_j}$$

x_j : j -th feature

μ : mean / average value of x_j

s : stand deviation

or range (max - min)

Topic 4: clustering / kmeans

Given a set of observations (x_1, x_2, \dots, x_n) , where each observation is a d -dimensional vector.

x_1 : n -dimensional feature vector

x_2 : n -dimensional

\vdots

x_n : n -dimensional

kmeans clustering aims to partition the n observations into k ($k \leq n$) sets $S = \{S_1, S_2, \dots, S_k\}$ to minimize the within-cluster sum of

squares:

$$\arg \min_S \sum_{i=1}^k \sum_{x \in S_i} \|x - \mu_i\|^2$$

Where μ_i is the mean point / vector in S_i .

Standard Algorithm:

① Initialize K means/clusters randomly

Repeat until convergence:

① Assignment step: Assign each observation to the "nearest" mean/cluster

② Update step: Calculate the new means of the re-assigned observations in the new clusters

* This algorithm has converged when the assignments no longer change.

* This algorithm does not guarantee to find the optimum.

How to define K ? Some heuristic algorithms:

$$\text{Loss}(k) = \frac{d_{\text{intra}}}{d_{\text{inter}}}$$

where d_{intra} is the average distance between each observation to the center within each cluster

d_{inter} is the average distance of each inter-cluster center.

$\underset{(k)}{\operatorname{argmin}} \text{Loss}(k)$: eg., Find the minimum loss when $k=2, \dots, 10$.
for example.