

W1-Instrcution

1. Welcome

1.1 Machine learning example:

- You probably use it dozens of times a day without even knowing it:
 - >1. Each time you do a web search on Google or Bing, that works so well because their machine learning software has figured out **how to rank what pages**;
 - >2. When Facebook or Apple's photo application recognizes your friends in your pictures, that's also machine learning;
 - >3. Each time you read your email and a spam filter saves you from having to wade through tons of spam, again, that's because your computer has learned to distinguish spam from non-spam email.

1.2: Machine learning definition:

A science of getting computers to learn without being explicitly programmed.

e.g.: getting robots to tidy up the house: having the robot watch you demonstrate the task and learn from that. The robot can then watch what objects you pick up and where to put them and try to do the same thing even when you aren't there.

1.3: AI application: reasons excited about AI, or artificial intelligence problem, is **building truly intelligent machines**, we can do just about anything that you or I can do.

Many scientists think the best way to make progress on this is through learning algorithms called **neural networks**, which **mimic how the human brain works**, and I'll teach you about that, too. In this class, you learn about machine learning and get to implement them yourself.

Machine Learning

- Grew out of work in AI
 - New capability for computers
- Examples:
- Database mining
Large datasets from growth of automation/web.
E.g., Web click data, medical records, biology, engineering
 - Applications can't program by hand.
E.g., Autonomous helicopter, handwriting recognition, most of Natural Language Processing (NLP), Computer Vision.
 - Self-customizing programs
E.g., Amazon, Netflix product recommendations
 - Understanding human learning (brain, real AI).



Andrew Ng

2. Instruction on what is MT

2.1 MT: older definiaton: field of study give computes the ability to learn witout being explicitly programmed.

Checker player algorithm by Arthur Samuel: algorithm play tens of thousands of games against itself, and by watching what sorts of board positions tended to lead to win/losses, the checker palyser algorithm learned over time what are good board position and bad positions.-->learned to play checker better.
note: computer have patient to play ten of thousand to learn experiance, while human have no such patience.

2.2 MT latest definition: a computer program is said to learn from experience E with respect to some task T and some performance measure P, if its performance on T , as mearused by P, improves with experience E.

2.2: MT algorithms

- **Supervised learning:** teach computer how to do something;
- **Unsupervised learning:** computer learn by itself.
- Others: reinforcement learning, recommender systems.

Need to pay attention: Practical advice for applying learning algorithms.

Machine Learning definition

- Arthur Samuel (1959). Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.
- Tom Mitchell (1998) Well-posed Learning Problem: A computer program is said to *learn* from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E.

Machine learning algorithms:

- Supervised learning
- Unsupervised learning

Others: Reinforcement learning, recommender systems.

Also talk about: Practical advice for applying learning algorithms.

3. Supervised learning_Insturciton

3.1 Supervised learning defination:

Give the algorithm a data set, in which 'right answers' were given, **task of algorithm was to produce more of these right answers**.
In every sample in our data set, we are told what is the 'correct answer', that we would have quite liked the algorithm have predicted on that example.

3.2 Regression problem: trying to predict continuous valued output (连续值输出)

3.3: Classification problem: trying to predict a discrete valued output (离散值输出, e.g: 0 or 1)

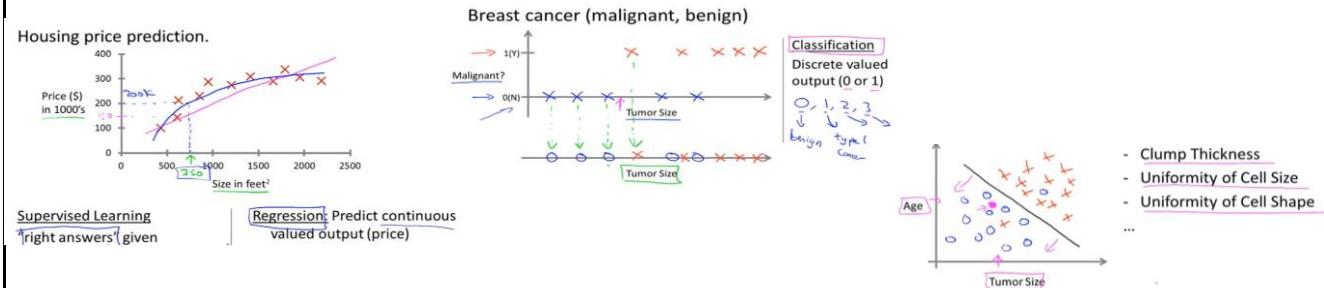
Note: classificaiton could have more than two possible values for the output

e.g:

- house price prediction: supervised learing, also regression problem
- Breast cancer: supervised learing, classificalon problem:
 - in this case: use only feature (tumor size) to define classification type;
 - have different plot draw to show: two axile (x tumor size, Y, output 1, 0) with same shape plots; and only 1 axile (tumor size) with different plots represent different result.
 - breast cancer : two feature (tumor size and age); could have more feature;

Note: learning algorithm can deal with not only 1, 2, 3 features, but an **infinite number of features**.

with infinite features, computer will run out of memory-->support vector machine: there will be neat mathmatical trick, allow computer to deal with an infinite number of features. How -->kernal, then have m features.



4. unsupervised learning_Insturciton

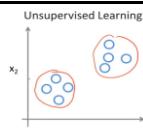
4.1: training data:

supervised learning : training data have right answer;
unsupervised learning: training data doesn't have any labels , **or all has the same label** (e.g: anomaly detection), or really no labels.

4.2: fucntion , task:

We are given a data set, and not told what to do with it, not told what each data point is.--> just told , here is a data set, can you find some structure in the data?

4.3 Clustering Algorithm



- Clustering algorithm

Algorithm break data set into two/more separate clusters.

Application e.g:

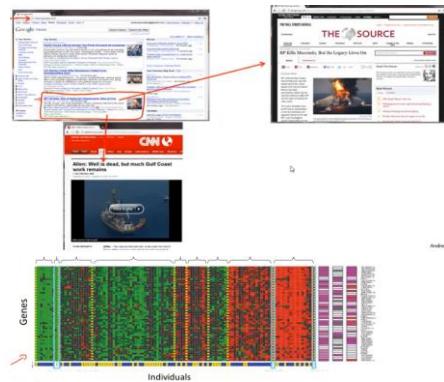
1. google news (news.google.com): everyday google news goes and look at hundreds of thousands of new stories on the web, and it groups them into cohesive news stories: automatically cluster them together, so the new stories that all about the same topic get displayed together;

2. Understanding genomics:

- Given DNA microarray data: have a group of different individuals, and for each of them, measure how much they do or do not have a certain gene-technically measure how much certain genes are expressed.
colors red, green, gray and so on, show the degree to which different individuals do or do not have a specific gene.

- Run clustering algorithm to group individuals, into different categories or into different types of people.

This one is an unsupervised learning program as do not tell in advance the 'right type' to certain genes. Instead only saying here is a bunch of data, do not know what's the data, do not know who is in what type, even do not know what the different types of people are, --> can you/algorithm automatically find structure in the data, automatically cluster the individuals into these types that i do not know in advance?



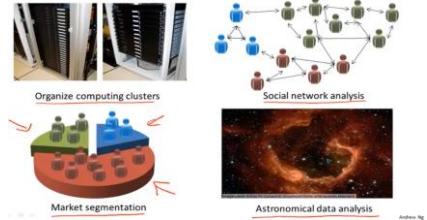
4.4 Applications_clustering algorithm

- Organize large computer clusters: try to find out which machines tend to work together (if put them work together, can make data center work more efficiently)

- Social network analysis: Given knowledge about which friends you email/facebook / circle friends the most, automatically classify which are cohesive groups of friends, which are groups of people that all know each other.

- Market segmentation: have huge databases of customer information, and automatically discover market segments and automatically group customers into different market segments. (do not know in advance market segments and customer group)

- Astronomical data analysis: clustering algorithms give useful theory of how galaxies are formed.

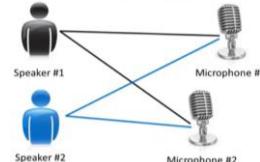


4.5 Cocktail party problem

1. e.g:

- two people in a room, talking at the same time;
- two microphones in the room: each microphone at different distance from two people.
- Each microphone records: a different/overlapping combination of these two speaker voices (maybe speaker 1 is a little bit louder on microphone1, and speaker 2 louder in microphone 2).
--give two microphone recordings to an unsupervised learning algorithm-cocktail party algorithm, and tell algorithm find structure in this data.
--cocktail party algorithm: listen to these audio recordings and say it sounds like two audio recordings that are being summed together; and algorithm will separate out these two audio sources that are being summed together to form other recording.

Cocktail party problem



Cocktail party problem algorithm

`[W,s,v] = svd((repmat(sum(x.*x1),size(x,1),1).*x)*x');`

2. Code_cocktail party problem

only one line code: `[W, s, v]=svd((repmat(sum(x.*x1, size(x,1),1).*x)*x');`

svd function: singular value decomposition (奇异值分解) .turn out to be a linear algebra routine (线性代数常规函数) that is just built into Octave.

If do with c+++o or Java or python, need many codes linking complex C++, or Java libraries. much more complicated to do so in those languages.

4.6 Octave programming environment

use in MT programming;

Octave: free open source software.

Use Octave or Matlab many learning algorithm become just a few lines of code to implement.

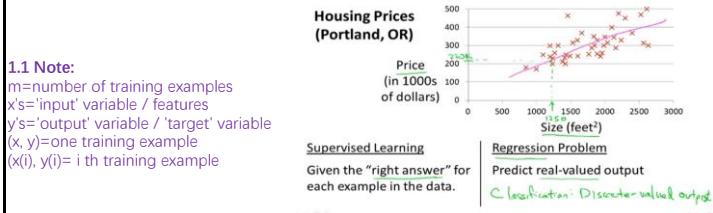
Normal process:

1. First prototype learning algorithm/ software in Octave;
Because software in Octave makes it incredibly fast to implement these learning algorithms.
2. Only after gotten it to work, then might migrate it to C++ or Java or whatever.

-->In this way much faster than if starting out in C++.

W2-Linear regression with one variable

1. Model-representation_Linear regression with one variable



Training set of housing prices (Portland, OR)

Size in feet ² ($x^{(i)}$)	Price (\$ in 1000's)($y^{(i)}$)
2104	232
1416	315
1534	178
852	460
...	...

Notation:

- m = Number of training examples
- x 's = "input" variable / features
- y 's = "output" variable / "target" variable
- $(x^{(i)}, y^{(i)})$ = one training example
- $(x^{(i)}, y^{(i)})$ = i th training example

1.2 Estimation process

train model: training set -> learning algorithm -> prediction function h
predict with trained h : new input -> h -> prediction y

hypothesis h : maps from x 's to y 's

1.3 hypothesis h : Linear regression

predict h is linear function: $h_{\theta}(x) = \theta_0 + \theta_1 x$

could be more complicated non-linear function, just start with this example first of fitting linear function, and build on this eventually more complex models and more complex algorithms.

2. Cost function_Linear regression with one variable

2.1: Parameter: θ_0, θ_1

choose parameter to make prediction function h correspond to a good fit to the training set (prediction close to training set).

IDEA: Choose parameters at least fit training set: give the x 's in the training set, make reasonable accurate predictions for the y values.

2.2 COST function

minimize square of(prediction output - label output), over θ_0, θ_1 .

$1/m \rightarrow$ average square, $1/2m$: do not change object, just makes some of the math a little easier..

Overall objective function-Squared error function:
finding parameter θ_0, θ_1 minimize the cost $J(\theta_0, \theta_1)$

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

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

Note: there are other method cost function, but square error function is the most commonly used.

3. Cost function intuition_linear regression with one variable

Prediction function $h_{\theta}(x)$: for fixed θ , this is function of x

Cost function $J(\theta)$: function of the parameter θ .

In plotting pic, horizontal axis is labeled parameter θ .

Summary:

- each parameter θ_i corresponds to a different prediction function: line fit;
- each prediction fit line, get the cost result $J(\theta_i)$;
- plot pic based on pair (cost result J , θ_i)

- Object: choose parameter, minimize cost J (corresponding to find a straight line that fit training set well).

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

Parameters:
 θ_0, θ_1

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

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

Simplified

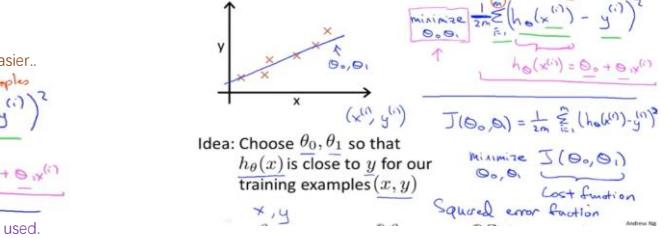
$$h_{\theta}(x) = \theta_1 x$$

$$\theta_0 = 0$$

$$\theta_1$$

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

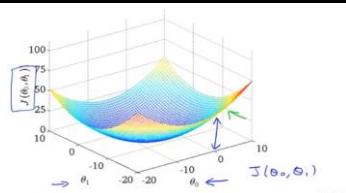
$$\text{minimize}_{\theta_1} J(\theta_1)$$



4. Cost function intuition II_linear regression with one variable

4.1 3D Visualize cost function: with all parameters (θ_0, θ_1)

Note: cost function is bowl shape over each parameter θ_0, θ_1 separately, and still bowl shape over parameter pair (θ_0, θ_1) at same time.
E.g.: right 3D, high is J .

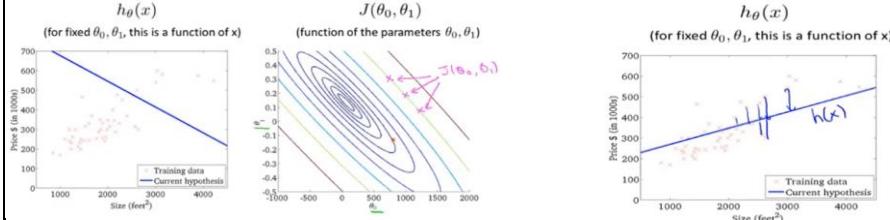
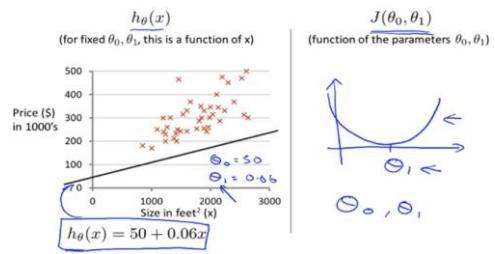


4.2 Contour plots Visualize cost function: with all parameters (θ_0, θ_1)

Contour plots: axis are parameter θ_0 and θ_1 , each of these oval/ellipses show a set of points that takes on the same value for cost $J(\theta_0, \theta_1)$

Like: Project the 3D bowl to 2D parameter space vertically, the middle small circle of bowl bottom, each oval corresponds to same high position in bowl (same cost J).

understanding: cost function value J how corresponds to different hypothesis h , and how better hypotheses may corresponds to points that are closer to the minimum of this cost function J .



5. Gradient descent_ Linear regression with one variable

Gradient descent used in all ML, and for cost function with as many number parameters.

5.1 Out line_gradient descent:

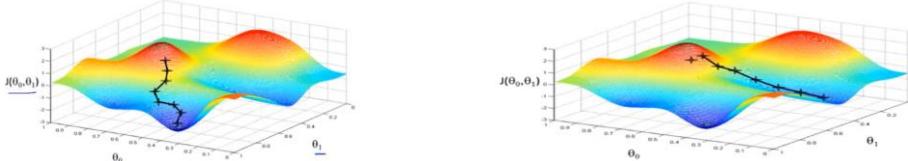
- start with initial parameter θ_0 : whatever value can be is ok
- keep changing parameters a little bit, to reduce cost function value $J(\theta)$, until hopefully end up at a minimum, or maybe a local minimum.

E.g:
standing at one point on the hill, and look all around (360 degree), find the best direction to take a little step downhill.
--> at new point on the hill, find a direction in order to take a little baby step downhill.
--> and so on, until converge to local minimum

Note: if take different initializing point:

repeat above step: take a little step in the direction of steepest descent, -->so on, may lead to the second local optimum, right pic

slightly different initial parameter may lead to very different local optimum.



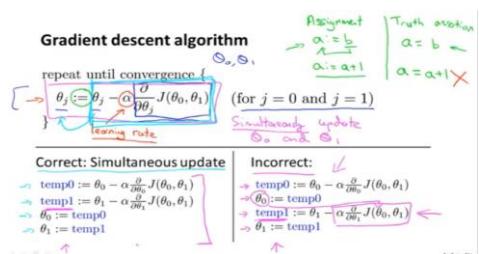
5.2: Gradient descent algorithm

'a:=b': assignment: take b and use it to overwrite whatever the value of a;

'a=b': true assertion/claiming-->1 true, 0 fail

a: learning rate: how big step take downhill with gradient descent. a large-->very aggressive gradient descent procedure, a small-->little baby steps downhill

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



more natural to implement the simultaneous update, incorrect update (unsimultaneous) in below might work , but have strange properties, not commonly use

e.g: in the right side, when calculating θ_1 , new updated θ_0 has been used in derivative of $J(\theta_0, \theta_1)$ over θ_1 .

Correct: Simultaneous update <pre> $\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$ $\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$ $\theta_0 := \text{temp0}$ $\theta_1 := \text{temp1}$ </pre>	Incorrect: <pre> $\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$ $\theta_0 := \text{temp0}$ $\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$ $\theta_1 := \text{temp1}$ </pre>
--	--

6. Gradient descent intuition_linear regression with one variable

1. Derivative $J'(\theta_i)$

$\theta = a * J'(\theta_i) \rightarrow$

1.1. $J'(\theta_i) > 0 \rightarrow \theta$ will decrease, direction to reduce J ;

1.2. $J'(\theta_i) < 0 \rightarrow \theta$ will increase, also direction to reduce J ;

$J'(\theta_i)$, make sure θ change in gradient descent , always in direction of reducing J .

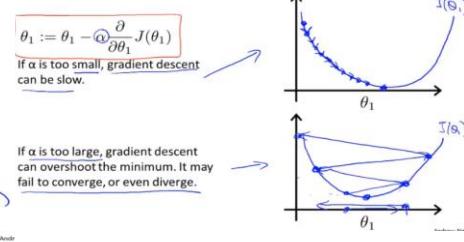
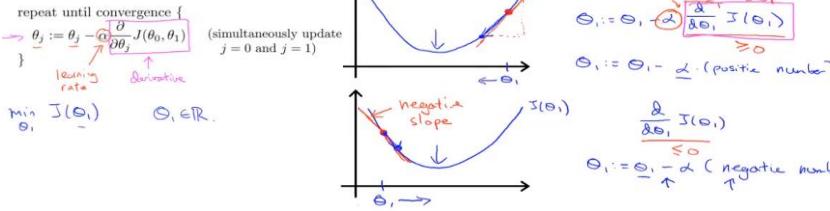
2. Learnign rate

- a too small, gradient descent can be slow.

To take a lot of mini baby steps to get to the minimum.

- a too large, gradient descent can overshoot the minimum. It **may fail to converge**, or even diverge.

Gradient descent algorithm



3. Gradient descent in local minimum

if parameter is already at a local minimum, one step of gradient descent actually does nothing, does not change parameter, $J(\theta_1) = 0$.

this is good as keep solution at local optimum-local optimum is good enough?

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

as approach a local minimum, gradient descent will automatically take smaller steps (derivative smaller approaching to 0). So, no need to decrease learning rate over time.

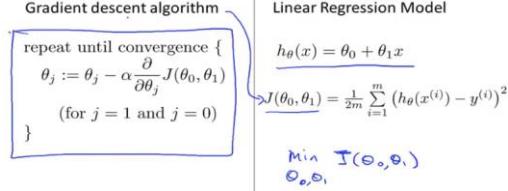
derivative will become smaller and smaller (less steep) as approaching to the local optimum which derivative is 0. --> gradient step: α *derivative become smaller.

when parameter take one step of gradient step, in the new position, the new gradient is less steeper.

Gradient descent can use to decrease any function J, not only linear regression.

7. Gradient descent for linear regression

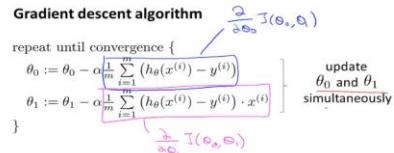
1. Linear regression cost function



$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{2}{2m} \cdot \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \\ &= \frac{2}{2m} \cdot \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2 \\ \theta_0, j = 0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ \theta_1, j = 1 : \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}\end{aligned}$$

2. Gradient descent algorithm

update all parameters θ_i simultaneously



Note: loss function definition: factor $1/2m$ --> for getting right derivative

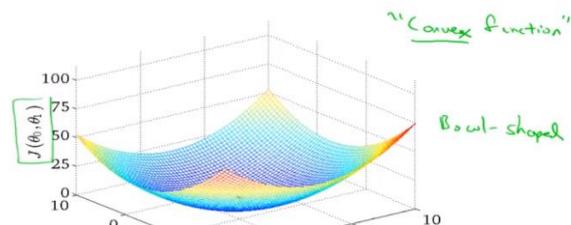
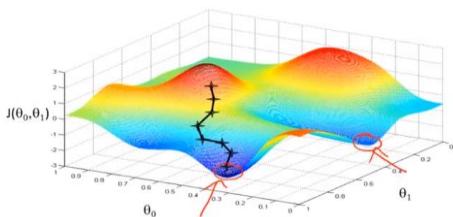
derivative defined: $1/m * \text{sum of each example loss } l(i) \text{'s derivative over parameter} = 1/m * [d_l(i) / d_{\theta_0}] \text{ sum over } i$ (so as to make overall cost of all examples decrease)
--> total cost function using factor $1/2m$ --> to get right derivative representation:

3. Convex function_linear regression function

susceptible to local optimum (容易陷入局部最优) : different initialization may end up different local optimum

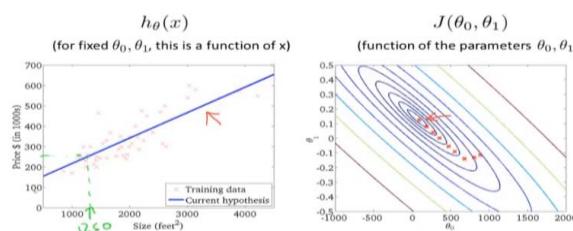
yet cost function of linear regression is always bowl-shaped-convex function: does not have any local optima, except one global optima.

Gradient descent for linear regression cost function, always converge to global optima



4. Linear regression cost function-gradient descent

global optima corresponds to best fit hypothesis h to training set.



5. 'Batch' gradient descent

above gradient descent also call batch descent, refer to the fact, in every step of gradient descent, we are looking at all of the training examples:
 calculate derivative: $J'(\theta)$. $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$
 go through all training examples in every step derivative calculation ($y^{(i)} - h_\theta(x^{(i)}) = h_\theta(x^{(i)}) - y^{(i)}$) --> J decrease in every step
 (note: if learning rate is reasonable--no overshooting)

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \\ &= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2 \\ \Theta_0 j = 0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ \Theta_1 j = 1 : \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}\end{aligned}$$

"Batch" Gradient Descent

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

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

6. Minimum cost function:

normal equation method: numerically solve the minimum issue of cost function.
 Yet gradient descent scale better to large data set, than normal equations method.

Note: derivative $dJ/d\theta$ (θ derivative vector) = $X \cdot (\Theta T \cdot X - Y)^T$
 $X = (x_1, x_2, \dots, x_m)$
 $Y = (y_1, \dots, y_m)$
 make $dJ/d\theta = 0 \rightarrow \text{get } \theta$

W3-Linear Algebra review (线性代数回顾)

3_14. Matrices and vectors _Linear Algebra review

1. Matrix: Rectangular array of numbers; use uper case for matrix (e.g.: A, B, ..)

Dimension of matrix: number of rows x number of columns

2. Matrix Elments:

A_{ij} = 'i,j entry' in i_th row, j_th colum

Matrix: Rectangular array of numbers:

$$\begin{bmatrix} 1402 & 191 \\ 1371 & 821 \\ 949 & 1437 \\ 147 & 1448 \end{bmatrix} \quad \begin{array}{l} \uparrow \quad \uparrow \\ 4 \times 2 \text{ matrix} \end{array} \quad \rightarrow \quad \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \quad \begin{array}{l} \uparrow \quad \uparrow \quad \uparrow \\ 2 \times 3 \text{ matrix} \end{array}$$

$\mathbb{R}^{2 \times 2}$

Dimension of matrix: number of rows x number of columns

Matrix Elements (entries of matrix)

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

A_{ij} = "i, j entry" in the ith row, jth column.

$A_{11} = 1402$

$A_{12} = 191$

$A_{32} = 1437$

$A_{41} = 147$

A_{12} = undefined (error)

3. Vector: n x 1 matrix (special matrix: only one column) ; use low case for vector (e.g. a, b c)

vector space: \mathbb{R}_n (n: vector dimension number)

4. Vector element: y_i = i_{th} element

5. Vector index: could start from 1 or 0

start from 1: $y = [y_1, y_2, \dots, y_n]$; (more common in math)

start from 0: $y = [y_0, y_1, \dots, y_{n-1}]$ (more common in learning algorithm)

unless otherwise specified:

use 1-index vector throughout rest of videos on linear algebra review.

use 0-index in machine learning applicaitons

Vector: An n x 1 matrix.

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix} \quad \leftarrow 4\text{-dimensional vector.} \quad \mathbb{R}^4$$

y_i = ith element

$y_1 = 460$

$y_2 = 232$

$y_3 = 315$

$y_4 = 178$

$\rightarrow [A, B, C, X]$

$\rightarrow [a, b, x]$

1-indexed vs 0-indexed:
 $y[1] = [y_1, y_2, y_3, y_4]$ $y = [y_0, y_1, y_2, y_3]$
 ↓ ↓ ↓ ↓
 1-indexed 0-indexed 1-indexed 0-indexed

Andrew Ng

3_15. Additon and scalar multiplication _Linear Algebra review

2.1 Matrix Addition

method: add up the elements of these matrices one at a time

$C = A + B \rightarrow$

$$C_{ij} = A_{ij} + B_{ij}$$

note: only same dimension matrices could add together.

Matrix Addition

$$\begin{array}{c} \downarrow \\ \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} + \begin{bmatrix} 4 & 0.5 \\ 2 & 5 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 5 & 0.5 \\ 4 & 10 \\ 3 & 2 \end{bmatrix} \\ \text{3x3} \quad \text{3x3} \quad \text{3x3} \end{array}$$

$$\begin{array}{c} \downarrow \\ \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} + \begin{bmatrix} 4 & 0.5 \\ 2 & 5 \\ 2 & 5 \end{bmatrix} = \text{error} \\ \text{3x3} \quad \text{3x3} \quad \text{3x3} \end{array}$$

Scalar Multiplication

$$\begin{array}{c} \downarrow \text{ real number} \\ 3 \times \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 6 & 15 \\ 9 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} \times 3 \end{array}$$

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

Andrew Ng

Combination of Operands

$$\begin{array}{c} \downarrow \text{ Scalar multiplication} \\ \begin{bmatrix} 3 \times \begin{bmatrix} 1 & 0 \\ 4 & 2 \end{bmatrix} + \begin{bmatrix} 0 & 5 \\ 0 & 5 \end{bmatrix} - \begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix} / 3 \end{bmatrix} \quad \text{Scalar division} \\ = \begin{bmatrix} 1 & 0 \\ 12 & 5 \end{bmatrix} + \begin{bmatrix} 0 & 5 \\ 0 & 5 \end{bmatrix} - \begin{bmatrix} 1 & 0 \\ 0 & \frac{2}{3} \end{bmatrix} \quad \text{matrix subtraction / vector subtraction} \\ = \begin{bmatrix} 2 & 5 \\ 12 & 5 \end{bmatrix} \quad \text{matrix addition / vector addition} \\ = \begin{bmatrix} 2 & 5 \\ 12 & 5 \end{bmatrix} \quad \text{scalar multiplication} \end{array}$$

2.3: Combination of operands.

3_16. Matrix-vector multiplication _Linear Algebra review

$3 \times 5 \times 2$ $3 \times (5 \times 2) = (3 \times 5) \times 2$ "Associative"

$3 \times 10 = 30 = 15 \times 2$

$A \times B \times C$

Let $D = B \times C$. Compute $A \times D$.

Let $E = A \times B$. Compute $E \times C$.

$(A \times B) \times C$

$A \times (B \times C)$

Some answer.

3. Identity Matrix: I or $I_n \times n$

special matrix, similar to real number 1: $1 \times z = z \times 1 = z$

Identity matrix property: 'comutative'

for any matrix A: $A (m,n) \times I (n,n) = I (m,m) \times A (m,n) = A (m,n)$

Identity Matrix 1 is identity. $1 \times z = z \times 1 = z$

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

Examples of identity matrices:

$$\begin{matrix} 1 \times 1 & \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \\ 2 \times 2 & \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \\ 3 \times 3 & \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \\ 4 \times 4 & \end{matrix}$$

Informally:

$$\begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix}$$

For any matrix A,

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

$m \times n \quad m \times n \quad m \times n \quad m \times n$

$I_{n \times n}$

Note: $AB \neq BA$ in general

$AI = IA \checkmark$

3_19. Matrix inverse and transpose _Linear Algebra review

1. Matrix Inverse

if A has an inverse: $A * A^{-1} = A^{-1} * A = I$

Singular / Degenerate matrices: matrices that do not have an inverse.
like zero matrix.

- Not all numbers have an inverse, not all matrix have an inverse;
- Only square matrix, have possibility to have inverse
- Many open source could use to compute matrix inverse
e.g: In Octave, $\text{pinv}(A) \rightarrow \text{inverse } A$.

1 is "identity" $3 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = 1$ $12 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = 1$

Not all numbers have an inverse.

Matrix inverse: square matrix ($n \times n = n \times n$) A^{-1}

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

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

$\begin{matrix} 3 & 4 \\ 2 & 1 \end{matrix} \begin{matrix} 1 & -0.4 \\ -0.5 & 0.25 \end{matrix} = \begin{matrix} 1 & 0 \\ 0 & 1 \end{matrix} = I_{2 \times 2}$

Matrices that don't have an inverse are "singular" or "degenerate".

2. Matrix Transpose A (m,n),

Method:

change each row to column position; take 45 degree axis-go through first value, and then flipping matrix along this 45 degree axis

Transpose $B = A^T$

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

Matrix Transpose

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

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

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

Then B is an $n \times m$ matrix, and

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

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

$B_{32} = 9 \quad A_{23} = 9$.

W5 - Linear Regression with multiple variables

5_27. Multiple features_Linear Regression with multiple variables

Multiple variable: have more than one features use to predict output y .

1. Notation: Multiple features

n =number of features

$x(i)$ = input (feature vector) of i _th training example

$x(i)_j$ = value of feature j in i _th training example

Multiple features (variables).

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

Notation:
 $\rightarrow n = \text{number of features}$
 $\rightarrow x^{(i)} = \text{input (features) of } i^{\text{th}} \text{ training example.}$
 $\rightarrow x_j^{(i)} = \text{value of feature } j \text{ in } i^{\text{th}} \text{ training example.}$

Hypothesis:

Previously: $h_{\theta}(x) = \theta_0 + \theta_1 x_1$

$$h_{\theta}(x) = \hat{\theta}_0 + \hat{\theta}_1 x_1 + \hat{\theta}_2 x_2 + \hat{\theta}_3 x_3 + \hat{\theta}_4 x_4$$

$$\text{e.g. } h_{\theta}(x) = \underline{\underline{\theta}}_0 + \underline{\underline{\theta}}_1 x_1 + \underline{\underline{\theta}}_2 x_2 + \underline{\underline{\theta}}_3 x_3 + \underline{\underline{\theta}}_4 x_4$$

$$\rightarrow h_{\theta}(x) = \underline{\underline{\theta}}_0 + \underline{\underline{\theta}}_1 x_1 + \underline{\underline{\theta}}_2 x_2 + \dots + \underline{\underline{\theta}}_n x_n$$

For convenience of notation, define $\underline{\underline{x}}_0 = 1$ ($x_0^{(i)} = 1$)

$$\begin{aligned} x &= \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} & \underline{\underline{\theta}} &= \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1} \\ h_{\theta}(x) &= \underline{\underline{\theta}}_0 x_0 + \underline{\underline{\theta}}_1 x_1 + \dots + \underline{\underline{\theta}}_n x_n & \underline{\underline{\theta}}^T &= [\theta_0, \theta_1, \dots, \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \\ &= \underline{\underline{\theta}}^T \underline{\underline{x}} \end{aligned}$$

Multivariate linear regression.

2. Hypothesis:

sum of all features:

$$h_{\theta}(x_i) = \theta_0 + \theta_1 x_{i,1} + \theta_2 x_{i,2} + \theta_3 x_{i,3} + \theta_4 x_{i,4} + \dots + \theta_n x_{i,n}$$

for convinient take $x(i)_0=1$ -->

$$X(i) = [x(i)_0, x(i)_1, \dots, x(i)_n]^T$$

$$\theta = [\theta_0, \theta_1, \theta_2, \dots, \theta_n]^T$$

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

5_28. Gradient descent for multiple variables_Linear Regression with multiple variables

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

Parameters: $\theta_0, \theta_1, \dots, \theta_n$ $\underline{\underline{\theta}}$ $n+1$ -dimensional vector

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

Gradient descent:

$$\begin{aligned} \text{Repeat } \{ & \\ \rightarrow \theta_j &:= \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n) & \boxed{J(\theta)} & \\ \} & \quad (\text{simultaneously update for every } j = 0, \dots, n) \end{aligned}$$

Gradient descent : feature $n=1$ vs $n>1$

feature number $n=1$ ----> parameter θ_0, θ_1

feature number $n>1$ ----> parameter $\theta_0, \theta_1, \dots, \theta_n$:

derivative of θ_0, θ_1 , formula same as $n=1$, (only now h function involve more features).
Same rule could apply to other parameter derivative.

Gradient Descent

Previously ($n=1$):

Repeat {

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

$$\boxed{\frac{\partial}{\partial \theta_0} J(\theta)}$$

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

(simultaneously update θ_0, θ_1)

}

$$\begin{aligned} \text{New algorithm } (n \geq 1): & \\ \text{Repeat } \{ & \\ \rightarrow \theta_j &:= \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} & \boxed{\frac{\partial}{\partial \theta_j} J(\theta)} \\ & \quad (\text{simultaneously update } \theta_j \text{ for } j = 0, \dots, n) & \boxed{\theta_j} & \\ \} & \\ \rightarrow \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} & \boxed{\theta_0} & \\ \rightarrow \theta_1 &:= \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)} & \boxed{\theta_1} & \\ \rightarrow \theta_2 &:= \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)} & \boxed{\theta_2} & \end{aligned}$$

5_29. Gradient descent in practice: feature scaling_Linear Regression with multiple variables

1. Skewed elliptical shape cost function: feature scale differe large

E.g: only two features for examples, x_1 (0-2000 range), x_2 (1-5 range);

-->Cost function $J(\theta)$ is very skewed elliptical shape (2000 to 5 ratio)

-->if run gradient descent on this cost function, may end up taking a long time and can oscillate back and forth and take a long time before it can finally find its way to the global minimum.

more skewed cost function, gradient is more harder taking it's way, meandering around, take a long time to find it's way to the global minimum.

(note: learning rate is same for every parameter, -->feature range differ, corresponding weights also differ)

-->each weight gradient step: a * derivative:

when weights range differ big-->weight derivative differ big-->gradient step differ big-->some weight update towards global minimum, while others may overshooting-->total cost may decrease or increase : cost oscillate

note: if learning rate is unique for each weight -->no need feature scaling, will make cost decrease in every step.)

Feature Scaling

Idea: Make sure features are on a similar scale.

E.g. x_1 = size (0-2000 feet²)

x_2 = number of bedrooms (1-5)

$$\theta_2$$

$$\theta_1$$

$$\theta_0$$

$$J(\theta)$$

$$x_2$$

$$x_1$$

$$x_0$$

$$J(\theta)$$

$$\theta_2$$

$$\theta_1$$

$$\theta_0$$

$$J(\theta)$$

$$x_2$$

$$x_1$$

$$x_0$$

$$J(\theta)$$

2. Feature scaling: feature / value range

Idea: Make sure features are on a similar scale:

different feaure take on similar ranges of values, then gradient descent can converge more quickly.

By feature scaling, contur of cost function can become much less skewed, **conturs may looks like more circles**.

-->run gradient descent on cost function like this, could find **more direct path to the global minimum (weight updating less overshooting: cost less oscillate with iteration steps)**, rather than taking a much more convoluted path, trying to follow a much more complicated trajectory, to go the global minumum.

3. Feature scaling rule

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

E.g:
 $x_1 < 1$, this range is ok, close to $[-1, 1]$;
 $-2 < x_2 < 0.5$, this range is ok, close to $[-1, 1]$;
 $-100 < x_3 < 100$, this range too big, no ok,
 $-0.0001 < x_4 < 0.0001$, this range too small, no ok

Note:

1. Normally below scale is ok: e.g : $[-3, 3], [-1/3, 1/3]$;
2. Do not require all features are on exactly same range, but as long as they are all close enough to this range $[-1, 1]$, **gradient descent** should work ok.

(note: scalling for gradient decent : -->weight updating more efficiently, leading cost to minimum)

Feature Scaling

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

$$\begin{aligned} x_0 &= 1 \\ 0 \leq x_1 \leq 3 &\quad \checkmark \\ -2 \leq x_2 \leq 0.5 &\quad \checkmark \\ -100 \leq x_3 &\quad \text{100} \quad \times \\ -0.0001 \leq x_4 &\quad \text{0.0001} \quad \times \end{aligned}$$

$-3 \rightarrow 3 \quad \checkmark$
 $-\frac{1}{3} \rightarrow \frac{1}{3} \quad \checkmark$

Mean normalization

Replace x_i with $\frac{x_i - \mu_i}{s_i}$ to make features have approximately zero mean
(Do not apply to $x_0 = 1$).

$$\begin{aligned} \text{E.g. } x_1 &= \frac{\text{size} - 1000}{2000} & \text{Range: } s_{\text{size}} \approx 100 \\ x_2 &= \frac{\#bedrooms - 2}{5} & \text{1-5 bedrooms} \\ (-0.5 \leq x_1 \leq 0.5) &\quad (-0.5 \leq x_2 \leq 0.5) \\ x_1 &\leftarrow \frac{x_1 - \mu_1}{s_1} & \text{avg value} \\ \text{range } (\max - \min) &\quad \text{in training set} \\ \text{(or Standard deviation)} & \quad | \quad x_2 \leftarrow \frac{x_2 - \mu_2}{s_2} \end{aligned}$$

4. Scaling method: mean normalization

except divide feature range size, could also use mean normalization

method:

$x(i)_j := x(i)_j - \bar{x}_j / s_j$
replace $x(i)_j$ with $x(i)_j - \bar{x}_j$, to make features have approximately **zero mean** (do not apply to $x(i)_0=1$)

\bar{x}_j : average of feature j in the training sets : average of all i ; $x(i)_j \cdot x(i)$: i example in training set.

s_j : range of feature j in the training set, max - min: $\max x(i)_j - \min x(i)_j$

Note: s_j could also use deviation of j feature in trainin set. here max -min to get range is also ok.

Summary:

any value get feature into anything close to these sorts of ranges $[-1,1]$ will be fine.

Feature scaling does not have to be extract, in order to get gradient descent to run quite a lot faster (and converge in a lot fewer iterations).

5_30. Gradient descent in practive II: learning rate_Linear Regression with multiple variables

1. 'Debugging' : make sure gradient descent is working correctly

Plot draw: cost function after each iteration step in gradient descent (not pramter θ)

1.1. after certain steps of iteration, -->get new parmater θ -->calculate cost function J of θ .

1.2. Cost function should descent after every iterations.

give warning if gradient descent does not work correctly.

1.3. plot draw could help judge whether or not gradient descent has converged or not.

from this plot could find, after certain time steps, cost functioin does not go down much-like has flattened out here. Seems gradient descent kind of converged because cost function isn't going down much more.

Note:

1. converge steps for different applications could be big different.

2. difficult in advance predict how many steps needed to converge, **usually using plot draw of cost function with iterations to see if converged or not**.

(note: plot draw of cost function with iterations:
-->1. debug: if gradient descent works;
-->2. check if converged or not
-->3. choosing learning rate
)

2. Automatically Converge test

another method except plot draw to check if gradient descent converged or not.

Idea: declare convergence if cost function decreases by less than some small value ϵ in one iteration.

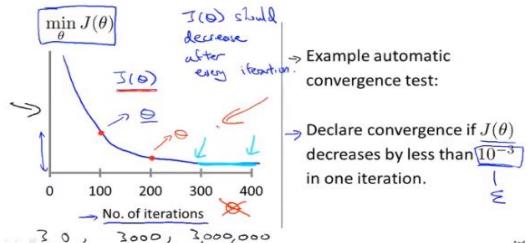
But choosing threshold ϵ is pretty difficult, prefer to look at plot-draw of cost function with iterations rather than rely on automatically converge test.

Gradient descent

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

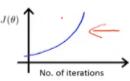
- "Debugging": How to make sure gradient descent is working correctly.
- How to choose learning rate α .

Making sure gradient descent is working correctly.



3. Uncorrect gradient descent

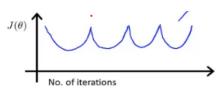
3.1: cost function single increase



If in plot-draw of cost function vs iterations, cost functions singal grow like above pic, then usually is learning rate too big:-->use smaller learning rate if no bug in coding
if learning rate too big, gradient descent may overshoot the minimum to a new position and then overshoot again to another position , keep on overshooting if learning rate too big, then end up getting worse and worse, getting higher value of cost function of θ .

3.2: cost function : go down for a while and then go up , then go down for a while and go up again.

-->use smaller learning rate also.



Summary:

- if learning rate too small: slow convergence;
- if learning rate too large: cost function J may not decrease on every iteration; may not converge.
Slow convergence also possible.

-->use plot draw of cost function with iterations, to figure out what's wrong.

To choose learning rate, try: ..., 0.001, 0.01, 0.1, 1, ...

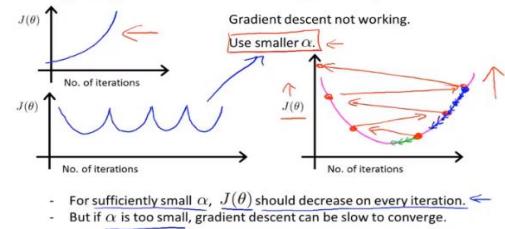
Recommend:

run with different learning rates (rafty 3 times bigger of previous value, till find the possible/reasonable largest value and smallest value), plot cost function vs iterations, and choose one learning rate that seems to be causing cost function decrease rapidly.

(note: cost function decrease rapidly:

decrease: learning rate should be reasonably small
rapidly: learning rate should be reasonably large
)

Making sure gradient descent is working correctly.



Andrew Ng

Note:

1. Mathematically proven, if learning rate is sufficiently small, cost function J decrease on every iteration.
So if gradient descent does not decrease, try smaller learning rate.
2. if learnig rate too small, gradient descent can be slow to converge.

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)

To choose α , try

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

5_31. Features and polynomial regression_Linear Regression with multiple variables

sometimes by defining new features , may get better model.

1. Polynomial regression:

is the idea of choosing features

e.g. 1.1 Reduce feature number by creating new feature:

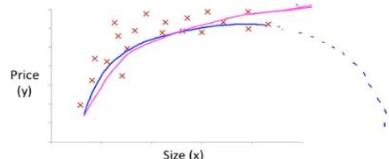
predict house price: $x_1 = \text{width}, x_2 = \text{length}$ $h = \theta_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2$
could create new feature x_1 (area)= width * length----> $h = \theta_0 + \theta_1 \cdot x_1$

e.g 1.2: build linear regression by creating new feature

predict house size: $x_1 = \text{size}$, to better fit training data, $h = \theta_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_1^2 + \theta_3 \cdot x_1^3 + \dots$
creat new feature: $x_2 = x_1^2, x_3 = x_1^3, \dots \rightarrow h = \theta_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 + \theta_3 \cdot x_3$

Note: need to take care of feature scaling after creating new feature . in above example 2, feature x_1, x_2, x_3 may has big range different, need to scaling

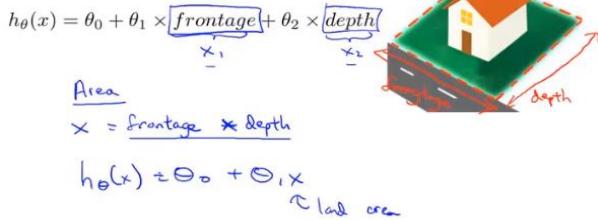
Choice of features



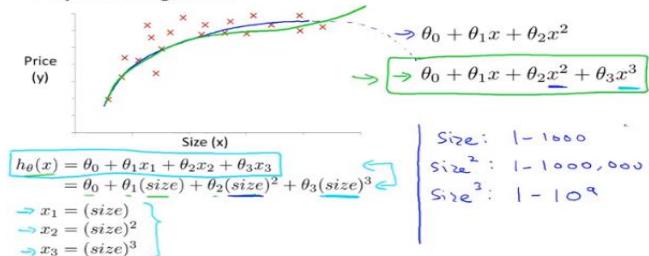
$$\begin{aligned} h_{\theta}(x) &= \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2 \\ h_{\theta}(x) &= \theta_0 + \theta_1(\text{size}) + \theta_2\sqrt{(\text{size})} \end{aligned}$$

K ↗

Housing prices prediction



Polynomial regression



Andrew Ng

Summary:

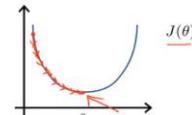
1. Have choice what features to use;
2. By defining different features, can fit more complex functions to data

5_32. Normal equation_Linear Regression with multiple variables

1. Normal equation: method to solve for parameter analytically

while gradient descent come to minimum by many gradient descent steps, normal equation solve optimal value for θ all at one go-in one step, get to the optimal value.

Gradient Descent



Normal equation: Method to solve for θ analytically.

2. Normal equation intuition:

set partial derivative of cost function over each parameter = 0, --> get parameter that minimize cost function J .

Note:

2.1 parameter dimension is $n+1$ (n: input feature number)

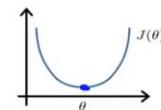
2.2: need go through all partial derivative of J over each parameter, and calculate separately.
--> too long and too involved,

Intuition: If 1D ($\theta \in \mathbb{R}$)

$$\rightarrow J(\theta) = ab^2 + b\theta + c$$

$$\frac{\partial}{\partial \theta} J(\theta) = \dots \stackrel{\text{set } 0}{=} 0$$

Solve for θ



$$\theta \in \mathbb{R}^{n+1} \quad J(\theta_0, \theta_1, \dots, \theta_m) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots \stackrel{\text{set } 0}{=} 0 \quad (\text{for every } j)$$

Solve for $\theta_0, \theta_1, \dots, \theta_n$

3. Normal equation calculation

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

$$x^{(i)} (n+1, 1) = [x^{(i)}_0, x^{(i)}_1, \dots, x^{(i)}_n]^T; \quad X (m, n+1) = [x^{(1)}_0, x^{(2)}_0, \dots, x^{(n)}_0]; \quad y (m, 1) = [y_1, y_2, y_3, \dots, y_m]^T$$

m examples $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$; n features.

$$\begin{matrix} \begin{bmatrix} x_0^{(1)} \\ x_1^{(1)} \\ x_2^{(1)} \\ \vdots \\ x_n^{(1)} \end{bmatrix} & \in \mathbb{R}^{n+1} \end{matrix} \times \begin{matrix} \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \end{matrix} = \begin{matrix} \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \end{matrix}$$

E.g. If $x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$ $\times = \begin{bmatrix} 1 & x_1^{(i)} \\ 1 & x_2^{(i)} \\ \vdots & \vdots \\ 1 & x_n^{(i)} \end{bmatrix} \times \begin{bmatrix} y_1^{(i)} \\ y_2^{(i)} \\ \vdots \\ y_m^{(i)} \end{bmatrix} = \begin{bmatrix} y^{(i)} \end{bmatrix}$

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

Examples: $m = 4$.

	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

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

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

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

4. Calculation in Octave:

pinv : calculate inverse

Note: feature scaling no necessary in normal equation, as feature scaling is for faster converge in gradient descent

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

$(X^T X)^{-1}$ is inverse of matrix $X^T X$.

$$\text{Set } A = X^T X$$

$$(X^T X)^{-1} = A^{-1}$$

Octave: $\text{pinv}(X^T X) * X^T y$

$$\text{pinv}(X^T X) * X^T y$$

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

$$X' \quad X^T$$

Feature scale
 $0 \leq x_1 \leq 1$
 $0 \leq x_2 \leq 1000$
 $0 \leq x_3 \leq 10$

Andrew Ng

5. Comparison with gradient descent

m training examples, n features

Gradient descent:

- need to choose learning rate;
- needs many iteration;
- works well even when n is large; still reasonably efficient.

Normal equation:

- No need to choose learning rate;
- Do not need to iterate;
- Need to compute $(X^T X)^{-1}$: $X^T X$ nxn dimension, computing time $O(n^3)$
- Slow if n is very large;

Note: if feature number n = 100, 1000--->normal equation;
n=10000, may wonder which to use, may start to use gradient descent.

m training examples, n features.

Gradient Descent

- Need to choose α .
- Needs many iterations.
- Works well even when n is large.

$$n = 10^6$$

Normal Equation

- No need to choose α .
- Don't need to iterate.
- Need to compute $(X^T X)^{-1}$ $n \times n$ $O(n^3)$
- Slow if n is very large.

$$n = 1000$$

$$n = 10000$$

$$\dots$$

$$n = 100000$$

summary:

1. so long as the **number of features** is not too large, normal equation gives us a great alternative method to solve for the parameter θ .

So long as feature number is less than 1000, use normal equation.

2. for more complicated algorithm, like classification algorithm, linear regression algorithm, normal equation actually do not work for those more sophisticated learning algorithm, have to resort to gradient descent.

logistic regression could use same equation, as cost J partial derivative is same as square error cost of linear algorithm.??
for this specific model of linear regression, normal equation can give an alternative that can be much faster than gradient descent.

5_33. Normal equation and non-invertibility_Linear Regression with multiple variables

$X^T X$: issue of being non-invertible should happen pretty rarely.

1. Method for non-invertible matrix $X^T X$

- Octave: $\text{pinv}(X' * X)$: actually do right thing..

Pinv: pseudo inverse (伪逆) : could calculate the parameter θ wanted even if $X^T X$ is non-invertible.

inv: inverse: 逆

(note: $X^T X$: invertible --> $r(X) = \min(n, m) = n$ -->

1. $m > n$;

2. $r(X) = n$: feature x_1, \dots, x_n 线性无关

$X = (x(1), \dots, x(m)) = U \Lambda V^T$

>> U: nxn: 单位正交矩阵

>> Λ : nxm: 对角矩阵, 对角值为 X 的特征值, 从大到小排序;

>> V: nxn: 单位正交矩阵

$X^T X = U \Lambda^2 U^T$, U^T : --> invertible --> $\Lambda^2 = \Lambda$, Λ^T : nxn: 对角矩阵, 特征值的平方, invertible -->

X n个特征值不为0 --> $r(X) = n$

)

Normal equation

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

$$X^T X$$

- What if $X^T X$ is non-invertible? (singular/degenerate)

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



2. Causes for non-invertible $X^T X$

2.1; Redundant features (linearly dependent)

e.g.: $x_1 = \text{size in feet}^2, x_2 = \text{size in m}^2$.

2.2 Too many features ($m \leq n$): try to find $n+1$ parameter to fit m example

- delete some features or use regularization (fit lot of parameters into small training data)

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

- Redundant features (linearly dependent).

E.g. $x_1 = \text{size in feet}^2$ $1m = 3.28 \text{ feet}$
 $x_2 = \text{size in m}^2$ $1m^2 = 10.76 \text{ feet}^2$
 $x_1 = (3.28)^2 x_2$ $\rightarrow m = 10 <$
 $\rightarrow n = 100 <$

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

- Delete some features, or use regularization.

$$\Theta \in \mathbb{R}^{100 \times 1}$$

\downarrow 100

Summary: if $X^T X$ is non-invertible:

1. check feature: if there is redundant features, delete it until no redundant features.: feature 线性无关

2. Then check if have too many features ($m < n$): delete features if bare to use fewer features, or consider use regularization: $n < m$

6_36. Basic operations_Octave Tutorial

4.6 Octave programming environment

use in ML programming;

Octave: free open source software.

Use Octave or Matlab many learning algorithms become just a few lines of code to implement.

Normal process:

1. First prototype learning algorithm/ software in Octave;

Because software in Octave makes it incredibly fast to implement these learning algorithms.

2. Only after gotten it to work, then might migrate it to C++ or Java or whatever.

---> In this way much faster than if starting out in C++.

Matlab works well too, but too expensive.

Recommend to do exercises in this course in Octave

W7-Logistic Regression

7_44. Classification_Logistic Regression

1. Linear regression for classification

e.g:

$h_{\theta}(x) >= 0.5$, predict 'y=1';
 $h_{\theta}(x) < 0.5$, predict 'y=0'.

when training data is small and close, linear regression seems ok for classification.
but when training data is larger, like right pic. right point, linear regression changes to fit all data (for minimizing distance of prediction & training set), threshold x corresponding to $h(x)=0.5$ change to larger ->classification for smaller x does not work.

- Linear regression use for classification task is not good idea;
- Linear regression prediction $h_{\theta}(x)$ can be > 1 or < 0 , even training data is 0 or 1.

question: logistic regression decision boundary is $\theta^T x = 0$ -> still linear boundary, how could fit right example better?

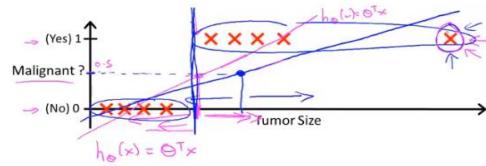
->linear boundary: as long as meet boundary spec, does not matter beyond how much.-->output is 1, cost is 0.

Linear regression: when used in classification problem, not only need to meeting the prediction line, but also the distance to prediction line is cost, need to be small.-->**prediction line forced to be close to each training example.**

while logistic model **only meet linear boundary, cost is 0** cost is $-\log(p) - \log(1-p)$

-->different cost function/optimizing object-->learn different decision boundary

(eg: take $\theta=(\theta_0, \theta_1)$)-->plot labeled y vs x (only one feature x_1)
 \Rightarrow logistic classification: optimizing object: $\theta_0 + \theta_1 x_1 > 0$ for $y = 1$
 \Rightarrow Linear regression: optimizing object: $\theta_1 x_1 + \theta_0$: close to labeled y 1/0
-->get different decision boundary



→ Threshold classifier output $h_{\theta}(x)$ at 0.5:

- If $h_{\theta}(x) \geq 0.5$, predict "y = 1"
- If $h_{\theta}(x) < 0.5$, predict "y = 0"

Classification: $y = 0$ or 1
 $h_{\theta}(x)$ can be > 1 or < 0
Logistic regression: $0 \leq h_{\theta}(x) \leq 1$
Classification

Summarize:

classification problem: variable y is discrete value.
Classification: positive or negative definition does not matter that much.

Two-class problem: result only two, 1, or 0;
Multi-class problem: result more than 2.

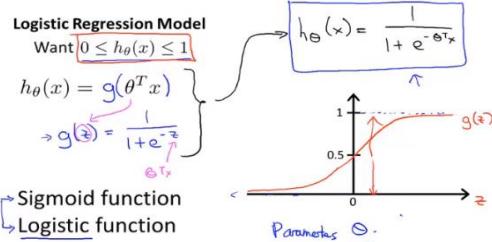
7_45. Hypothesis Representation_Logistic Regression

1. Logistic regression model: want $0 \leq h_{\theta}(x) \leq 1$

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

sigmoid function / logistic function: $g(z) = 1 / (1 + e^{-z})$

use this function, fit parameters θ to training data. Training data-->choose parameter-->use sigmoid/logistic function make predictions.



2. Interpretation of hypothesis Output

$h_{\theta}(x)$ = estimated probability that $y=1$ on input X

$h_{\theta}(x) = P(y=1 | x; \theta)$: 'probability that $y=1$, given x , parameterized by θ '

as only have two result $y=1, y=0$, $\Rightarrow P(y=0 | x; \theta) = 1 - P(y=1 | x; \theta)$

Interpretation of Hypothesis Output
 $h_{\theta}(x)$ = estimated probability that $y=1$ on input x
Example: If $x = [x_0 \ x_1] = [1 \ tumorSize]$
 $h_{\theta}(x) = 0.7$ $y=1$
Tell patient that 70% chance of tumor being malignant

$$h_{\theta}(x) = P(y=1 | x; \theta) \quad \text{"probability that } y=1 \text{, given } x, \text{ parameterized by } \theta"$$

$$y = 0 \text{ or } 1$$

$$\Rightarrow P(y=0 | x; \theta) + P(y=1 | x; \theta) = 1$$

$$\Rightarrow P(y=0 | x; \theta) = 1 - P(y=1 | x; \theta)$$

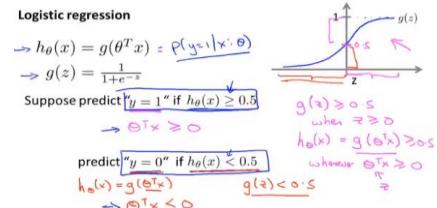
7_46. Decision boundary_Logistic Regression

1. Decision intuition

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

suppose predict :

$y=1$ if $h_{\theta}(x) >= 0.5$ --> $\theta^T x \geq 0$
 $y=0$ if $h_{\theta}(x) < 0.5$ --> $\theta^T x < 0$

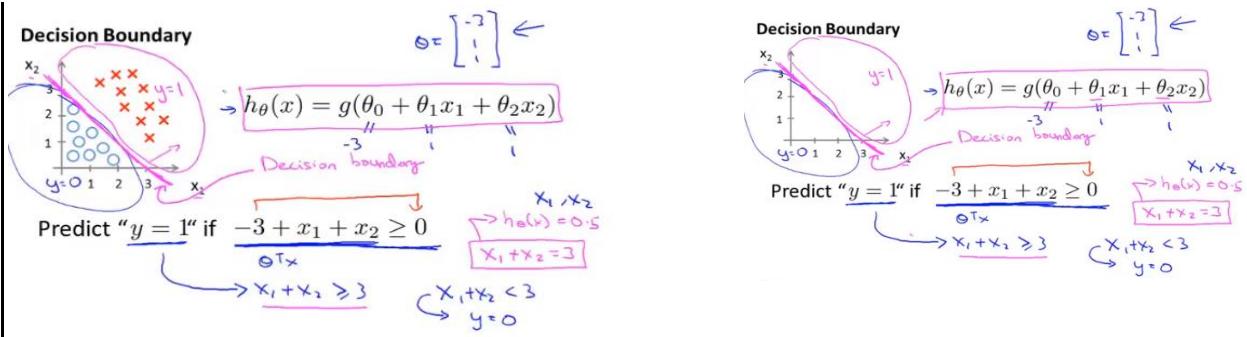


2. Decision Boundary_Linear

decision boundary: for chosen parameters θ : take $\theta^T x = 0$, this line is decision boundary.

Note:

-Decision boundary is property of $h_{\theta}(x) = g(\theta^T x)$ and parameter, not training data.
once parameter is chosen, decision boundary is defined. training data is not used to define decision boundary, but may be used to fit the parameters θ



3. Decision Boundary_Non-linear

when have right pic. Training data, could:

add extra higher order polynomial terms to the features, Like In linear regression.

Once parameter is chosen, make take $\theta^T \cdot x = 0$; this line is the decision boundary

Note:

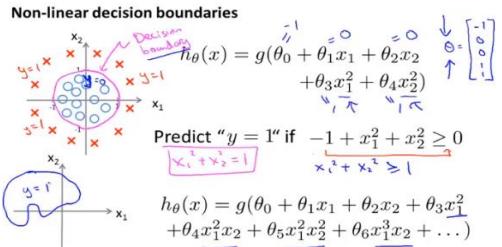
-Decision boundary is property of $h_\theta(x) = g(\theta^T \cdot x)$ and parameter, not training data. once parameter is chosen, decision boundary is defined. training data is not used to define decision boundary, but may be used to fit the parameters θ

Linear regression for classification: could not get good decision boundary.

maybe linear regression also get good decision boundary, by adding new feature??

(no: as optimizing object is different; linear regression forced decision boundary to be close to training set;

while logistic regression optimizing object : is maximize probability)



7_47. Cost function_Logistic Regression

Fit parameter for logistic regression; define cost function used to fit parameters.

1. Cost function _linear cost function?

define cost $(h_\theta(x), y) = 1/2 * (h_\theta(x) - y)^2$

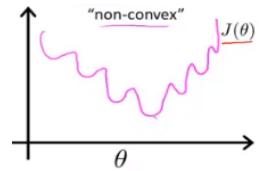
Linear regression: $J(\theta) = 1/m \text{ (sum of cost } (h_\theta(x_i), y_i) \text{ over all examples)}$

cost $(h_\theta(x), y) = 1/2 * (h_\theta(x) - y)^2$,

does not work well for logistic regression -> cost function J is non-convex (cause by intro. of sigmoid function), gradient descent could not guarantee convergence. Need find another kind of cost function '_convex shape', so to guarantee gradient descent would converge to the global optima.

(note: cost function purpose is to make gradient step/weight updating step bigger when prediction is different from label;

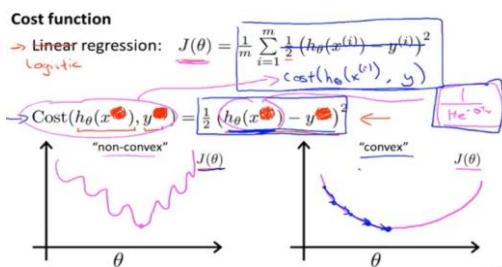
e.g: linear regression: square error: --> weight derivative $dJ/d\theta = (h_\theta(x) - y) \cdot x$, will be bigger if h difference big from label y if logistic regression, use square error: cost = $(h_\theta(x) - y)^2$, --> $dJ/d\theta = (1-h) \cdot h \cdot (h-y) \cdot x$: when $h-y$ big, derivative may not big)



Note: for one signal example

if use squared error--> $d\theta = (h-y) \cdot h(1-h) \cdot x$: when prediction h is far from labeled y , derivative θ is not big due to factor: $h \cdot (1-h)$

if use log: $I = -(y \log(h) + (1-y) \log(1-h))$ --> $d\theta = (y-h) \cdot x$, when prediction h is far from labeled y , derivative θ is also big.



Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

m examples $x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$ $x_0 = 1, y \in \{0, 1\}$

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

How to choose parameters θ ?

2. Cost function_Logistic regression

Cost $(h, y) =$

if $y=1$, cost = $-\log(h)$; (maximum probability)

if $y=0$, cost = $-\log(1-h)$;

2.1: if $y=1$, cost = $-\log(h)$

if $y=1$, $h=1$ --> cost=0:

$h=1 \rightarrow (y=1 | x, \theta)=1$, match label data ' $y=1$ ', cost =0.

if $y=1$, $h=0$, cost -->huge:

$h=0 \rightarrow P(y=1 | x, \theta)=1$, does not match label data ' $y=1$ ', cost =huge. penalize learning logarithm by a very large cost.

2.1: if $y=0$, cost = $-\log(1-h)$

if $y=0$, $h=1$ --> cost=huge:

$h=1 \rightarrow P(y=0 | x, \theta)=1$, does not match label data ' $y=0$ ', cost =huge. penalize learning logarithm by a very large cost.

if $y=0$, $h=0$, cost =0:

$h=0 \rightarrow P(y=0 | x, \theta)=1$, match label data ' $y=0$ ', cost =0.

(note: cost function define:

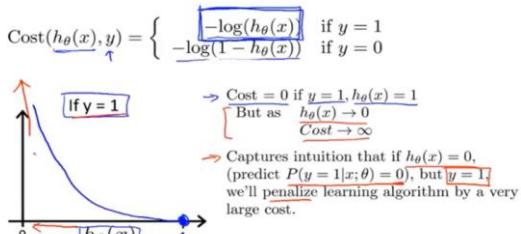
object -->make derivative $dJ/d\theta = (h-y) \cdot x$

--> $dJ/dh \cdot dh/dz \cdot dz/d\theta = dJ/dh \cdot h \cdot (1-h) \cdot x = (h-y) \cdot x$

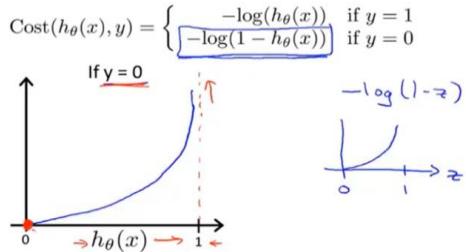
--> $dJ/dh = (h-y)/h/(1-h)$

--> $J = - (y \log(h) + (1-y) \log(1-h))$

Logistic regression cost function



Logistic regression cost function



7_48. Simplified cost function and gradient descent_Logistic Regression

Note: $y=0$ or 1 always (in training data)

1. Simplified cost function

Purpose:

combine into one equation more convenient to write out cost function and derive gradient descent

single example cost:

$$\text{cost_i} = -y \log(h) - (1-y) \log(1-h)$$

for each single example i:

$$\text{cost_i} = -\log(h) \text{ if } y=1, =-\log(1-h) \text{ if } y=0;$$

the combined equation is same for single example cost_i;

for all training set, cost

cost = sum(cost_i), since single cost_i is same whether in combined equation or separated equation-->overall cost also same in these two methods.

Note:

- could use this combination, key point is training data y is always 0 or 1.
- factor used here is $1/m$:
- $J=1/m(\text{sum of cost}_i)$
- for linear regression, $\text{cost}_i=1/2*(h-y)^2$
- for logistic regression, $\text{cost}_i=y \log(h) - (1-y) \log(1-h)$

- Use this function found parameter minimizing cost function J, then calculate output h using chosen parameter. Output is the probability $P(y=1|x, \theta)$

Logistic regression cost function

$$\rightarrow J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

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

Note: $y=0$ or 1 always

$$\rightarrow \text{Cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$$

If $y=1$: $\text{Cost}(h_\theta(x), y) = -\log(h_\theta(x)) = 1$

If $y=0$: $\text{Cost}(h_\theta(x), y) = -\log(1-h_\theta(x))$

Logistic regression cost function

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

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

To fit parameters θ :

$$\min_{\theta} J(\theta) \quad \text{Cost} \Theta$$

To make a prediction given new x :

$$\text{Output } h_\theta(x) = \frac{1}{1+e^{-\theta^T x}} \quad p(y=1|x, \theta)$$

2. Reason for choosing this particular function

- This function can be derived from statistics (统计) using the principle of **maximum likelihood estimation** (最大似然), which is the idea in statistics for how to efficiently find parameters for different models.

- This function is CONVEX

Note: another view of cost function, derived from likelihood estimation:

(后验概率: 已知y, 预测最可能的x分布) max probability: $P(y=1|x, \theta) = p(y=1|x, \theta)^y * (1-p(y=1|x, \theta))^{1-y}$
maximize probability $\log p = y \log(h) + (1-y) \log(1-h) \rightarrow$
minimize cost: $J = -\log(p)$

3. Gradient descent

for each parameter θ_j : $\theta_j := \theta_j - \alpha * \text{derivative of } J \text{ over } \theta_j$,
simultaneously update all θ_j .

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

Note:

- parameter gradient descent algorithm looks identical to linear regression.
- Only difference is: hypothesis function h . for linear regression $h=\theta^T X$, logistic regression $h=1/(1+e^{-\theta^T X})$;
- could use vectorization to update all parameters in one time:
Gradient with parameter vector $\theta=[\theta_0, \dots, \theta_n]$
 $\theta := \theta - \alpha * \sum ((h(X) - Y) * X, x=1)$
 $\theta := \theta - \alpha / m * X \cdot (H-Y)$

Gradient Descent

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

Want $\min_{\theta} J(\theta)$: $\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$

Repeat {

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

} (simultaneously update all θ_j)

Algorithm looks identical to linear regression!

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

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

4. Monitor Gradient descent

- same method as linear regression could apply on logistic regression to make sure converge correctly.
- feature scaling could also apply on logistic regression to make gradient descent converge faster.

7_49. Advanced optimization_Logistic Regression

Purpose: to run much more quickly with gradient descent.

1. Optimization algorithm

Gradient descent:

given parameter, have code that compute J , and derivative of J , than plug in parameter gradient descent, which can then try to minimize the cost function.

Note:

for gradient descent, only compute derivative of J is ok;
for monitoring cost function convergence, need compute cost function J also.

(note: for batch gradient descent: compute cost J periodically as need to scan all example, computation cost is expensive;
while for stochastic could compute cost on every iteration, and average over like 1000 examples, then plot the averaged cost;
mini-batch: could compute cost on every iteration and plot—also average over certain iterations??
)

Optimization algorithm

Cost function $J(\theta)$. Want $\min_{\theta} J(\theta)$.

Given θ , we have code that can compute

$\rightarrow -J(\theta)$ (for $j = 0, 1, \dots, n$)

Gradient descent:

Repeat {

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

}

2. other sophisticated algorithm than gradient descent to minimize the cost function:

- Conjugate gradient;
- BFGS;
- L-BFGS

Advantages:

- No need to manually pick learning rate;
Have a clever inner-loop-line search algorithm, automatically try out different values for the learning rate and automatically pick a good learning rate. Even could pick different learning rate for each iteration.

- Converge often faster than gradient descent.

Disadvantages:

- More complex

Optimization algorithm

Given θ , we have code that can compute

$\rightarrow -J(\theta)$ (for $j = 0, 1, \dots, n$)

Optimization algorithms:

- Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS

Advantages:

- No need to manually pick α
- Often faster than gradient descent.

Disadvantages:

- More complex

3. Example_other sophisticated algorithm

- Octave:

fminunc: function minimization unconstrained in Octave (无约束最小化函数) .

Note:

- parameter dimension have $>= 2$ to use fminunc function;
- parameter vector start from θ_1 , not θ_0 .

- normally use when have large algorithm to run.

[optTheta, functionVal, existFlag]...= fminunc (@costFunction, initialTheta, options);
feedback:
optTheta: parameters chosen;
functionVal: cost function value;
existFlag: converged or not; if is 1, converged.

Example: $\min \theta$
 $\rightarrow \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$ $\theta_1 = 5, \theta_2 = 5$.
 $\rightarrow J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$
 $\rightarrow \frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5)$
 $\rightarrow \frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$
 \rightarrow options = optimset('GradObj', 'on', 'MaxIter', '100');
initialTheta = zeros(2,1);
[optTheta, functionVal, exitFlag] ...
= fminunc(@costFunction, initialTheta, options);

```
function [jVal, gradient] = costFunction(theta)
    jVal = (theta(1)-5)^2 + ...
            (theta(2)-5)^2;
    gradient = zeros(2,1);
    gradient(1) = 2*(theta(1)-5);
    gradient(2) = 2*(theta(2)-5);

```

theta = $\begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$ $\theta_0 = \text{theta}(1)$
 $\theta_1 = \text{theta}(2)$
 \vdots
 $\theta_n = \text{theta}(n+1)$

```
function [jVal, gradient] = costFunction(theta)
    jVal = [code to compute  $J(\theta)$ ];
    gradient(1) = [code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$ ];
    gradient(2) = [code to compute  $\frac{\partial}{\partial \theta_1} J(\theta)$ ];
    :
    gradient(n+1) = [code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$ ];

```

7_50.Multi-class classification: One-vs-all_Logistic Regression

e.g: output have more than two classification

1. Training data example:

for multi-class classification, training data example like right pic.

Multiclass classification

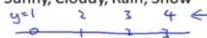
Email folder/tagging: Work, Friends, Family, Hobby



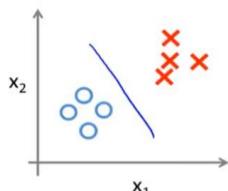
Medical diagrams: Not ill, Cold, Flu



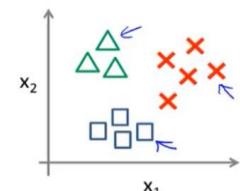
Weather: Sunny, Cloudy, Rain, Snow



Binary classification:



Multi-class classification:



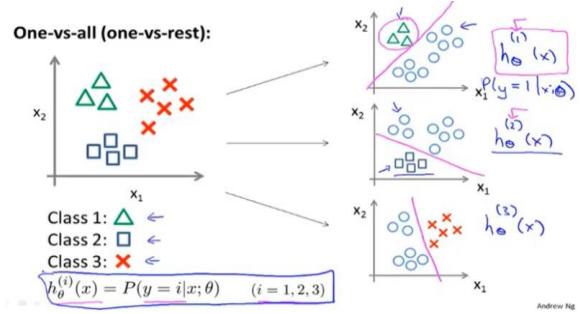
2. multi-classification Idea: one-versus-all/rest classification

one-versus-all/rest: for multi classification:

take these training set into separate binary classification problems

- create a new training set: only one positive/class and negative output-all data except positive.
-->get hypothesis function h_1 and decision boundary: $p(y=1|x, \theta_1)$
 θ_1 : parameter vector for class 1.
- create new training set: another positive/class and negative output-->
-->get hypothesis for this positive data: h_2 and 2nd decision boundary: $p(y=2|x, \theta_2)$
- Repeated, till find all the hypothesis h and decision boundary for all class. $p(y=3|x, \theta_3)$

for given example X , run $h(x)$ with all parameter for different class separately, and choose the highest possibility as output class.
(note: why not just use softmax..)



Summary:

- Train a logistic regression classifier $h_{\cdot i}(x)$ for each class i to predict the probability that $y=i$
- On a new input x , to make a prediction, pick the class i that maximizes $\max_i h_{\cdot i}(x)$: run all classifiers $h_{\cdot i}(x)$ on the input x , then pick the class i that maximizes among all the classes.

One-vs-all

Train a logistic regression classifier $h_{\theta}^{(i)}(x)$ for each class i to predict the probability that $y = i$.

On a new input x , to make a prediction, pick the class i that maximizes

$$\max_i h_{\theta}^{(i)}(x)$$

W8-Regularization

8_52.The problem of overfitting_Regularization

1. Example_Linear regression

- underfit / high bias: hypothesis model does not fit training data well

call bias: if fitting a straight line to the data, then as if the algorithm has a very strong preception or very strong bias, that housing prices are going to vary linearly with their size. And despite the data to the contrary, despite the evidence of the contrary, its preconceptions still or bias still causes it to fit a straight line.

- overfit / high variance: fitting such a high order polynomial (高阶多项式) almost as if it can fit almost any function, this face of possible hypothesis is just too large, it's too variable. do not have enough data to constrain it, to give us a good hypothesis.

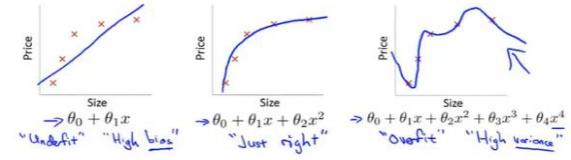
(note: training set have too many features-->hypothesis high order polynomial, no enough constraint-training set.)

overfitting: if have too many features (high order polynomial), the learned hypothesis may fit the training set very well (cost function close to 0), but fail to generalize to new examples.

try too hard to fit training data, then fail to fit new samples.

Generalize: means how well a hypothesis applies even to new examples-have not seen in training set.

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

2. Example_Logistic regression

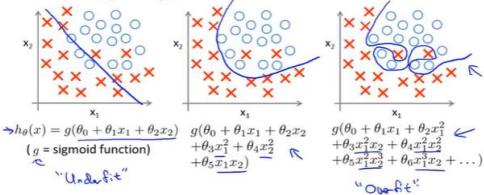
underfit / high bias: decision boundary does not fit training data well.

overfit / high variance: high order polynomial for decision boundary: generate high order polynomial terms of features.

algorithm try really hard to find a decision boundary to fit training data, or go to great lengths to contort itself to fit every training data well.(make every examples' loss =0)

check decision boundary fitting with training data.

Example: Logistic regression



Summary:

check fitting status for both linear regression /logistic regression:

check fitting with training data:

Linear regression: fitting of hypothesis function with training data;

Logistic regression: fitting of decision boundary with training data.

3. Addressing

3.1: Reduce training set feature numbers:

- Manually select which features to keep.
- Keep important features.
- Model selection algorithm : model automatically decide which feature to keep. ??

Reducing feature number works well for overfitting.

But disadvantage: is by throwing away some features is also throwing away some of the information you have about the problem.

(note: same reason for using reduce_dimension on training set: do not use dimension reduction unless have strong evidence current solution not working/too slow)

Addressing overfitting:

Options:

1. Reduce number of features.
 - Manually select which features to keep.
 - Model selection algorithm (later in course).
2. Regularization.
 - Keep all the features, but reduce magnitude/values of parameters θ_j
 - Works well when we have a lot of features, each of which contributes a bit to predicting y .

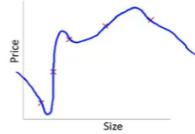
3.2 Regularization

- keep all features, but reduce magnitude / values of feature parameters works well when have a lot of features, each of which contributes a bit to predicting y .

(note: root for overfitting: too many data features-->algorithm high polynomial, and no enough training set to constrain these parameters)

Addressing overfitting:

- x_1 = size of house
- x_2 = no. of bedrooms
- x_3 = no. of floors
- x_4 = age of house
- x_5 = average income in neighborhood
- x_6 = kitchen size
- :
- x_{100}



3.3: Plotting hypothesis(decision boundary)_does not work for addressing overfitting

when have only 1 or two features, we can just plot the hypothesis(decision boundary), and select the appropriate degree polynomial.

Plotting hypothesis could be one way try to decide what degree polynomial to use. but does not often work. more often we may have problems where we just have a lot of features, -->not just a matter selecting polynomial, and when have many features also becomes much harder to plot the data.-->**harder to visualize it to decide what features to keep or not.**

when have lots of feature and small training data, overfitting can be a problem.

(note: estimate different hypothesis on training & dev set: cost function J_{train} , J_{dev} of each type hypothesis vs iteration-->choose feature number:

1. just fit training set , no bias, no overfitting: J_{train} small

2. and generalize to dev set well: J_{dev} small)

8_53. Cost function_Regularization

1. Intuition:

Reduce feature weight by adding in cost function

-->when weight is small enough, is kind of remove that corresponding feature.

(note: penalizing weight while considering label y value)

2. Regulation Idea:

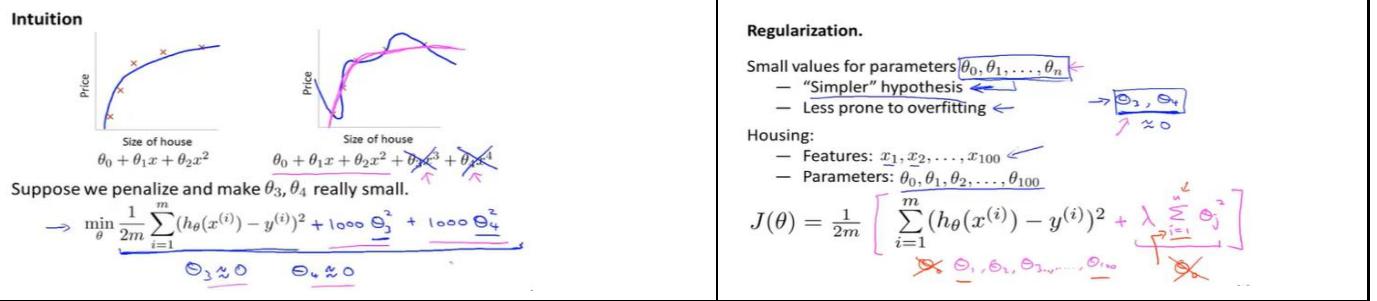
small values for parameters $\theta_0, \theta_1, \theta_2, \dots, \theta_n$ -->

- 'Simpler' hypothesis:

smaller parameters corresponds to usually smoother functions also simpler. also less prone to overfitting.

e.g.: $\theta_0, \theta_1, \dots, \theta_n$ close to 0, then h nearly only have 2 features.

-Less prone to overfitting

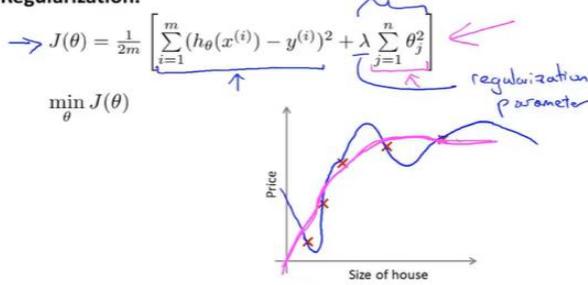


3. Regularizaiton:

as no idea which feature weight to shrink, modify cost function: adding regulation term in the end to shrink every single of parameter

Note:
 θ_0 shrink $\theta_1, \theta_2, \theta_n$.
 θ_0 not included here, but very little influence on result.
(note: make $d/d\theta_k = 0 \rightarrow$
 $\theta_k = a \times k / (a + x_k \cdot k^2)$; a : real value
1. $\lambda = 0$:
feature x_k big: \rightarrow weight θ_k small
feature x_k very small: \rightarrow weight big
2. $\lambda \neq 0 \rightarrow$ reduce weight optima

Regularization.



Regularization.

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

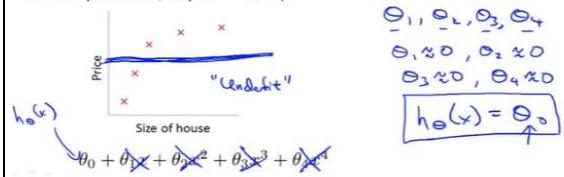
$\min_{\theta} J(\theta)$

regularization parameter

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 $\lambda = 10^{10}$)?



4. Regulation parameter λ :

Function: is a controls a trade off between two different goals.

first goal: without regulation term, would like to fit the training data well.

second goal: regulationterm: try to keep paramters small.

λ by balancing this two goal, to make hypothesis relatively simple to avoid overfitting.

e.g: when parameter is high, \rightarrow high order polynomial hypothesis, fitting traing data well, (first goal), yet regulation term is high \rightarrow reduce parameter, fitting error may increase, yet regulation term reduce.

(\rightarrow regulation will make first goal bigger: fit training set less well $\rightarrow J_{\text{train}}$ increase, yet less overfitting)

need have good chosen for regulation parameter λ :

when λ is too big \rightarrow all parameter (except θ_0) close to 0, hypothesis is a flat straight line, could not fit training set \rightarrow underfitting.

8_54. Regularized linear regression_Regularization

Regularized linear regression

1. Object

find parameter θ to minimizing cost function J .

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

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

2. Gradient descent with regulation

2.1 Gradient descent formula

write gradient descent formula θ_0, θ_j separately, as θ_0 is not involved in regulation term.

Gradient descent

$$\text{Repeat } \{$$

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

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

2.2: Interpreter of gradient descent

- First part: factor add for $\theta_j : 1 - \alpha \lambda / m < 1$,
A little bit less than 1, shrinking the parameter a little bit.
(α is small, while m is very large \rightarrow a little bit smaller than 1)

- Second part: exactly same as the formula before add regulation term.

$$\theta_j := \theta_j * (1 - \alpha \lambda / m) - \alpha / m * [\text{sum} (h_i - y_i) * x_j, \text{over } i \text{ to all example } m]$$

Gradient descent

```

Repeat {
     $\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$ 
     $\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad (j = 1, 2, 3, \dots, n)$ 
}
 $\theta_j := \theta_j (1 - \alpha \lambda / m) - \alpha / m * [\text{sum} (h_i - y_i) * x_j, \text{over } i \text{ to all example } m]$ 
 $|1 - \alpha \lambda / m| < 1$ 

```

Andrew Ng

3. Normal equation regulation

Method:

Create new matrix X and labeled output vector y. --> find min J by making derivative of J over each parameter = 0

$$X = [x_0 \ T, \dots, x_m \ T] \ T. \text{mx} (n+1) \quad \rightarrow \min_{\theta} J(\theta) \quad \frac{\partial}{\partial \theta_j} J(\theta) = 0 \quad m \times n \\ x_i = [x_{i,0}, \dots, x_{i,n}] \quad \Rightarrow \Theta = (X^T X + \lambda \begin{bmatrix} 0 & 1 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix})^{-1} X^T y \\ \text{E.g. } n=2 \quad \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (n+1) \times (n+1)$$

3.1: add new matrix (n+1, n+1) in the parameter:

Normal equation

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \text{IR}^m \\ \rightarrow \min_{\theta} J(\theta) \quad \frac{\partial}{\partial \theta_j} J(\theta) = 0 \quad m \times (n+1) \\ \rightarrow \Theta = (X^T X + \lambda \begin{bmatrix} 0 & 1 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix})^{-1} X^T y \\ \text{E.g. } n=2 \quad \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (n+1) \times (n+1)$$

3.2 non-invertibility

- without regulation:

when have fewer example than features $m \leq n$, matrix $X^T X$ is non-invertible / singular / degenerate. (pinv could calculate fake invertible anyway)

- with regulation:

as long as regulation parameter λ is > 0, $X^T X + \lambda * \text{new matrix}$ is invertible.

Non-invertibility (optional/advanced).

Suppose $m \leq n$, $\begin{pmatrix} \# \text{examples} & \# \text{features} \end{pmatrix}$

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

If $\lambda > 0$,

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

Summary:

adding regulation could fix overfitting when have lots of features with small training set.

8_57. Regularized logistic regression_Regularization

1. Regularized logistic regression

method: add regulation term in cost function. Penalizing parameters when they are too big --> make decision boundary more reasonable.

E.g. when decision boundary is high order polynomial, --> overfitting.

By penalizing parameters in cost function with regularization term, --> reduce parameter, --> make decision boundary more reasonable for separating positive and negative examples.

Regularized logistic regression.

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2^2 + \dots)$$

Cost function:

$$\rightarrow J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \quad | \theta_1, \theta_2, \dots, \theta_n$$

Gradient descent

Repeat {

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \rightarrow \theta_j := \theta_j - \alpha \underbrace{\left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]}_{\theta_1, \dots, \theta_n} \quad | \theta_1, \dots, \theta_n \\ \frac{\partial J(\theta)}{\partial \theta_j} \quad h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

2. Gradient descent with regulation

2.1 Gradient descent formula

write gradient descent formula θ_0, θ_j separately, as θ_0 is not involved in regulation term.

Formula looks identical with linear regression, except only hypothesis function is different.

3. Advanced optimization algorithm

- need to define costFunction (θ) and gradient; then pass to fminunc

$$\theta = [\theta_0, \dots, \theta_n] \rightarrow [\theta_0, \dots, \theta_n]$$

Note:

- parameter dimension have $>= 2$ to use fminunc function;
- parameter vector start from θ_1 , not θ_0 .

In coding:

- update cost function jVal with regulation term.
- update gradient of each parameter:
- gradient (1) is derivative of J over θ_0 .
- gradient (2): update with regulation term over θ_1 .

Advanced optimization

\rightarrow function [jVal, gradient] = costFunction(theta)
 \rightarrow jVal = [code to compute $J(\theta)$];
 \rightarrow $J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log (h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$
 \rightarrow gradient(1) = [code to compute $\frac{\partial}{\partial \theta_0} J(\theta)$];
 \rightarrow $\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$
 \rightarrow gradient(2) = [code to compute $\frac{\partial}{\partial \theta_1} J(\theta)$];
 \rightarrow $\left(\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)} + \frac{\lambda}{m} \theta_1 \right)$
 \rightarrow gradient(3) = [code to compute $\frac{\partial}{\partial \theta_2} J(\theta)$];
 \vdots
 \rightarrow $\left(\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)} + \frac{\lambda}{m} \theta_2 \right)$
 \rightarrow gradient(n+1) = [code to compute $\frac{\partial}{\partial \theta_n} J(\theta)$];

W9-Neural Networks: Representation

9_58. Non-linear hypotheses-Neural Networks: Representation

1. Machine learning _feature problem:

when feature number is large, feature number to consider

->>1. all quadratic features:

like 100, then the number of quadratic features grows (二次项的个数: $x_1 \times x_2, x_1 \times x_3 \dots$) : rafly $n^2/2$

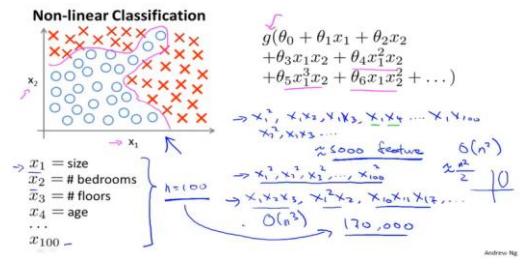
including all the quadratic features maybe not a good idea, may end up **overfitting** and computation cost is high for working with so many features.

- >>2. quadratic features: but only $x_1^2, x_2^2, \dots, x_n^2$, only n quadratic features.

-->not enough features, **could not fit interesting hypothesis** (decision boundary simpler), can not fit the training set well

- >>3. quadratic features $x_1^2, x_2^2, \dots, x_n^2$, and cubic / third order polynomial features: $x_1 \times x_2 \times x_3, x_1^2 \times x_2 \dots \rightarrow O(n^3)$ -->cubic features number about 170,000

-->including these high order polynomial features when original features is large, will dramatically **blows up feature space** and is not a good way to come up with additional features with which to build non-linear classifiers when n is large.



yet for many machine learning problems, n is pretty large.
E.g. in computer vision:

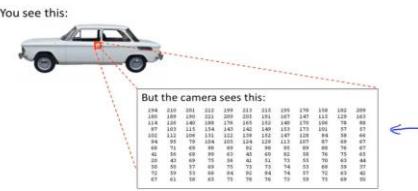
2. Example_Machine learning problem

e.g: car classification

2.1. Car door handle classification

Algorithm deal with below datas (pixel data , representing a certain position in pic.) to recognize this represents a door handle

What is this?



2.2 computer vision: car detection:

Normal process:

use machine learning to build a car detector: come up with a label training set-->learning algorithm trained a classifier-->test new image, see performance

Computer Vision: Car detection



2.3: understanding_why need non-linear hypothesis

> In the training set, pick a couple of pixel location in the image, -->plot the car at the location , at a certain point (depending on intensity of pixel one and pixel two),

>Do same thing with other training set image: look at the same two pixel locations. that image have different intensity of pixel one and different intensity for pixel two,-->end up different location in the figure.

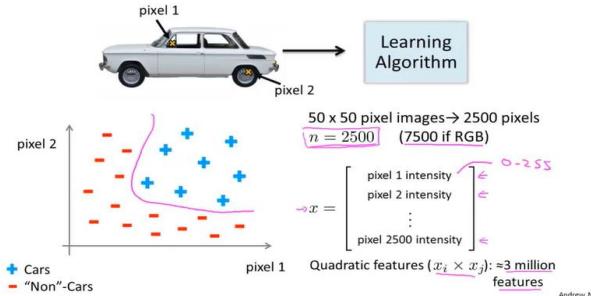
>plot some negative examples (not car)

>will find car and non car lying in different regions of the space.-->need non-linear hypothesis to separate out the two classes:

dimension space:

- suppose image is 50 x 50 pixel-->2500 pixels , dimension of feature size will be $n=2500$ (7500 if RGB).

- if including all the quadratic features in learning algorithm, --> $n^2/2 \sim 3$ million features
too large to be reasonable. computation cost will be very expensive to find and represents all of these three million features per training example while feature vector is the list of all the pixel intensity.



Summary:

simple logistic regression, together with quadratic or the cubic features, is not a good way to learn complex nonlinear hypotheses, **when n is large**. Because just end up with too many features.

Neural network is a good way to learn complex nonlinear hypothesis , even when input features is large.

9_59. Neurons and the brain-Neural Networks: Representation

Neural network works well for different modern ML problem and AI.

1. Neural networks status:

Origins: algorithm that try to mimic the brain

Recent resurgence: state-of-the-art technique for many applications.

Reason: neural network computionally somewhat more expensive algorithm, it only because computes became fast enough, to really run large scale neural networks.

Neural Networks

- Origins: Algorithms that try to mimic the brain.
- Was very widely used in 80s and early 90s; popularity diminished in late 90s.
- Recent resurgence: State-of-the-art technique for many applications

2. 'One learning algorithm' hypothesis

Hypothesis: brain do different things, is not worth like a thousand different programs, instead the brain does it is worth just a single learning algorithm

2.1 Proven for this hypothesis:-neuro-rewiring experiments

e.g 1: auditory cortex (听觉区) :

normal process:

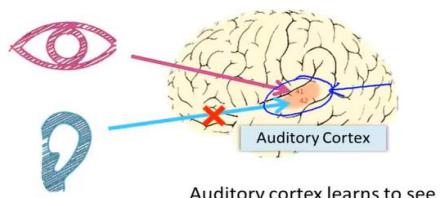
ear pick up sound signal --> route sound signal to auditory cortex --> understand words/vocie;

new process: cut wire from ear to auditory cortex, and rewire to eye
signal to the optic nerve, eventually gets routed to the auditory cortex -->**auditory cortex will learn to see (in every single sense of the word see as we know: perform visual discrimination task: look at image and make decision by auditory cortex).**

Summary:

if the same piece of physical brain tissue can process sight, sound or touch, then maybe there is one learning algorithm that can process sight or sound or touch. instead of implement a thousand different programs/algorithms to do the thousand wonderful things that brain does.
maybe what we need to do is to find some approximation or to whatever the brain's learning algorithm is and implement that and let the brain **learned by itself: how to process these different types of data.**

The "one learning algorithm" hypothesis



Auditory cortex learns to see

3. Large extent of brain

seems we can plug in any sensor to any part of the brain, within the reason the brain will learn to deal with data.

More e.g:

- Seeing with tongue: help blind people to see.
Camara capture pic, pix-->transfer tonue by wire: each pixel tensity reflected by voltage-->learned to see
- Human echolocaton (sonar): help blind people to location
learn to interpret the pattern of sounds bouncing off environment.
- Haptic belt: Direction sense
similar with how birds can sense where north is.
- Implanting a 3rd eye: in frog.
frog will learn to use that eyse as well.

Sensor representations in the brain



9_60. Model representation I- Neural networks: Representation

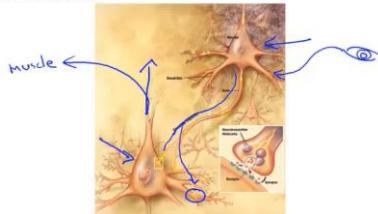
1. Neuron in the brain:

Simulating neurons or networks of neurons in the brain.

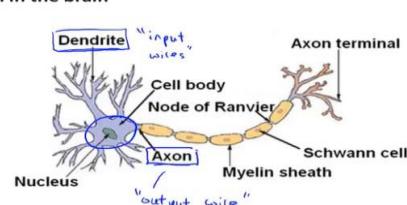
1.1 Neuron unit:

input: Dendrite ---> Nucleus ---> Output Axon (computation), send message to other neurons (by sending little pulse of electricity), and will be input of next neurons.

Neurons in the brain



Neuron in the brain



2. Neuron model: logistic unit (singal neuron)

>Input: X (included x_0 or not)

>Nucleus: computation

>output-->activation function

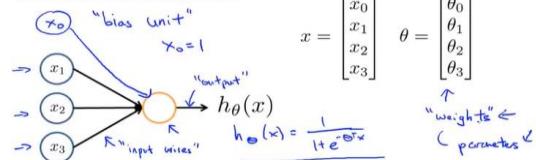
>parameter θ , also called weight.

Note:
 x_0 is added into X vector sometimes, as bias unite. $X_0=1$;
 $h_{\theta}(x)=1/(1+e^{-\theta^T x})$

sigmoid (logistic) activation function:

an artificial neuron with a sigmoid or a logistic activation function. The activation function is another term for that non-linearity function.

Neuron model: Logistic unit



Sigmoid (logistic) activation function.

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

3. Neuron network

A group of different neurons strung together.
Final Output: $h(x)$: hypothesis of input X.

Input layer: X;

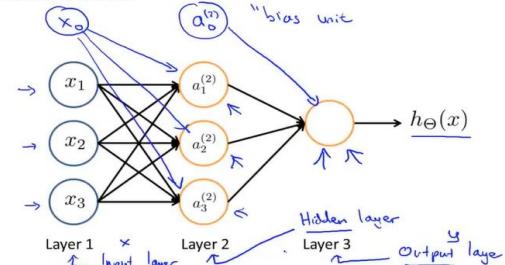
Hidden layer: do not get to observe in the training set;
in supervised learning, get to see input and correct output. whereas the **hidden layer** are values do not get to observe in the training set;
layers that are not input layer , output layer, all called hidden layer.
output layer: $h(x)$

Note:

>Input x_0 could add or not.

> a_0 could also add as an extra bias unite there: $a_0=1$. (always is 1)

Neural Network



4. Neuron network_computational step

Activation: $a(j)_i$

$a(j)_i$ = 'activation' of unit i in layer j
by activation, means the value that is computed by and that is output by a specific.

Weight matrix:

$\Theta(j)$ = matrix of weights controlling function mapping from layer j to layer j+1.

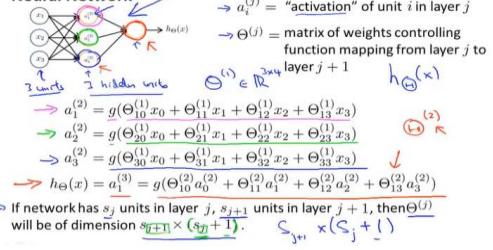
weight matrix dimension: $\Theta(j) : s_{(j+1)} \times (s_j + 1)$ (bias)

feature number in next layer x feature number in current layer

s_j : units/feature number in layer j

$s_{(j+1)}$: feature number in layer j+1.

Neural Network



9_61. Model representation II- Neural networks: Representation

How carry out computation efficiently -> vectorization implementation

why neural network representation is a good idea and how they can help to learn complex non-linear hypotheses

1.2: Vectorizing forward propagation:

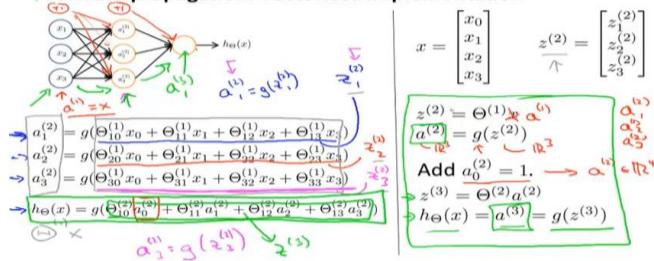
$$x(i) = [x(i)_0, \dots, x(i)_n]$$

$$x(i)_0 = 1$$

- define extra term Z: for layer i
- $z(i) = \Theta_{(i-1)}^T * a(i-1)$
- $a(i) = g(z(i))$: a & z have same dimension, g(z) apply sigmoid function to every element of Z.
- Take care of bias unit: add $a(i)_0 = 1$

Vectorizing is a relatively efficient way of computing h of x.

Forward propagation: Vectorized implementation



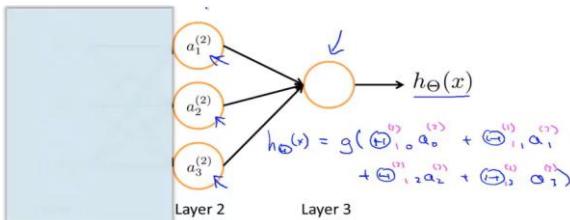
2. Neural network learning its own features

how could neural network help to learn interesting nonlinear hypothesis

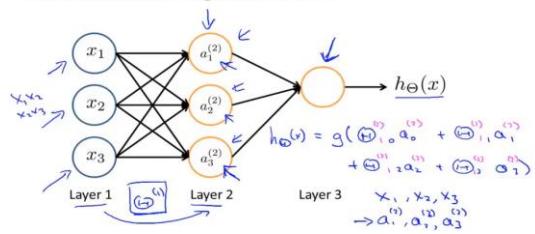
'Standard logistic regression':

for each neuron unit: sort of logistic regression: using logistic regression unite to make a prediction h of x: $h = g(\Theta T a)$

but only the features fed into logistic regression, are these values computed by the hidden layer.



Neural Network learning its own features



Summary:

Neural network is just like logistic regression, except that rather than using the original features x_1, x_2, \dots, x_n , is using these new features a_1, a_2, \dots, a_l . while a_1, a_2, \dots, a_l they themselves are learned as functions of input. Concretely the function mapping from layer 1 to layer 2 is determined by some other set of parameters $\Theta_{(1)} (\Theta_{(1),10}, \dots, \Theta_{(1),11}, \Theta_{(1),12}, \dots, \Theta_{(1),20}, \dots)$.

-->neural network, instead of being constrained to feed the original feature to logistic regression, it gets to learn its own features $a(i)$ ($a(i)_1, a(i)_2, \dots$) to feed into the logistic regression. and depending on what parameters it chooses for $\Theta_{(1)}$ can learn interesting and complex features. and therefor can end up with a better hypothesis (better than considering only original features x_1, x_2, \dots or even polynomial features $x_1 \cdot x_2, \dots$).

Neural network have flexibility to learn whatever feature want. and then feed into last unit: logistic regression for output.
(note: previous L-1 layer, learn new features based on original features, as input to the final logistic regression unit.)

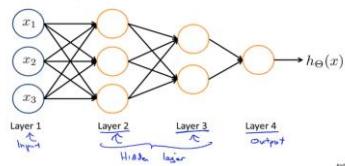
3. Other network architectures

architecture: refers to how the different neurons are connected to each other.

Right pic.: input feature feed into layer 1 to compute complex feature and feed to layer 2 , compute new complex feature and feed to layer 3... till output.

hidden layer: anything not an input layer or output layer

Other network architectures



Summary:

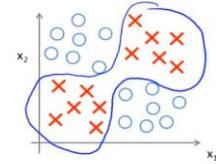
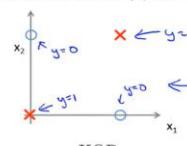
feed forward propagation in a neural network: start from the activations of the input layer, and forward propagate that to the first hidden layer, then the second layer, and then finally the output layer. While in this process: vectorize computation.

9_63. Example and intuitions I- Neural networks: Representation

How neural network compute complex features, logistic function

Non-linear classification example: XOR/XNOR

$\Rightarrow x_1, x_2$ are binary (0 or 1).



1. Non-linear classification example: XOR / XNOR

Function: output result of XOR inputs

x_1, x_2 are binary (0 or 1) \rightarrow function $y=x_1 \text{ XOR } x_2, X_1 \text{ XNOR } x_2$ (Not $x_1 \text{ XOR } x_2$)

2. Simple example : AND function

x_1, x_2 are binary (0,1), $x_0=1$; $y=x_1 \text{ AND } x_2$

\rightarrow neural network logistic regression:
 $h(x)=g(\theta_0 * x_0 + \theta_1 * x_1 + \theta_2 * x_2)$

Note:

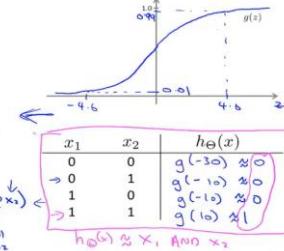
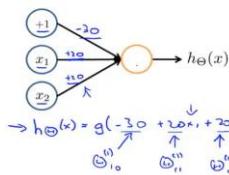
for sigmoid function, $g(z) = 1/(1+e^{-z})$, when $Z=4.6 \rightarrow g(z)=0.99 \approx 1$; when $Z=-4.6, g(z)=0.01 \approx 0$

by choosing parameter $\theta \rightarrow h(x) \approx x_1 \text{ AND } x_2$; right pic.: $\theta=[-30; 20; 20]$

Right pic. pink table shows how neural network computes:

Simple example: AND

$\Rightarrow x_1, x_2 \in \{0, 1\}$
 $\Rightarrow y = x_1 \text{ AND } x_2$

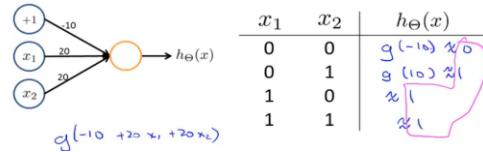


Example: OR function

3. Example: OR function

Similar as AND function, choose different parameter $\theta \rightarrow h(x) \approx x_1 \text{ OR } x_2$

$\theta=[-10; 20; 20]$



9_65. Example and intuitions II- Neural networks: Representation

How neural network compute complex features

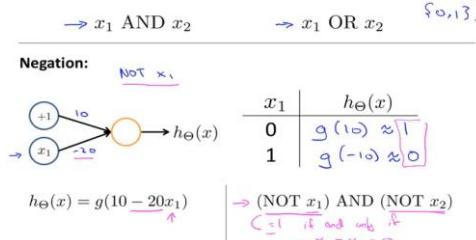
1. Negation_not x1

x_1 is binary(0,1)

Idea: put a large negative weight for variable want to negate
Parameter: $\theta=[10; -20]$

2. Function (NOT x1) AND (NOT x2)

\rightarrow only if $x_1=x_2=0$,
could choose parameter Parameter: $\theta=[10; -20; -20]$



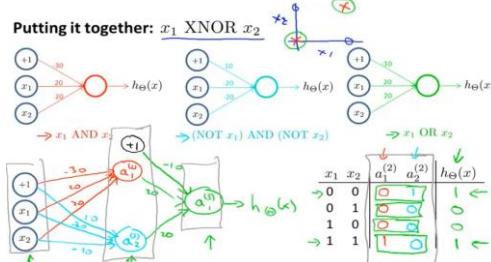
3.Putting it together: x1 XNOR x2

use two hidden layer:

- first hidden layer_2nd layer a(2):
 $a(2)_1 = x_1 \text{ AND } x_2$ function; binary (0,1)
 $a(2)_2 = (\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$ function; binary (0,1)
 $a(2)_0 = 1$

- 2nd hidden layer_2nd layr a(3)
 $a(3)_1 = x_1 \text{ OR } x_2$; binary (0,1)

(note: kind of breakdown problem-computing new feature , into computing many simple tasks/simple new features, whose computation based on previous new feature computed.)



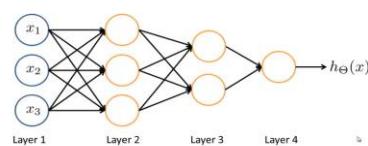
4. Neural network intuition:

hidden layer have slightly more complex functions. \rightarrow and yet add another layer, end up with an even more complex nonlinear function. \leftarrow why neural network can compute pretty complicated functions.

When have multi layers, have relative simple function of the inputs the second layers, but the 3rd layer can build on that and compute even more complex functions, and the layer after that can compute even more complex functions.

(note: deeper layer , compute more complex functions/features)

Neural Network intuition



Handwritten digit classification

5. Example_Handwritten digit classification

neural network application that captures above intuition of deeper layers computing more complex features.

Input: pic: handwritten character shown to network;

1st hidden layer: different features: edges, and lines and so on detected.

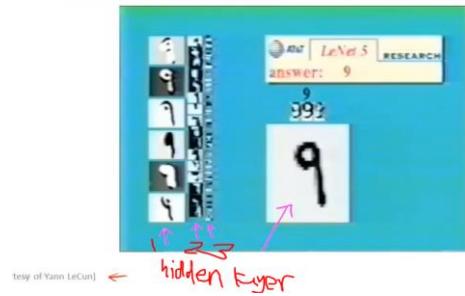
right pic.-a visualization of the features computed by sort of the first layer of the network.

2nd hidden layer: harder to understand the visualization of this layer.

3rd hidden layer: even harder to understand;

-->probably have a hard time seeing what's going on much beyond the first hidden layer, but finally all the learned features get fed to the output layer, and show the final answer.

classifier work well: different type, noise , flip,could look into video in this course



9_66. Multi-class classification- Neural networks: Representation

handwritten digit classification is actually multi-class classification problem, have 10 possible categories.

1. Multi output units: One-vs-all method

try to recognize categories of input: how many categories in input image.

Category number of input: c

1.1 Output:

- neural network output vector dimension is $(c,1)$: corresponding to classification result of each category for one input image.

- output value expected $h(x) = [0, \dots, 1, 0, \dots, 0]$: 1 for categories in image (may have multi category 1 in image), 0 for categories not in image.

e.g: right pic, have four categories-->**have four logistic regression classifiers $a(4)_1, a(4)_2, a(4)_3, a(4)_4$, each of which is trying to recognize one of the four classes that we want to distinguish amongst.**

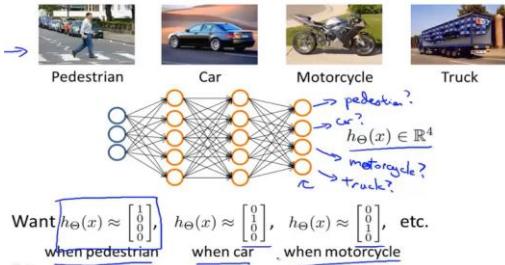
$a(4)_i$: result value is 0-1, corresponding to the possibility of being i category with given image.

(note:

1. output: may have multi object: -->multi classfire-logistic regression units in output layer

2. output: have only one object of multi possible class-->multil unit/classfire: softmax)

Multiple output units: One-vs-all.



2. Training set representation

2.1 Training set:

pair: $(x(1), y(1)), \dots, (x(m), y(m))$;

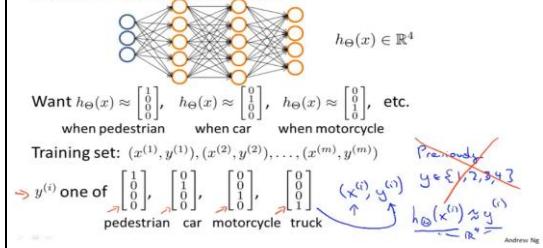
label $y(i) = [0, \dots, 1, 0, \dots, 0]$ ($c \times 1$ dimension): depending on what the corresponding image $x(i)$ is.

2.2 Object:

find $h(x(i)) (c,1) \sim y(i) (c, 1)$

$h(x)$ and label y have same dimension

Multiple output units: One-vs-all.



W10-Neural networks: Learning

10_67. Cost function_Neural networks: Learning

Learning algorithm for fitting the parameters of the neural network given the training set.

(note: neural network: hard to find decision boundary; as the learned new feature: $A[i-1]=f(x_{i-1}, \dots, x_n)$ to complicated to compute)

1. Neural Network (classification problem)

Training set:

pair: $(x(1), y(1)), \dots, (x(m), y(m))$:

$x(i)$ is.

K = class number

L = total layer number in network (including input and output layer)

s_l = units number in layer l , not counting bias unit ($a(l)$). $0=1$.

Output:

- Binary classification: $y = 0 / 1$ ($K < 3$)

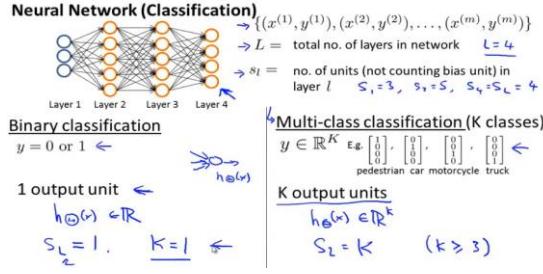
- only one output unit ($s_L / K=1$). $h(x)$ is an gender

- Multi-class classification (K classes, $K > 3$)

$y \in \mathbb{R}^K$

label $y(i)=[0, \dots, 1, 0, \dots, 0]$ ($c \times 1$ dimension): depending on what the corresponding image

$-K$ output units $s_L=K$, $h(x) \sim \mathbb{R}^K$



2. Cost function:

2.1 the generalization of the one used for logistic regression:

- θ_0 not included in regularization

- first part: sum of all example standard cost without regularization;

- second part: regularization : **sum of all parameters** $\lambda/2 * \sum \theta_j^2$

Logistic regression:

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

Cost function

Logistic regression:

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

Neural network:

$$\Rightarrow J(\Theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\Theta(x^{(i)})_k) + (1 - y_k^{(i)}) \log(1 - h_\Theta(x^{(i)})_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_j^{(l)})^2$$

2.2 Neural network cost function:

one classifier cost function is same as standard logistic regression cost function-->

for all classifier cost function:

>>> 1. sum of each classifiers' cost function: over all K classifier . for one example

>>> 2 average of sum one example cost over all training example,

>> second part: regulation: consider all parameters in neural net work: sum of each layer parameters' regulation. $\Theta_{-l}(s_{-l+1}, s_{-l})$: parameter matrix in layers l

Note: do not consider Θ_0 either: $\Theta_{-l}(j, i)$: start from 1, to $s_{-l}(i+1)$, i from 1 to s_{-l}

$\Theta_{-l}(j, 0)$: multiply $a_{-l}(0) (=1)$ -->kind of like bias unit, same as logistic regression regularization, will not sum over into regularization term (because we do not want to regularize them and string the values 0).

but if consider $\Theta_{-l}(j, 0)$ in regularization, it will work out same, does not make a big difference.

$h(x)$: K dimensional output vector;

$h(x)_j$: selects out the j -th element of the vector $h(x)$

$y(i)$: labeled output y for i -th example/input, K dimensional

$y(i)_j$: select out the j -th element of vector $y(i)$

10_68. Backpropagation algorithm_Neural networks: Learning

either use gradient computation or advanced algorithm, need code to compute cost function J and partial derivative of J over each parameter

1. Gradient computation

Partial derivative J over each parameter

given one training example (x, y)

1.1 Forward computation : given the input, what the hypothesis actually output

$a(1) = X$;

$Z(2) = \Theta(1)^T \cdot x$

$a(2) = g(Z(2))$ (add $a(2)_0$)

$Z(3) = \Theta(2)^T \cdot a(2)$

$a(3) = g(Z(3))$ (add $a(3)_0$)

$Z(4) = \Theta(3)^T \cdot a(3)$

$a(4) = g(Z(4)) = h(X)$, Θ

forward computing by vectorizing implementation, allow to compute activation values for all of the neurons.

Gradient computation

$$\Rightarrow J(\Theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\Theta(x^{(i)})_k) + (1 - y_k^{(i)}) \log(1 - h_\Theta(x^{(i)})_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_j^{(l)})^2$$

$$\Rightarrow \min_{\Theta} J(\Theta)$$

Need code to compute:

$$\Rightarrow \frac{\partial J(\Theta)}{\partial \Theta_{ij}^{(l)}} J(\Theta) \leftarrow \Theta_{ij}^{(l)} \in \mathbb{R}$$

2. Backpropagation algorithm_Gradient computation

$\delta(l)_j$: activation of the j unit/neural duo in layer l

$\delta(l)_j$: 'error' of node j in layer l : capture error in the activation of that neural do.

For each output unit

$\delta(L)_j = a(L)_j / h(x)_j - y_j$

y_j = j -th element in the labeled output vector

For all output unite: K dimension

$\delta(L) = a(L) / h(x) - y$

Note:

>1. No $\delta(1)$: as first layer corresponding to input layer, and that's features observed in training set, so that does not have any error associated with that. - we do not really want to try to change those values.

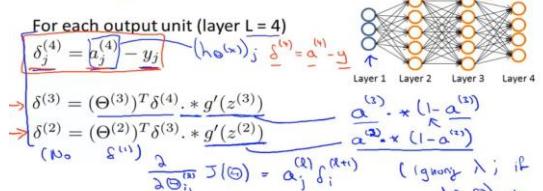
>2. Name backpropagation comes from the fact:

Start to computing delta term for output layer, then go back a layer and compute the delta terms for this layer, and then go back another step to compute delta: backpropagating error of the output layer to previous layer.

>3. Using these delta term, could compute derivative terms for all of parameters.

Gradient computation: Backpropagation algorithm

Intuition: $\delta_j^{(l)}$ = "error" of node j in layer l .



公式推导

为了施展 反向传播 的魔法，我们首先要引入一个中间变量 δ ，定义为：

$$\delta_j^l = \frac{\partial J}{\partial z_j^l}$$

其中 l 为第几层， j 表示第 l 层的第几个神经元， z 为上次课程提到的中间变量（为了让式子看上去更清晰，反向传播 中的公式上标不使用括号）。 δ 被称为第 l 层第 j 个神经元的误差。

反向传播 就是先计算每一层每个神经元的误差，然后通过误差来得到梯度的。

首先来看输出层的误差：

$$\delta_j^L = \frac{\partial J}{\partial z_j^L}$$

对它使用 梯式法则 得到：

$$\delta_j^L = \sum_{k=1}^K \frac{\partial J}{\partial a_k^L} \frac{\partial a_k^L}{\partial z_j^L}$$

而只有当 $k == j$ 时，右边部分才不为 0，所以：

$$\delta_j^L = \frac{\partial J}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^L} = \frac{\partial J}{\partial a_j^L} g'(z_j^L) = a_j^L - y_j^L \quad (1)$$

对于其它层的误差：

$$\delta_j^l = \frac{\partial J}{\partial z_j^l}$$

使用 梯式法则：

$$\delta_j^{l+1} = \sum_{k=1}^{n+1} \frac{\partial J}{\partial a_k^{l+1}} \frac{\partial a_k^{l+1}}{\partial z_j^l} = \sum_{k=1}^{n+1} \delta_k^l \frac{\partial a_k^{l+1}}{\partial z_j^l}$$

其中：

$$a_k^{l+1} = \sum_{p=1}^{s_l} \Theta_{kp}^l g(z_p^l) + b_k^l$$

求偏导得：

$$\frac{\partial a_k^{l+1}}{\partial z_j^l} = \Theta_{kj}^l g'(z_j^l)$$

所以：

$$\delta_j^l = \sum_{k=1}^{n+1} \Theta_{kj}^l \delta_k^{l+1} g'(z_j^l) = \sum_{k=1}^{n+1} \Theta_{kj}^l \delta_k^{l+1} a_j^l (1 - a_j^l) \quad (2)$$

同样地使用 梯式法则 有：

$$\frac{\partial J}{\partial \Theta_{ij}^l} = \sum_{k=1}^{n+1} \frac{\partial J}{\partial a_k^{l+1}} \frac{\partial a_k^{l+1}}{\partial \Theta_{ij}^l} = \sum_{k=1}^{n+1} \delta_k^l \frac{\partial a_k^{l+1}}{\partial \Theta_{ij}^l}$$

3. Backpropagation algorithm_put all together

All training examples involved.

$\Delta(l)ji=0$ (initializing) (use to compute derivative of J over parameter)
 $(\Delta(l))ij = dJ/d\Theta(l)_{ij}$

process: look through training set.

>Loop : repeat through all training examples:

for one example ($x(i), y(i)$):
 >>>Set $a(1)=x(i)$
 >>>Perform forward propagation to compute activationi value $a(l)$
 >>>Using $y(i)$, compute $\delta^{(L)}=a(L)-y(i)$
 >>>Compute $\delta^{(0)}$: $l=1, 2, \dots, L-1$
 >>> Compute error in wheel neural network for one training example:
 $\Delta(l)=\Delta(l)+\delta(l+1) * a(l)^T$ (accumulate these partial derivative terms): update of all values of ij position.

>out of loop

derivatione J over $\Theta(l)$:

$D(l).ij := 1/m \Delta(l).ij + \lambda/m * \Theta(l).ij$ if $j \neq 0$;
 $D(l).ij := 1/m \Delta(l).ij$ if $j = 0$ (corresponding to bias term)

由于：

$$z_k^{l+1} = \sum_{p=1}^{s_l} \Theta_{kp}^l g(z_p^l) + b_k^l$$

只有当 $k = i \cdot p = j$ 时留下一项：

$$\frac{\partial J}{\partial \Theta_{ij}^l} = g(z_j^l) \delta_i^{l+1} = a_j^l \delta_i^{l+1} \quad (3)$$

反向传播

有了 (1) (2) (3) 式，就可以来完成 反向传播 算法了（需要注意的是刚才所推导的式子都是针对一 组训练数据而言的）。

1. 对于所有的 l, i, j 初始化 $\Delta_{ij}^l = 0$

2. 对于 m 组训练数据， k 从 1 取到 m :

· 令 $a^l = x^{(k)}$

· 前向传播，计算各层激活向量 a^l

· 使用 (1) 式，计算输出层误差 δ^L

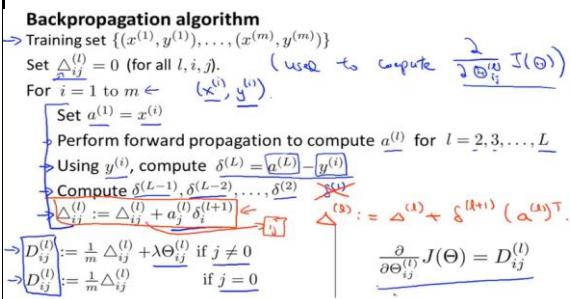
· 使用 (2) 式，计算其它层误差 $\delta^{L-1}, \delta^{L-2}, \dots, \delta^2$

· 使用 (3) 式，累加 Δ_{ij}^l 。 $\Delta_{ij}^l := \Delta_{ij}^l + a_j^l \delta_i^{l+1}$

3. 计算梯度矩阵：

$$D_{ij}^l = \begin{cases} \frac{1}{m} \Delta_{ij}^l + \frac{\lambda}{m} \Theta_{ij}^l & \text{if } j \neq 0 \\ \frac{1}{m} \Delta_{ij}^l & \text{if } j = 0 \end{cases}$$

4. 更新权值 $\Theta^l := \Theta^l + \alpha D^l$



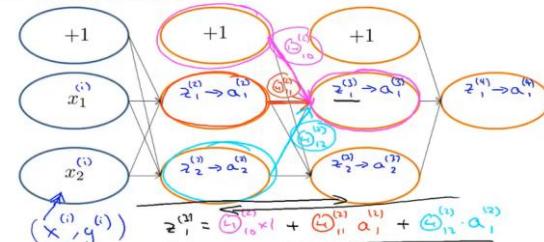
10.69. Backpropagation Intuition_Neural networks: Learning

1. Forward:

input: $(x(i), y(i))$
 $x(i) \rightarrow [z(1) \rightarrow a(1)] \rightarrow [z(2) \rightarrow a(2)] \rightarrow [z(3) \rightarrow a(3)] \rightarrow [z(4) \rightarrow a(4)]$

e.g:
 $z(3)_1 = \Theta(2)_10 + \Theta(2)_11 * a(2)_1 + \Theta(2)_12 * a(2)_2$
 $z(3)_2 = \Theta(2)_20 + \Theta(2)_21 * a(2)_1 + \Theta(2)_22 * a(2)_2$

Forward Propagation



3. Forward propagation:

$a(l)_{ij}$: activation of the j unit/neural duo in layer l
 $\delta(l)_{ij}$: 'error' of cost for $a(l)_{ij}$ (unit j in layer l)
 'error' of node j in layer l : capture error in the activation of that neural duo.
 $\delta(l)_{ij}$: partial derivative of cost function over $z(l)_{ij}$:
 If $z(l)_{ij}$ changes-->change output function -->end up change cost function.
 Measure: how much we want to change neural network weights, in order to affect the intermediate values (z) of the computation, so as to affect the final output that the neural network h(x), and therefore affect the overall cost.

2. What is backpropagation doing?

single output cost function as right:

$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\Theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\Theta(x^{(i)})) + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2$$

focusing on signel example, one output unit, no regularization:
 could take it as square error, approximate way of measuring distance of prediction and labeled value.

$$\text{cost}(i) = y^{(i)} \log h_\Theta(x^{(i)}) + (1 - y^{(i)}) \log h_\Theta(x^{(i)})$$

What is backpropagation doing?

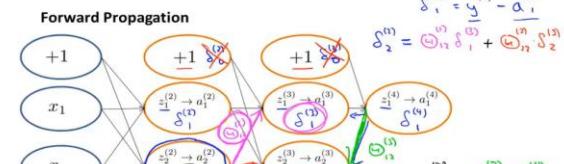
$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\Theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\Theta(x^{(i)})) + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2 \quad (x^{(i)}, y^{(i)})$$

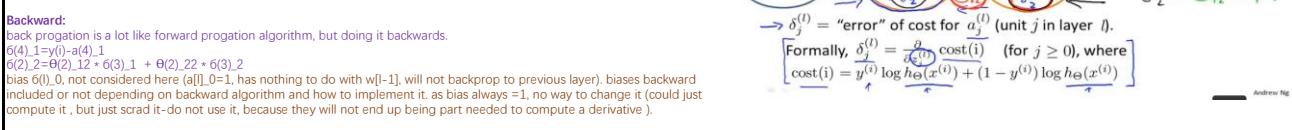
Focusing on a single example $x^{(i)}, y^{(i)}$, the case of 1 output unit, and ignoring regularization ($\lambda = 0$),

$$\text{cost}(i) = y^{(i)} \log h_\Theta(x^{(i)}) + (1 - y^{(i)}) \log h_\Theta(x^{(i)}) \quad \text{(Think of cost}(i) \approx (h_\Theta(x^{(i)}) - y^{(i)})^2)$$

i.e. how well is the network doing on example i ?

Forward Propagation





forward:
 $x(3)_1 = \Theta(2)_1 \cdot 10 + \Theta(2)_2 \cdot 12 * a(2)_1 + \Theta(2)_3 \cdot 12 * a(2)_2$
 $x(3)_2 = \Theta(2)_1 \cdot 20 + \Theta(2)_2 \cdot 21 * a(2)_1 + \Theta(2)_3 \cdot 22 * a(2)_2$

10_70. Implementation note: Unrolling parameters_Neural networks: Learning

Unrolling matrix into vector, needed in order to use the advanced optimization routines.

1. Advanced optimization

>costFunction:
 inputs parameter: 'theta', 'vector'
 output: cost function value 'JVal', and derivatives 'gradient'-vector

>then pass this value to an advanced algorithm like fminunc
 Note: fminunc is not the only one advanced algorithm, there are also other advanced authorization algorithms.
 >>> take costFunction , and initial parameterTheta (vector).

when using logistic regression, this works fine as theta and gradient is vector. but for a full neural network, we will have parameter matrix $\Theta(1)$, $\Theta(2)$, $\Theta(3)$, and gradient matrix $D(1)$, $D(2)$, $D(3)$.

Advanced optimization

```
function [JVal, gradient] = costFunction(theta)
    ...
optTheta = fminunc(@costFunction, initialTheta, options)
```

Note: fminunc is not the only one advanced algorithm, there are also other advanced authorization algorithms.
 >>> take costFunction , and initial parameterTheta (vector).

Neural Network ($L=4$):
 $\rightarrow \Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)} - \text{matrices } (\Theta_{1,1}, \Theta_{1,2}, \Theta_{1,3})$
 $\rightarrow D^{(1)}, D^{(2)}, D^{(3)} - \text{matrices } (D_1, D_2, D_3)$
 "Unroll" into vectors

2. Unroll into vectors

thetaVec = [Theta1 (:); Theta2 (:); Theta3 (:)]
 Dvec = [D1 (:); D2 (:); D3 (:)]

D / Theta (:); take all the element in matrix;
 and unroll them and put all the elements into a big long vector thetaVec, Dvec.

3. Vector -->Matrix

Theta1 = reshape(thetaVec(1:110), 10, 11)

take vector thetaVec first 110 elements , and use reshape command to get 10x11 matrix.

Example

$s_1 = 10, s_2 = 10, s_3 = 1$
 $\rightarrow \Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$
 $\rightarrow D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$
 $\rightarrow \text{thetaVec} = [\Theta_{1,1}; \Theta_{1,2}; \Theta_{1,3}; \dots; \Theta_{3,1}; \Theta_{3,2}; \Theta_{3,3}]$
 $\rightarrow \text{DVec} = [D_{1,1}; D_{1,2}; D_{1,3}; \dots; D_{3,1}; D_{3,2}; D_{3,3}]$

Theta1 = reshape(thetaVec(1:110), 10, 11);
 Theta2 = reshape(thetaVec(111:220), 10, 11);
 Theta3 = reshape(thetaVec(221:231), 1, 11);

4. Learning algorithm

>have initial parameter
 >unroll to get 'initialTheta' to pass to fminunc

Implement costFunction:
 >get inputs thetaVec : theta1 unrolled vector
 >get back original matrix from thetaVec by using reshape function: $\Theta_1, \Theta_2, \Theta_3$
 >Use forward prop/back prop to compute D_1, D_2, D_3 and cost function J:
 >Unroll D_1, D_2, D_3 to get gradientVec, which can return by costFunction-a vector of these derivatives.

Learning Algorithm

> Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$.
 > Unroll to get initialTheta to pass to
 > fminunc (@costFunction, initialTheta, options)

```
function [JVal, gradientVec] = costFunction(thetaVec)
    ...
    From thetaVec, get  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$  reshape
    Use forward prop/back prop to compute  $D^{(1)}, D^{(2)}, D^{(3)}$  and  $J(\Theta)$ 
    Unroll  $D^{(1)}, D^{(2)}, D^{(3)}$  to get gradientVec.
```

Summary:
matrix representation advantage: when parameters are stored as matrices, it is more convenient when doing forward propagation and back propagation, and easier to take advantage of vectorized implementations.

Vector representation advantage: when using advanced optimization algorithm, algorithms tend to assume all parameters unrolled into a big long vector.

10_71. Gradient checking_Neural networks: Learning

Unfortunate property:

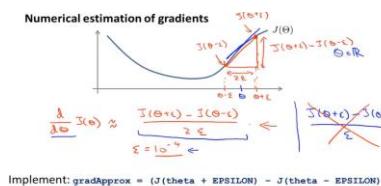
there are many ways to have subtle bugs in back prop, if run with gradient or some other algorithm, it could actually look like it's working. Cost function J of theta may end up decreasing on every iteration of gradient descent. (how this comes???)

Might just wind up with neural network that has a high level of error, than you would with a bug-free implementation. and you might just do not know it's subtle bugs that giving this performance.

Gradient check eliminate all of these problem.

(note:
 1. linear regression or logistic regression:
 >>1. could check plotting of cost vs iterations-->debug J_train: decrease with every iteration if learning rate is reasonable small
 Out put = $W \cdot X \cdot \text{sigmoid}(W \cdot X)$: ->if gradient descent is wrong-->J_train plotting will be unnatural

2. Neural network:
 >> debug: by gradient check:
 >> Output = $g(w[i] \cdot f(w[i-1] \cdot f(\dots)))$: -->cost = f(output); d/d parameter: is too complicated, if computation wrong, plotting J_train/dev may still decrease



1. Numerical estimation of gradients

Derivative of cost function J over one specific parameter $\Theta(l)_j$ - right pic blue line, could check by slope of it's two side point $\Theta(l)_j + \epsilon$ and $\Theta(l)_j - \epsilon$

ϵ could take 10^{-4} , if too small may cause numerical problem.

Note: Two-sided differences gives a slightly more accurate estimate, rather than one-sided difference estimate.

2. Parameter vector (unrolled matrix into vector)

do two-sided difference for each parameter in parameter vector.

Implement in Octave: use 'for loop' for each parameter in vector.

>check two-sided difference calculated above with Dvec (from back prop)
 If very close , have more confidence that computing derivatives correctly in back prop /other sophisticated algorithm.-->hope code run correctly and do a good job optimizing J of theta.

Parameter vector θ

$\rightarrow \theta \in \mathbb{R}^n$ (E.g. θ is "unrolled" version of $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$)
 $\rightarrow \theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n]$
 $\rightarrow \frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{J(\theta_1 + \epsilon, \theta_2, \theta_3, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \theta_3, \dots, \theta_n)}{2\epsilon}$
 $\rightarrow \frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \theta_2 + \epsilon, \theta_3, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \theta_3, \dots, \theta_n)}{2\epsilon}$
 \vdots
 $\rightarrow \frac{\partial}{\partial \theta_n} J(\theta) \approx \frac{J(\theta_1, \theta_2, \theta_3, \dots, \theta_n + \epsilon) - J(\theta_1, \theta_2, \theta_3, \dots, \theta_n - \epsilon)}{2\epsilon}$

`for i = 1:n,
 thetaPlus = theta;
 thetaPlus(i) = thetaPlus(i) + EPSILON;
 thetaMinus = theta;
 thetaMinus(i) = thetaMinus(i) - EPSILON;
 gradApprox(i) = (J(thetaPlus) - J(thetaMinus)) / (2 * EPSILON);
end;`

Check that $\text{gradApprox} \approx \text{DVec}$

From back prop.

3. Implementation Note:

>Implement backprop to compute DVec (unrolled D1, D2, D3...)
 >Implement numerical gradient check to compute gradApprox:
 >Make sure they give similar values.
 >Turn off gradient checking. Using backprop code for learning

Important:

>Be sure to disable gradient checking code before training classifier.
 If run numerical gradient computation on every iteration of gradient descent, or in the inner loop of costFunction, code will be very slow as:

gradient checking is very computationally expensive, is a very slow way to try to approximate the derivative.. While backprop is much more computationally efficient way of computing the derivatives. So once checked backprop is doing right, should turn off gradient checking.

Implementation Note:

- > Implement backprop to compute DVec (unrolled $D^{(1)}, D^{(2)}, D^{(3)}$).
- > Implement numerical gradient check to compute gradApprox.
- > Make sure they give similar values.
- > Turn off gradient checking. Using backprop code for learning.

Important:

- > Be sure to disable your gradient checking code before training your classifier. If you run numerical gradient computation on every iteration of gradient descent (or in the inner loop of `costFunction(...)`) your code will be very slow.

10_72. Random initialization_Neural networks: Learning

when run gradient descent or other advanced algorithm, need pick some initial value for parameter theta.
 Then can slowly take steps to go downhill, using gradient descent, to minimize the function J of theta.

1. Zero initialization:

Initialzing all paramters to 0 is ok for logistic regression. (note: error do not transfered by weight)

But does not work when traing neural network. (note: for layer 2- L-1: cost error transferred by weight=0)
 E.g: right pic 3layer, two feature in fist lyer, second layer, 1 feature in output
 $\theta_{in_1}=0, \theta_{in_2}=0$
 $X=[1; x1; x2]$

1.1 Forward

> $Z(2)=\theta_{in_1}^T X = 0+X = 0$ vector
 $A(2)=a(2)_1; a(2)_2=g(Z(2))=0$ (all unite values in layer 2 is same)
 -> $a(2)_1=a(2)_2=\dots=a(2)_u$: u: unite number
 $Z(3)=\theta_{in_2}^T a(2)+A(2)=0$ vector
 $A(3)=g(Z(3))=g(0) vector=0$ (all unite values in layer 3 is same)
 -> $a(3)_1=\dots=a(3)_u$: u: unite number

-if $\theta_{in_l}=0$:
 unites in each layer , have same value $A(l)=g(\theta_{in_l}) + A(l-1)=g(0) vector$
 (generated by $g(0)$ vector);
 all $Z(l)$ in each layer is 0 vector (due to parameter $\theta_{in_l}=0$), activation $A(l)$ have same value in each element

>if $\theta_{in_l}=0$, rest parameter in each layer $\theta_{in_l}=a_{in_l}$ (same value for parameters in same layer)
 $A(l)=a_{in_l}$: all activation values in same layer are same.

1.2 Backprop

1st iteration:
 $\delta(L)=h(x)-y$ ($!=0$, != same value)
 $\delta(L-1)=\delta(L)\cdot\theta_{in_L}(L-1)\cdot g'(Z(L-1))$
 ...
 $\delta(2)=\delta(3)\cdot\theta_{in_2}(2)\cdot g'(Z(2))$
 $\theta_{update_l}=\theta_{in_l} - \alpha \cdot \delta(l)$

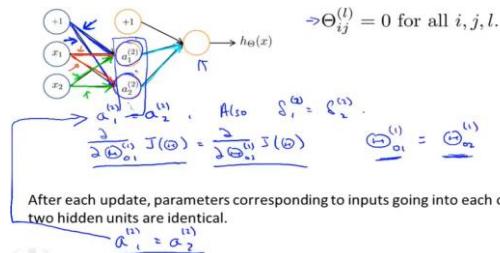
>if $\theta_{in_l}=0$
 -> $\delta(2)=\delta(3)=\dots=\delta(L-1)=0$
 -> $\theta_{update_l}=0$ (parameters could not be updated if set all initial paramters to 0)

>if only 1 output (classifer, $\theta(L)$ is same in all uintes), $\theta_{in_1}=0$, rest parameter in each layer $\theta_{in_l}=a_{in_l}$ (same value for parameters in same layer)
 -> $\delta(l)=0$
 $\theta_{update_l}=\theta_{in_l} - \alpha \cdot \delta(l)$ = same value for all unites in θ_{update_l}

>even after one gradient descent update, θ_{update_l} all values are still same
 >>neural network really can not compute very interesting functions: all the hidden units in one layers are computing the exact same feature-this is highly redundant representation, as final logistic regression only gets to see one feature.

Note: for logistic regression:
 only two layer: input $X \rightarrow Z=W.X, A=Y$
 if all parameter =0, as $\delta(l)=h(x)-y$ ($!=0$, != same value)-> $\theta_{update_l}=\theta_{in_l} - \alpha \cdot \delta(l) = -\alpha \cdot 0 + A(l) = -\alpha \cdot 0 + A(l)$
 parameter update ok, as $\delta(2)$ not =0
 (note: but if mult class output: multi logistic regression classifier: -->output layer multi unit; then all output layer unit: have same weight value,-->only one classifier)

Zero initialization



Initial value of Θ

For gradient descent and advanced optimization method, need initial value for Θ .

`optTheta = fminunc(@costFunction, initialTheta, options)`

Consider gradient descent

Set initialTheta = zeros(n,1) ?

from course:
 $\delta(2)=\delta(2)_2$
 >derivative of δ over parameter $\theta_{in_1}=0$ = over parameter $\theta_{in_2}=0$.
 derivative of δ over $\theta_{in_1}=0$ = $\delta(2)_1 \cdot T \cdot \theta_{in_1} \cdot g'(Z(1))$
 >updated after 1 iteration: $\theta_{in_1}=0=\theta_{in_1}-0$

after each update, paramters corresponding to inputs going into each two hidden unites are identical

2. Random initialization : symmetry breaking

Random initialization could solve above problem.

Initial each θ_{ij} to a random value in $[-\epsilon, \epsilon]$:
 >rand (10,11):
 generate a random 10 x 11 dimensional matrix: all of the values are between 0 and 1.
 >compute rand(10,11)*2 -<-<->matrix values between $[-\epsilon, \epsilon]$.

Note: ϵ here is different thing with ϵ used in grad checking

(note:
 1. feature: 0 mean ans scaling to (-1,1) range, by $(x-u)/s$
 s: could be max-min or variance of feature x_k
 2. weight: initialize: $[-\epsilon, \epsilon]$; close to 0)

Summary:
 To train neural network, randomly initialize the weights to small values close to 0 between $[-\epsilon, \epsilon]$, then do backprop and use either gradient descent or one of the advanced optimization algorithm to try to minimize cost function J of theta, as a function of the parameter theta.

Starting from just randomly chosen initial value for the parameters, and by doing symmetry braking, hopefully gradient descent or advanced algorithm will be able to find a good value of theta.

(note: symmetry braking: needed for optimizing method that use gradient descent)

10_73. Putting it together_Neural networks: Learning

1. Training a neural network

pick a neural network architecture.

Architecture: connectivity pattern between the neuron. How many unites in hidden layer, and how many hidden layers.

>Features decided by training set:
 >>input unites number: feature $x(i)$ dimension;
 >>Output unites number: number of class.
 Note: labeled y for choosing classes from 1,2,3,..., should rewrite $y=[0,0,...1,0,...]$

Random initialization: Symmetry breaking

> Initialize each $\Theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$
 (i.e. $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$)
 E.g.
 >> Theta1 = $\text{rand}(10,11) * (2 * \text{INIT_EPSILON}) - \text{INIT_EPSILON};$
 >> Theta2 = $\text{rand}(1,11) * (2 * \text{INIT_EPSILON}) - \text{INIT_EPSILON};$

Training a neural network

Pick a network architecture (connectivity pattern between neurons)

→ No. of input units: Dimension of features $x^{(i)}$

→ No. output units: Number of classes

Reasonable default: 1 hidden layer, or if >1 hidden layer, have same no. of hidden

> Reasonable default: 1 hidden layer (perfectly reasonable default)
 or if >1 hidden layer, have same number of hidden units in every layer (usually the more the better);
 but if have a lot of hidden units, will be very computation expensive. Very often, having more hidden units is a good thing.
 number of hidden layer units: maybe comparable to the dimension of input feature, or could be same number with input features or to maybe twice or three or four times of input feature number.

hidden units in each layer: input feature normally very large, --> still set hidden units to be comparable or times of input feature???

2. Training a neural network

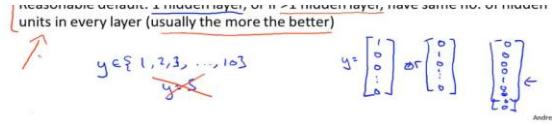
> 1. Randomly initialize weights; to small values near zero.
 > 2. Implement forward propagation to get output vector $h(x)$ with given x .
 > 3. Implement code to compute cost function J of theta.
 > 4. Implement backprop to compute partial derivatives of J over $\theta_{(l),ij}$.
 when do backprop, usually use for loop over the training example, there is advanced vectorization methods where don't have a for-loop over the m-training examples, but first time to implement backprop, should almost certainly be a for loop in the code when iterating over the examples;
 $x(1), y(1) \rightarrow$ forward prop --> backprop, then
 $x(2), y(2) \rightarrow$ forward prop --> backprop, then
 till $x(m), y(m)$

first time try backprop better use for loop instead of advanced algorithm.

```
for i=1:m
  %from forward prop and back prop using example x(i), y(i):
  %get activations a(l) and delta terms delta(l): l=2,...,L
  Delta(l)=Delta(l)+6*(l+1)*a(l)'*Delta(l)
}
compute these partial derivative terms: d J(theta)/theta_(l),ij;
also need to take regulation terms lambda as well.
```

> 5. Use grad checking to compare $d J(\theta)/\theta_{(l),ij}$ computed using backprop vs using numerical estimate of gradient of $J(\theta)$. Reassure backprop computation is correct, then disable gradient checking code.

> 6. Use gradient descent or advanced optimization method with backprop (for calculating partial derivatives $d J(\theta)/\theta_{(l),ij}$) to try to minimize $J(\theta)$ as a function over parameters θ .



Training a neural network

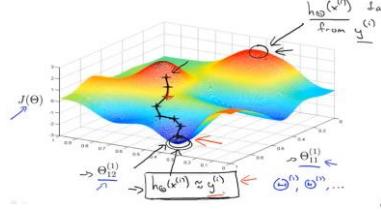
- 1. Randomly initialize weights
 - 2. Implement forward propagation to get $h_\Theta(x^{(i)})$ for any $x^{(i)}$
 - 3. Implement code to compute cost function $J(\Theta)$
 - 4. Implement backprop to compute partial derivatives $\frac{\partial}{\partial \Theta_{j,k}^{(l)}} J(\Theta)$
 - for $i = 1:m$ { $(x^{(i)}, y^{(i)})$, $(x^{(i)}, y^{(i)})$, ..., $(x^{(i)}, y^{(i)})$ }
 - Perform forward propagation and backpropagation using example $(x^{(i)}, y^{(i)})$
 - Get activations $a^{(l)}$ and delta terms $\delta^{(l)}$ for $l = 2, \dots, L$.
 - $\Delta^{(2)} := \Delta^{(2)} + \delta^{(2)}(a^{(1)})^T$
 - ...
 $\Delta^{(L)} := \Delta^{(L)} + \delta^{(L)}(a^{(L-1)})^T$
 - compute $\frac{\partial}{\partial \Theta_{j,k}^{(l)}} J(\Theta)$.
 - Then disable gradient checking code.
 - 6. Use gradient descent or advanced optimization method with backpropagation to try to minimize $J(\Theta)$ as a function of parameters Θ
- $J(\Theta)$ — non-convex.

Note:
 for neural network, cost function $J(\Theta)$ is non-convex. It can theoretically be susceptible to local minima. In fact gradient descent or other advanced algorithm can get stuck in local optima, but it turns out usually it is not a huge problem, as usually algorithm like gradient descent will do a great job minimizing cost function J of theta, and get a very good local minimum.

3. Intuition of backprop in neural network:

when implement backprop and use gradient descent or one of the advanced optimization methods, right pic shows what the algorithm is doing:
 start from some initial point, and repeatedly go downhill. Backprop is for computing the direction of the gradient, and gradient descent is for taking little steps downhill, until hopefully it gets to a good local optima.

try to find values of the parameters where the output values in the neural network closely matches the values of the labeled $y(i)$ (minimizing cost function J): low values of the cost function corresponds where J of theta is low and therefore corresponds to where the neural network happens to be fitting training set well. $h_\Theta(x,i) \approx y(i)$ is needed to be true in order for J of theta to be small.



10_74. Autonomous driving example_Neural networks: Learning

Autonomous driving: getting a car to learn to drive itself.

1. Visualization:

> down in the lower left: view seen by car, kind of see a road.
 > on top in the left:
 >> First horizontal bar: direction selected by human driver and is the location of this bright white band: left end-steering hard to the left; right end-side-steering hard to the right.
 White band is a little bit to the left: means human driver at this point is steering slightly to the left.
 >> Second horizontal bar: direction selected by learning algorithm: white band means direction selected is a just slightly to the left.
 before algorithm starts learning initially, the output is like a uniform grey band--corresponding to neural network having been randomly initialized, and initially having no idea how to drive the car, /which direction to select, and only after it's learned for a while and then to output like a solid white band in just a small part of the region corresponding to a particular steering direction.



2. Autonomous navigation experiment:

ALVINN: a system of artificial neural networks, that learns to steer by watching a person drive.

NAVLABII: car used for this experiment

neural network setup: Learn to steer by watching a person drive.

Device: Sensors, computer, and actuators

> Initial step in configuring this car: training a network to steer
 >> during training, human driver drive the car while car watches.
 >> Once every two seconds, ALVINN digitizes a video image of the road ahead, and records the person's steering direction.
 > training image is reduced in resolution to 30 x 32 pixels and provided as input to ALVINN's three-layer network.
 >> using backprop learning algorithm, ALVINN is trained to use same steering direction as the human driver for that image.
 >> Initially steering directly is random, after about 2 minutes training, the network learns to accurately imitate the steering reactions of the human driver.

>> Same training procedure used in other road type

> After training: neural network to drive

>> 12 times per second, ALVINN digitizes an image, and feeds it to its neural networks.

>> each neural network, work in parallel, produces a steering direction and a measure of its confidence in its response, steering direction from the most confident network.
 e.g. now the one-lane algorithm is used to control the vehicle.
 when approach intersection, confidence of the one-lane network decreases.-->
 after went through intersection, and the two-lane road ahead comes into view, the confidence of the two-lane network rises, then two-lane algorithm is selected to steer vehicle on the two-lane road.

W11-Advice for applying Machine learning

10_77. Deciding what to try next_Advice for applying Machine learning

Purpose:
While improve algorithm, what are the proxy avenues to try next.

1. Debugging a learning algorithm

Problem: hypothesis on new set , have large error in its predictions.-->improve algorithm

Method:

>Get more training examples.

Not always help

>Try smaller set of features

(very effective for overfitting, yet will throw away useful info., use regularization instead.)

>Try getting additional features:

e.g for predicting house price: get more info. about the pieces of land...and so on.

>>before doing above things, helpful to know in advance if this is going to help.

>Try adding polynomial features (x^1, x^2, x^3, x^4 , etc)

(help underfitting, make algorithm more powerful with more weights)

>Try decreasing regularization parameter λ

>Try increasing regularization parameter λ

each above method can scale to a long time project: spend may about 6 month on one term

2. Machine learning diagnostic

Diagnostic: A test that can run to gain insight what is /isn't working with a learning algorithm, and gain guidance as to how best to improve its performance.

Diagnostics can take time to implement, but doing so can be a very good use of your time.
when developing learning algorithms, as they will save you spending many months pursuing an avenue that could have found out much earlier what not going to be fruitful.

Debugging a learning algorithm:

Suppose you have implemented regularized linear regression to predict housing prices.

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

However, when you test your hypothesis on a new set of houses, you find that it makes unacceptably large errors in its predictions. What should you try next?

→ - Get more training examples

- Try smaller sets of features $x_1, x_2, x_3, \dots, x_{100}$

→ - Try getting additional features

- Try adding polynomial features (x_1^2, x_2^2, x_1x_2 , etc.)

- Try decreasing λ

- Try increasing λ

Techni easily to rule out many of these options.

Machine learning diagnostic:

Diagnostic: A test that you can run to gain insight what is /isn't working with a learning algorithm, and gain guidance as to how best to improve its performance.

Diagnostics can take time to implement, but doing so can be a very good use of your time.

10_78. Evaluating a hypothesis_Advice for applying Machine learning

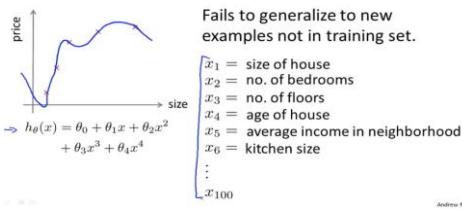
1. Evaluating hypothesis

Principle: how algorithm trained have low error in training set , generalize to new examples not in training set.

1.1 Method:

>Plotting hypothesis: to see overfit or not
very hard to plot when have many features.

Evaluating your hypothesis



1.2 Standard process:

>Split datat set into : training set (70%)-m_training, test set (30%)-m_test

Note: if there is any ordinary to the data, better split data randomly. Or randomly reorder the examples before splitting into training set and test set.

Evaluating your hypothesis

Dataset:

Size	Price		
2104	400		
1600	330		
2400	369		
1416	232		
3000	540		
1985	300		
1534	315		
1427	199		
1380	212		
1494	243		

Training set m_{train}

$(x^{(1)}, y^{(1)})$
 $(x^{(2)}, y^{(2)})$
 \vdots
 $(x^{(m)}, y^{(m)})$

$m_{\text{test}} = \text{no. of test example}$

$(x_{\text{test}}^{(1)}, y_{\text{test}}^{(1)})$
 $(x_{\text{test}}^{(2)}, y_{\text{test}}^{(2)})$
 \vdots
 $(x_{\text{test}}^{(m_{\text{test}})}, y_{\text{test}}^{(m_{\text{test}})})$

Andrew Ng

Training/testing procedure for linear regression

→ - Learn parameter θ from training data (minimizing training error J) 70% data

→ Compute test set error (on 30% data), could use cost function

take parameters learned from first step and plug in here and compute test set error.

>>>Error definition: for linear regression, could use cost as error: sum squared error over all test examples.

(note:
1. cost function: is used as optimizing object , to get weight more close to labeled y
2. Error: could be anything judging output good or not: could use
>>2.1 cost function measure distance of prediction-label
>>2.2 : accuracy, recall, precision..F1 etc.

→ - Learn parameter θ from training data (minimizing training error $J(\theta)$) 70%

- Compute test set error:

$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_\theta(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$$

3. Training / test procedure for logistic regression

same procedure as liner regression. only cost function defination is different.

> Learn paramter θ from training data (minimizing training error J) 70% data

>Compute test set error (on 30% data)
take parameters learned from first step and plug in here and compute test set error.

>>>Error definition for logistic regression:

>>Misclassification error (0/1 misclassification error):

Alternative test sets metric:

0/1: get an example right /wrong.

e.g: hypothesis mislabeled the example:

>>>error (h(x), y)=1:

If $h(x) > 0.5$, $y=0$;

or if $(h(x)<0.5$, $y=1$)

>>>error (h(x), y)=0: other cases.

Test error could define as: Test error = $1/m_{\text{test}} * \sum (\text{error}(h(x_i), y_i))$ over m_{test} data
why not use cost as error for logistic regression ?? (cost of logistic regression: is maximizing probability to optimize weight-->understanding/visualization' not good as accuracy here;

Training/testing procedure for logistic regression

→ - Learn parameter θ from training data

- Compute test set error:

$$J_{\text{test}}(\theta) = -\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} y_{\text{test}}^{(i)} \log h_\theta(x_{\text{test}}^{(i)}) + (1 - y_{\text{test}}^{(i)}) \log (1 - h_\theta(x_{\text{test}}^{(i)}))$$

- Misclassification error (0/1 misclassification error):

$$\text{err}(h_\theta(x), y) = \begin{cases} 1 & \text{if } h_\theta(x) \geq 0.5, y=0 \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Test error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_\theta(x_{\text{test}}^{(i)}), y_{\text{test}}^{(i)}).$$

10_79. Model selection and training / validation /test sets_Advice for applying Machine learning

Model selection problems:

decide what degree of polynomial / features (neural network: hidden layer number, hidden layer units number) to include to fit to a data set, or regularization paramter λ .

Overfitting example



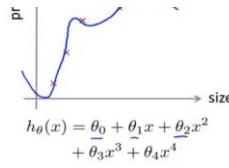
Once parameters $\theta_0, \theta_1, \dots, \theta_d$

1. Overfitting example

hypothesis parameters trained in training set, does well in training set does not mean how well hypothesis (or decision boundary-decided by parameters) will generalize to new examples not seen in the training set.

General principle:

Once parameter well fits to some set of data, then the error of hypothesis is measured on that same data set, that's unlikely to be a good estimate of actual generalization error. (how well hypothesis generalize to new samples).



were fit to some set of data (training set), the error of the parameters as measured on that data (the training error $J(\theta)$) is likely to be lower than the actual generalization error.

2. Model selection (Linear regression, logistic classification)

2.1 New parameter: polynomial degree d.

choose polynomial order /degree to fit data.: d

except parameter θ , have one more parameter d trying to determine using data set.

2.2. Process:

> choose one degree model;

> fit that model :

on training set, get parameter θ_d

> get some estimate of how well fitted hypothesis was generalize to new examples (dev set)..

compute test error.

take hypothesis corresponding to each trained ploynomial degree model, and measure performance on the test set.

>choose model (polynomial degree) has the lowest test set error.

note:

could d combine into hyperameter metrics, and search hyperameters to use (random choice-not grid, corase to fine method.)

2.3: choosed model generalization performance

One way: look at the test error this chosen polynomial hypothesis had done in test set (which used to choose polynomial degree).

Yet is not a fair way to estimate how well hypothesis generalizes, as we've fit this parameter d using test set - chose the d value that give the best possible performance on the test set.

-->performance of paramter θ (for the model d chosen) on the test set, that's likely to be an overly optimistic estimate.

hypothes is likely to do better on this test set than it would on new exampels it hasn't seen before.

test set and dev set are same distribution, -->dev set feature shoud same as test set, -->using dev error as test error still a problem?-->algorithm hyperparameter fitted on dev set, -->evalutation final choosed algorithm on test set.

3. Evaluation hypothesis

Dataset:

>Training set: 60% data-->m_training

>Cross validation set (cv/v set): 20% data-->m_cv

>Test set : 20% data-->m_test

Evaluating your hypothesis

Dataset:

Size	Price
2104	400
1600	330
2400	369
1416	232
3000	540
1985	300
1534	315
1427	199
2044	315
1380	212
1494	243

2044 is highlighted in red. Labels: Training set, Cross validation set (cv), Test set.

*(x⁽¹⁾, y⁽¹⁾)
(x⁽²⁾, y⁽²⁾)
⋮
(x^(m), y^(m))*

*(x_{cv}⁽¹⁾, y_{cv}⁽¹⁾)
(x_{cv}⁽²⁾, y_{cv}⁽²⁾)
⋮
(x_{cv}^(m_cv), y_{cv}^(m_cv))*

*(x_{test}⁽¹⁾, y_{test}⁽¹⁾)
(x_{test}⁽²⁾, y_{test}⁽²⁾)
⋮
(x_{test}^(m_test), y_{test}^(m_test))*

Andrew Ng

4. Train / validation / test error

cost function on different set: training error, validation error and test error.

Train/validation/test error

Training error:

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

Cross Validation error:

$$\rightarrow J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

Test error:

$$\rightarrow J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

5. Model selection:

use validation set instead of test set to select model (polynomial degree)

> for each potential chosen polynomial hypothesis di:
train on training set, Min J(θ)--> $\theta_{trained}$ -->compute $J_{cv}(\theta_{train_di})$

> pick model of the lowest $J_{cv}(\theta_{train_di})$: fit polynomial degree using validation set.

>Estimate generalization error for test set $J_{test}(\theta_{train_d})$
(note: generalization error: trained & choosed final algorithm performance on new examples not seen before)

Model selection

- 1. $h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \min_{\theta} \mathcal{J}(\theta) \rightarrow \theta^{(1)} \rightarrow \mathcal{J}_{cv}(\theta^{(1)})$
- 2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^{(2)} \rightarrow \mathcal{J}_{cv}(\theta^{(2)})$
- 3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \theta^{(3)} \rightarrow \mathcal{J}_{cv}(\theta^{(3)})$
- ⋮
- 10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)} \rightarrow \mathcal{J}_{cv}(\theta^{(10)})$

Pick $\theta_0 + \theta_1 x + \dots + \theta_4 x^4$

Estimate generalization error for test set $J_{test}(\theta^{(4)})$

Note:

many people report test error of hypothesis acutally on dev/test set (no dev set). If have massive test set, maybe not a terrible thing to do (less overfitting dev set). Better do not so do and practice to have seperate train , validation and test sets.

10_80. Diagnosing bias vs variance_Advice for applying Machine learning

If algorithm does not doing well, almost all the time it will be because have either a high bias problem, or a high variance problem: either underfitting or overfitting

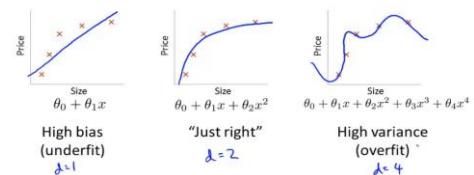
Bias/variance

1. Bias / variance intuition (from linear regression model /logistic classification model)

d: polynomial degree

High bias: underfit (when d is small)

High variance: overfit (when d is big)



2. Bias / variance definition

Train / validation / test error:

cost function on different set: training set, validation set and test set.

2.1 Training and dev error vs. polynomial degree

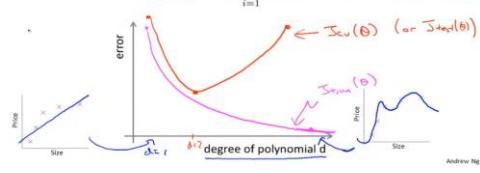
d is small --> underfitting: J_{train} is big, J_{cv} is big;
d is just right --> just right: J_{train} is small, J_{cv} is small
d is big --> overfitting: J_{train} is smaller, J_{cv} is big.

J_{train} will decrease as d increase: fitting training set better with higher polynomial hypothesis.
 J_{cv} will be high when d is too small or big: not able to fit well when overfitting or underfitting.
 J_{test} : similar with J_{cv}

Bias/variance

Training error: $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$

Cross validation error: $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$ (or $J_{dev}(\theta)$)



3. Diagnosing bias vs. variance

If algorithm performing less well: J_{cv}/J_{test} is high-->d is too small or big

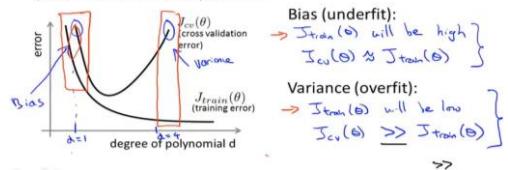
> High bias: d is too small:
 J_{train} is also high; J_{cv} is maybe slightly higher than J_{train} ;

> High variance: d is too big:
 J_{train} is low: fitting training set very well;
 J_{cv} is much bigger >> than J_{train}

(note: when diagnosing, J_{train} , J_{cv} : regularization item not involved in calculation)

Diagnosing bias vs. variance

Suppose your learning algorithm is performing less well than you were hoping. ($J_{cv}(\theta)$ or $J_{test}(\theta)$ is high.) Is it a bias problem or a variance problem?



10.81. Regularization and bias / variance_Advice for applying Machine learning

For preventing , regulation added to keep the values of the parameters small.

1. Linear regression with regularization

regularization parameter: λ

for chosen model (polynomial degree d):

> High bias (underfit):

large λ -->all parameters θ are heavily penalized, -->small parameter (θ close to 0)
---> hypothesis end up more or less a flat, constant straight line, close to θ_0 (do not include in regulation)

this hypothesis do not fitting well data set.

> High variance (overfit)

Small λ -->large parameter

if λ close to 0, given that we are fitting a high order polynomial, this is a usual overfitting setting.
given a high order polynomial, basically without regularization or very minimal regularization, end up with high variance overfitting setting.

Question: if $\lambda=0$, given that we are fitting a just right /low order polynomial-->result in underfitting?

2. Choosing the regularization parameter

Training / validation / test error: J_{train} , J_{cv} , J_{test}

cost function without regularization term on different set: training set, validation set and test set. Note: for these three error calculation, Without regularization term;

Cost function: $J(\theta, \lambda)$: consider regularization term.

Note:

cost function used to finding parameter to make hypothesis function h fit training set well, need to involve regularization into computation while training algorithm purpose: penalizing weights.-->increase J_{train}

while training /dev/test error is estimating trained hypothesis performance.
for linear regression, error also used squared error over all test examples, but without regularization, to reflect real error.

2.1 Selecting process:

>1. Try different λ for cost function $J(\theta, \lambda)$: start from 0, then 0.001 (then may step up in multiples of two until some larger value)

>> for each λ value: similar process with choosing polynomial d:

>>> min cost function $J(\theta)$ with regularization term on training set ---> $\theta_{train,\lambda}$

>>> Compute validation error $J_{cv}(\theta_{train,\lambda})$ on dev set .

Note: in validation error , regularization term not considered.

>>> pick the lowest $J_{cv}(\theta_{train,\lambda})$ -->hypothesis function $h(\theta_{train,\lambda})$

>Report test error: generalize selected hypothesis function to test set which hypothesis function did not see before..

similar to polynomial degree selecting, using training set, dev set and test set.

(note:

training set: -->using optimizing object, find weights, minimize cost

dev set: -->choose hyperparameter , based on algorithm dev error

test set: --> evaluate final fixed algorithm performance/test error.)

3. Bias / Variance as a function of the regularization parameter λ

how training error and dev error vary as vary the regularization parameter lambda λ

original cost function: $J_{o_}(\theta, \lambda)$: consider regularization term.

Training error and validation error: $J_{train}(\theta)$, $J_{cv}(\theta)$: without considering regularization term

> 1. λ is small: not using much regularization

>>> high risk of overfitting (given high order polynomial)

$J_{o_}$ is small?

J_{train} : is small; (as $J_{o_}$ and J_{train} fit training set well with high order polynomial)

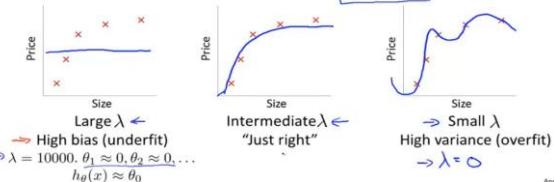
J_{cv} : is high

actually without regularization, so hypothesis performance do not change : underfitting or just right or overfitting?? decided by polynomial degree d

Linear regression with regularization

Model:
$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

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



Choosing the regularization parameter λ

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

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

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

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

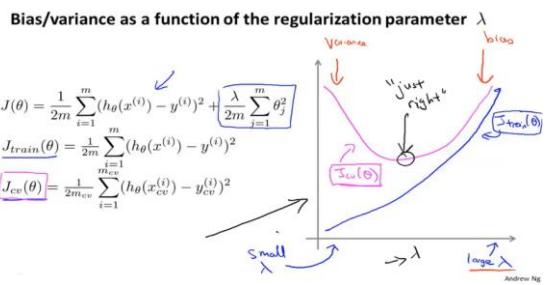
$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_\theta(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

Choosing the regularization parameter λ

Model:
$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

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

1. Try $\lambda = 0 \leftarrow \min J(\theta) \rightarrow \theta^{(0)} \rightarrow J_{cv}(\theta^{(0)})$
 2. Try $\lambda = 0.01 \leftarrow \min J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
 3. Try $\lambda = 0.02 \leftarrow \min J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
 4. Try $\lambda = 0.04 \leftarrow \min J(\theta) \rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$
 5. Try $\lambda = 0.08 \leftarrow \min J(\theta) \rightarrow \theta^{(4)} \rightarrow J_{cv}(\theta^{(4)})$
 - ⋮
 12. Try $\lambda = 10 \leftarrow \min J(\theta) \rightarrow \theta^{(10)} \rightarrow J_{cv}(\theta^{(10)})$
- Pick (say) $\theta^{(5)}$. Test error: $J_{test}(\theta^{(5)})$



Overfitting: - decided by polynomial degree d.

> 2. λ is big: heavy penalizing parameter θ

>>> high risk of underfitting (θ all close to 0, hypothesis nearly a flat straight line $\sim=\theta_0$ -which not involved in regularization)

J_o with regularization term=big, as underfitting.

(object is to minimize J_o ; while first part is J_{train} -high, second part regularization term-small as θ close to 0)

J_train is big: first part of J_o : underfitting, could not fitting training set well.

J_cv: is high.

4. Summary:

J_{cv} : always without regularization term, is high in both underfitting or overfitting;

J_{train} : always without regularization term, is high in underfitting, and low in overfitting;

J_o : consider regularization term if have, should be low both underfitting or overfitting?

J_o is only used to selecting parameter θ / train algorithm on training set. -> do not compare vs J_{cv} or J_{train} , which used for selecting regularization parameter or polynomial degree.

for some data set, could look at the **plot of the whole cross validation error**, then either manually, or automatically try to select a point, that minimize the cross validation error, and select the value of λ that corresponding to low cross validation error.

when selecting regularization parameter λ , plotting a figure like right, helps understanding better what's going on and helps verify that whether picked a good value λ or not.

10.83. Learning curves / variance_Advice for applying Machine learning

Ploting to check if algorithm is working correctly or to improve algorithm performance.
Very often used to diagnostic if algorithm is suffering from bias, variance or a bit both.

Which condition could have both bias and variance error?

note:

1. debug:

>> 1.1 linear regression/logistic classification: plot cost vs iteration: decrease or not-->reflect gradient descent computation right or not

>> 1.2 : neural network: check gradient numerically vs computed derivative by algorithm: as cost could not reflect when algorithm compute derivative wrong--may still decrease with iterations: cost function of parameter is not convex

2. diagnostic: bias, variance

J_{train}, J_{dev} : work both for neural network and linear regression/logistic classification

1. Learning curves, m number of training example

1.1 Plot:

1. J_{train} : Training error: average squared error for linear regression on training set. or logistic cost function(or use average error times ?) without regularization term.

2. J_{cv} : on validation set. average squared error for linear regression on validation set. or logistic cost function without regularization term.

vs: m-number of training example.

(note:

1. for choosing hyperparameter based on test error or dev set:

test error : could use cost , accuracy, recall, precision, F1 , etc: -->to judge algorithm performance we cared, better or worse

2. for diagnostic: bias, variance

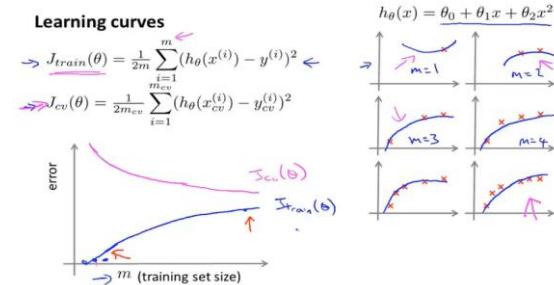
could also use same test error : accuracy, recall, ...

)

train/ validation / test error: $J_{train}, J_{cv}, J_{test}$

cost function without regularization term on different set: training set, validation set and test set. Note: for these three error calculation, Without regularization term;

Cost function for training algorithm: $J(\theta, \lambda)$: consider regularization term.



1.2. Ploting process: m number of training example

for a selected hypothesis function (e.g.: quadratic function like right pic):
like polynomial degree is already selected;

for each selected m : (take m to different value 1, 2, 3.etc.)

> train $h(x)$ on m numberof training set by minimizing cost function $J(\theta, \lambda)$, get trained parameter $\theta_{train,m}$:

>compute $J(\theta_{train,m})$ on validation set--> J_{cv}

and $J(\theta_{train})$ on training set--> J_{train} .

-->training error and validation error with only m number of training examples.

dev set should not be m size .

1.3. Ploting result: m number of training example

>1. if m is small--->over fitting

>>> J_{train} : small:

hypothesis function could fit training set perfect well when m is very small (eg.: m=1, 2, 3), with regularization term or not. $J_{train}(\theta_{train})$ is small (every easy to fit single training example when m is small)

>>> J_{cv} is big:

trained hypothesis function could not generalize well to data did not seen before.

>2. if m is large ---->not like underfitting, as J_{cv} is small

$H = \theta^T * (x_1, x_2, \dots, x_m)$:

example x_1, \dots, x_m used to constrain parameter θ . after certain size , new added training example could not add more constrain to parameter : do no more contribute to example space matrix (x_1, \dots, x_m) ; note $m > n$, matrix is singular.-->adding more training example do not help get parameter fit training set / dev set better could not add more constrain to algorithm parameters.

>>>1. J_{train} : larger

when m become larger, it hard for hypothesis function $h(x)$ to process through / fit all examples perfectly. J_{train} become larger.

-->actually J_{train} increase as m increase

>>>2. J_{cv} is become smaller:

only when get large training set, could the trained hypothesis function maybe fit data do not seen before better.

the more data you have, the better you do at generalizing to new examples.

-->actually J_{cv} decrease as m increase.

(more training set-->training set features space could cover dev& test set feature-->algorithm fitting dev/test well)

>>> 3. J_{cv}, J_{train} will flatten out with m increase: $J_{cv,f} \sim= J_{train,f}$

now training example space matrix (x_1, \dots, x_m) almost covers all dev/test set-sub example space, -->algorithm performance on these two example space almost same.--no matter undertaking, just right or overfitting--both large or small , or medium value.

note: even huge training size m , does not change the algorithm (polynomial degree) property of underfitting or overfitting for all training examples feature space.

with huge m, J_{cv}, J_{train} flatten out, and close to each other:

-->1. algorithm is undertaking: d is small--> $J_{train} \sim= J_{cv}$, end up with big error value

-->2. algorithm is just right: d is just right--> $J_{train} \sim= J_{cv}$, end up with small error value

-->3. algorithm is just overfit: d is very big-->overfit training set, dev set--> $J_{train} \sim= J_{cv}$, end up with smaller error value.

4. High bias (underfitting: example feature > hypothesis feature/degree d)

given hypothesis is linear functon as right:

Increase training example number m:

will not change much hypothesis function -straight line (May same line gotten with less training example).

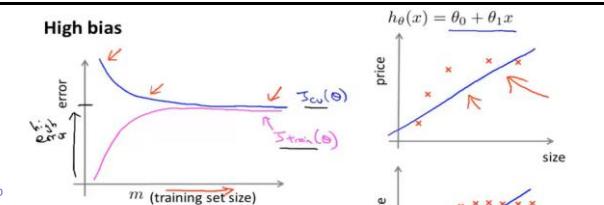
Just get the best straight line to fit training set, which could not fit this training data well.

--> J_{cv} do not decrease much: (when m=1, J_{cv} is large)

and when m come to certain value (could not add more feature/constraint to $h(x)$), you have almost fit the best possible straight line , which nearly do not change even m further go to very larger, --> J_{cv} plateau out/flatten out. ($h(x)$ nearly no change, then it's error on test set no change)

J_{train} : increase with m increase

J_{train} be small when m is small nearly 0, and then increase with m increase,-->in high bias case, J_{train} will end up to close to the $J_{cv,f}$ - $J_{cv} \sim= J_{train}$ both high error



If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much.

Note: in high bias case:

1. as $h(x)$ parameter number/feature/degree d is very small, when m is large, the performance on training set and validation set will be very similar (when m is large, training set features have big possibility including all test set features).--> $h(x)$ error on training features, will similar on test features)
2. In high bias case, the final flatten value: J_{cv}, J_{train} are high values. underfitting, $h(x)$ polynomial degree/computed new feature could not learn good functions to fit data set, could not fix by increasing data size.
3. Increasing training example, will not help much high bias issue. **increased training examples, will end up almost same $h(x)$ parameters which get in less training example.**

5. High variance

high variance: hypothesis have high polynomial degree d-->many parameters/features > huge data features.

$h(x)$ function like right pic: with very small regularization term--nearly no this term

>1. m is very small: overfitting.

$J_{train} \approx 0$: $h(x)$ fit training data very well with a function that overfit this

J_{cv} very high as is most overfitting.

>2. m increase a little bit: -still -overfitting

J_{train} increase:

still overfitting training data, but less overfitting- slightly harder to fit training data perfectly.--> J_{train} increase. but J_{cv} is small due to overfitting.

J_{cv} decrease much: as $h(x)$ less overfitting with m increase, generalize better to set error with m bigger.

but J_{cv} is big value due to overfitting, and now gap between J_{train} and J_{cv} is big.

(trained $h(x)$ only fitted to few data set feature included in training set, could not better generalize to test feature-->error gap between training error and test error is big)

>3. m increase - size including all example features (still less than $h(x)$ degree)-->overfitting

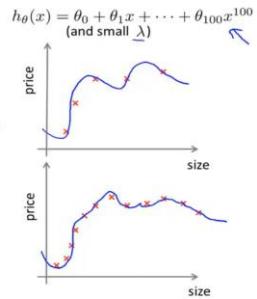
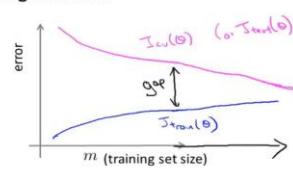
J_{train} increase:

still overfitting training data, but less overfitting- harder to fit training data perfectly.--> J_{train} increase.

J_{cv} decrease much: as $h(x)$ less overfitting with m increase, generalize better to dev set with m bigger.

(note: training set feature space: is cover more and more of dev set and test set)

High variance



If a learning algorithm is suffering from high variance, getting more training data is likely to help.

>4. m is large:-still overfit

increased training set could not add more constraints/features to hypothesis parameter, $h(x)$ parameter θ almost no change

--> J_{train} come to flatten value J_{strain_f} (not big as still overfit)

--> J_{cv} come to flatten value J_{cv_f} , close to J_{strain_f}

$h(x)$ could fit almost all data set features, overfitting. -->should also overfit test set--> J_{cv_f} is also small and close to J_{train_f} .

with m increase , J_{train} increase slow (J_{train_f} small as still overfit), J_{cv} decrease much (close to J_{train_f})

Summary :

in high variance

1. Getting more training data is likely to help.

esp. training data that could have new features/constraints to $h(x)$ parameter θ /degree d .

high variance: $h(x)$ parameter/feature/degree is very high, while training data volume is small/do not have enough constraint to each feature's weight. (data set feature number>> data volume-data that can add constraint to feature weight)

getting more training data-->mainly increase training data feature space: $n \times m$: -->extreme case: training set (x_1, \dots, x_n) 线性无关, --> $n \times n$ space, 基向量-->cover any further new example feature: $x(k) = [x(k)_1, \dots, x(k)_n]$

when training data including all data set feature, adding training data make no sense, training data space matrix (x_1, \dots, x_m) is singular after certain size / value of m .

Note:

real plotting curve of J_{cv}, J_{train} may be a little noise, messier than above ideal case. But plotting like this could tell if algorithm have bias or variance or both problem.

10.84. Decide what to try next_Advice for applying Machine learning

1. Debugging a learning algorithm

Problem: hypothesis on new set , have large error in its predictions (J_{cv}/J_{test} is big).-->improve algorithm

Method:

>Get more training examples --->only fix high variance

plot error pic and figure out have at least a bit of variance: means J_{cv} is quite a bit bigger than J_{train}

>Try smaller set of features in hypothesis--->only fix high variance

plot error pic and if find have a high variance, carefully select out small set of features to use.fewer features could help.

>Try getting additional features in hypothesis--->fix high bias (underfitting)

Current hypothesis is too simple (feature/degree to small), get additional feature in hypothesis to make it fit better to training set.

plot error pic and if find have a high bias

e.g. for predicting house price: get more info. about the pieces of land .and so on.

>>before doing above things, helpful to know in advance if this is going to help.

Debugging a learning algorithm:

Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis in a new set of houses, you find that it makes unacceptable large errors in its prediction. What should you try next?

- Get more training examples → fixes high variance
- Try smaller sets of features → fixes high variance
- Try getting additional features → fixes high bias
- Try adding polynomial features (x_1^2, x_2^2, x_1x_2 , etc) → fixes high bias.
- Try decreasing λ → fixes high bias
- Try increasing λ → fixes high variance .

>Try adding polynomial features (x_1^2, x_2^2, x_1x_2 , etc): fix high bias another way of adding features in hypothesis.

>Try decreasing regularization parameter λ -->fix high bias

decrease penalty on $h(x)$ parameter, so to fit better on training set and test set.

>Try increasing regularization parameter λ -->fix high variance

Note:

1. High variance problem:

>Definition: big gap between J_{cv} and J_{train} . Hypothesis function performance on training set and test set has big difference;

Root cause:

Training set feature number is not big enough, could not cover test features;

Too many data features considered in hypothesis (high degree), yet training set has no enough constraints to feature weights.

or algorithm captured additional fake features of training set , which not applicable on test set feature

(some features' weight on fit training set, which has no enough constraint.)

Solution:

>1. Increase training set (new feature vector):-->esp. with existing training set feature vector 线性无关-->increase feature

vector space-->adding more constraint to feature weights

>2. If increasing training set could not improve any more, then try regularization -reduce feature weights (esp. for the fake features)

(training set big enough--> feature vector-->feature matrix nxn, if still overfitting, means there are feature x, k 与 other feature 线性相关, feature space < nxn, feature matrix no enough constraint for feature weights

-->reduce fake features' weight

2. High bias problem:

>Definition: big J_{train} . Hypothesis could not even fit training set well.

Root cause:

Trained hypothesis model could not capture enough training set feature.

(data set feature number used in hypothesis is too small -->polynomial degree small, no enough to capture data interesting property)

Solution:

>1. Model architecture: try big polynomial degree, more layers, more units in layer;

>2. iterates more steps, get better $h(x)$ parameter θ .

2. Neural networks and overfitting

Neural network pattern chosen.

2.1: 'Small' Neural network: (few new created features) prone to underfitting

small network: few layers, and fewer units in each layer.-->few parameters, more prone to underfitting

computationally cheaper

2.2: 'Large' Neural network, (more features, prone to overfitting)

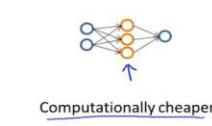
large network: more hidden layers or more hidden units in layers.

computationally more expensive.

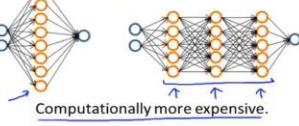
Normally the larger network, the performance better. could use regularization to address overfitting.

Neural networks and overfitting

→ "Small" neural network (fewer parameters; more prone to underfitting)



→ "Large" neural network (more parameters; more prone to overfitting)



Use regularization (λ) to address overfitting.

$J_{cv}(\theta)$

Note:

> Normally, Large neural network + regulizaiton more efficient than small neural network.

>>> Hidden layer number selecting:

same as selecting polynomial degree/ regulizaiton parameter, try neural network with different layers. for each chosen layer neural network: training set, validation set, test set.

train algorithm on training set-->calculate J_{cv} error on validation set --->selecting the lowest J_{cv} , and corresponding layer ---->report J_{test} for the chosen layer neural network error on test set.

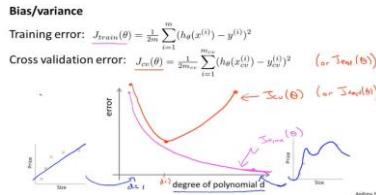
Summary_linear regression/logistic classification

for a given size training example m:

1. polynomial degree d influence:

algorithm performance will change from underfitting to overfitting , when d increase from small to large

if overfitting-->could reduce d: dropout,



for a given size training m, and selected polynomial degree d_s.

2. Regulation parameter lambda influence:

algorithm performance will change from 'overfitting' to underfitting with parameter increase (reducing parameter theta from huge large to 0, there is just right lambda to make just right J_cv).

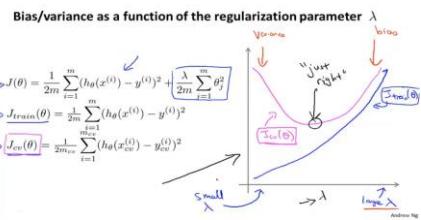
'overfitting': regulation parameter=0, no regularization term.

'underfitting': regulation parameter high;

but if d0 degree polynomial hypothesis is underfitting itself, then performance change should start from 'underfitting' to worse underfitting with lambda increase.

if overfitting (high polynomial degree d)-->increase regulation term (right lambda) could help high variance /reduce J_cv to just right point.

if underfitting-->decrease regulation term (right lambda) could help reduce high bias problem



if without regulation, h(x) with selected d_s, m training set,

>1. Underfitting: h(x) feature less than training set
>>add regulation term lambda_0, will penalize h(x) parameter, --> trained h(x) now worse fitting training set and test set-->add regulation term to undertake make J_train increase and J_cv increase (still underfitting)
J_cv decrease and J_cv decrease (still underfitting)
>>increase lambda_0, J_train increase and J_cv increase (still underfitting)

>2. overfitting: h(x) feature more than training set feature:

>>add/increase regulation term lambda_0, will penalize h(x) parameter, --> trained h(x) feature weight smaller, worse fit each training set all feature and but better fit on test set feature
>> J_train increase and J_cv decrease (if parameter is not 0, then end up still with a little bit overfit, J_train_f is very small, while J_cv_f close to J_train_f)
>>decrease lambda_0, more overfitting J_train decrease and J_cv increase (still overfitting, J_train not much)

for selected polynomial degree d_s.

d_o= just right polynomial degree for huge training example

3. influence of training example number m

J_train: increase from very small to large till nearly flatten , with m increase;
J_cv: decrease from high value continually then to flatten , with m increase

if d_s > d_o:

note:

when overfitting: increase training example could help.
as when overfit, example feature number involved in hypothesis is too many, not enough training set volume not big enough /feature matrix (x1..xm) no big enough to give constraint for every feature/hypothesis parameter.
Add training example number, especially example that add new features vector (add new constraint -与现有 training set/feature vector 线性无关), could help overfitting;

when underfitting: increasing training example could not help.
as underfit: example feature number hypothesis degree/feature too small, adding more training example may slightly increase J_train, decrease J_cv slightly , not much help

if d_s < d_o:

J_train: increase from very small to large till nearly flatten , with m increase;
J_cv: decrease from high value continually then to flatten , with m increase

>m: 1-just right number (constraint h(x) feature d):

J_train increase from nearly 0, J_cv decrease from high value

>m: continue increase: could not add more feature, both J_train and J_cv come to flatten value: these two values close, but normally J_cv_f slightly higher than J_train_f.

if d_s < d_o:

J_train: increase from very small to large till nearly flatten , with m increase;

>>m: 1-just right number (constraint h(x) feature d):

when m is very small, overfitting, training example is fitted perfectly J_train almost 0, when m number increase, less overfitted-fit training example less perfectly-->J_train increase.
When m increase to a certain value, that could just constrain d_s (polynomial features), get the trained h(x)_theta_o'

>>m: just right number-->large number: just have all example feature

m feature > h(x) feature, kind of underfitting, but compare to above h(x), current h(x) is learning to fit more example features, and fit less well on all training example than it did on previous smaller training example which have less features , but could fit better on test set features-->, J_train continue increase

>>m: continue increase:

at this time h(x) parameter theta nearly no change as new added example do not add more constraint: -->J_train come to nearly flatten value.

J_cv: decrease from high value continually then to flatten , with m increase

with m number increase, h(x) have been trained to fit 'huge training example' better-->error on test set is becoming less.

>> m: 1-just right number (constraint h(x) feature d);
J_cv decrease singly, as h(x) is fitting better on test set;

>>m: just right number-->large number: just have all example feature

m feature > h(x) feature, kind of underfitting, but compare to above h(x), current h(x) is learning to fit more example features, and fit less well on all training example than it did on previous smaller training example which have less features , but could fit better on test set features (test set could be any some features of data)-->, J_train continue increase, and J_cv continue decrease

Note: from nowon, h(x) with d_s degree is actually underfitting: could not solve underfitting problem with m increase:

increasing m will only add more example features, while h(x) learned to fit all features while without enough ability (h(x) degree not high enough)--->
J_cv will decrease but have bottom line (flatten value)-do not decrease much, and J_train will increase (also have top line)when increasing m.

>m: continue increase:

at this time h(x) parameter theta nearly no change as new added example do not add more constraint. h(x) error on test set nearly no change either- ->J_cv come to a flatten value : normally this value J_cv_f close to J_train_f, but slightly higher.

W12 - Machine learning system design

12_85. Prioritizing what to work on: Spam classification example_Machine learning system design

Foucs on Issues faced on Machine learing system design. Prioriting your time for what to work on.

Building a spam classifier

1. Building a spam classifier

using supervised learnign to distinguish between spam and non-spam based on training set (labeled as spam or not)

From: cheapsales@buystufffromme.com
To: ang@cs.stanford.edu
Subject: buy now!

Deal of the week! Buy now!
Rolex w4tches - \$100
Medicline (any kind) - \$50
Also low cost M0rgages available.

Spam (1)

From: Alfred Ng
To: ang@cs.stanford.edu
Subject: Christmas dates?

Hey Andrew,
Was talking to Mom about plans
for Xmas. When do you get off
work. Meet Dec 22?
Alf

Non-spam (0)

2. Feature represent_for training set

>first: represent x as features of the email (while $y=1$ (spam) or 0 (not spam)):

represent method:

e.g: choose 100 words indicative of spam/not spam.

When there are those words in email-->judged as spam/not spam.

E.g when there are words'now' in mail, judge as not spam.

>> create feature vocabulary/dictionay : dimenion is the feature words number/size

>> create feature vector X for each mail: dimension is same as feature dictionay;

$0 / 1$ in feature vector means corresponding, same positon feature word in feature dictionay , not occurre /occurre in email.

Note: for feature words choosen:

normally look at the training set and take most frequently occuring words (10,000 -50,000) in training set, rather than manually pick 100 words.

Building a spam classifier

Supervised learning. $x = \text{features of email}$. $y = \text{spam (1) or not spam (0)}$.
 Features x_i : Choose 100 words indicative of spam/not spam.

E.g. deal, buy, discount, andrew, now, ...

$x = \begin{bmatrix} 0 & \text{andrew} \\ 1 & \text{buy} \\ 1 & \text{deal} \\ 0 & \text{discount} \\ \vdots & \vdots \\ 1 & \text{now} \\ \vdots & \vdots \end{bmatrix}$

$x_i = \begin{cases} 1 & \text{if word } i \text{ appears} \\ 0 & \text{otherwise.} \end{cases}$

From: cheapsales@buystufffromme.com
To: ang@cs.stanford.edu
Subject: buy now!

Deal of the week! Buy now!

Note: In practice, take most frequently occurring n words (10,000 to 50,000)
 in training set, rather than manually pick 100 words.

3. Priority for making classifier having lower error

>1. Collect lots of data

e.g: 'honeypot' project:

creat fake email address and get them into the spammers hands, then collect tons of spam email.

>2. Develop sophisticated features based on emali routing info. (from email header).

spammer often will try to obscure email origins using maybe fake email addres or send through very unusual sets of compute service..etc. some of these info. will reflect in email header.-->new feature based on email header to caputer email routing info. is an efficient way to identify if it is spam or not.

>3. Develop sophisticatd features for message body.

e.g how about 'deal' and 'dealer', treated as same word? spam may have particular favor for word/punctuations.

>4. Develop sophisticatd algorithm to detect/correct misspellings

e.g: 'w4tches' in mail, spam classifier might not equate this as the same thing as the word 'watches' (which choosen as feature word), harder to realize this mail is spam with these deliberate imspellings. that's why spammers do it.

use error analysis , to tell where can try to have a more systematic way to choose among options of the many different things you might work on.-->More likely to select what is a good way to go for further step.

Building a spam classifier

How to spend your time to make it have low error?

- Collect lots of data
 - E.g. "honeypot" project.
- Develop sophisticated features based on email routing information (from email header).
- Develop sophisticated features for message body, e.g. should "discount" and "discounts" be treated as the same word? How about "deal" and "Dealer"? Features about punctuation?
- Develop sophisticated algorithm to detect misspellings (e.g. m0rtgage, med1cine, w4tches.)

12_86. Error analysis_Machine learning system design

1. Recommended approach

>1. Start with a simple algorithm that can implement quickly.

Implement it and test it on cross-validation data.

Usually spend 1 day to try get something quick and dirty , and implement it and test on dev set.

>2. Plot learning curve to decide if more data, mroe features, etc. are likely to help.

figure out if have high bias/ variance or both using plotted curve.-->then try corresponding method to improve.

when just start machine learning problem, it is very hard to tell in advance whether need more data or more feature in the beginning without learning curve.

by running even a very quick and dirty implementation, and by plotting learning curves, that helps make decisions where to go --->let evidence rather than feeling...guide on where to spend time

Recommended approach

- - Start with a simple algorithm that you can implement quickly.
 Implement it and test it on your cross-validation data.
- - Plot learning curves to decide if more data, more features, etc. are likely to help.
- - Error analysis: Manually examine the examples (in cross validation set) that your algorithm made errors on. See if you spot any systematic trend in what type of examples it is making errors on.

>3. Error analysis:

manually examine the examples (in cross validation set) that algorithm made error on.
 see if you spot any systematic trend in what type of examples it is making errors on.

-->inspirer to design new features, or tell whether current things /shortcomings of the system etc. give inspiration you need to come up with improvements to it.

2. Error analysis

e.g: $m_{cv}=500$, algorithm make error on 100 validation example.

-->

Manually examine the 100 errors, and categorize them based on:

>Email type

counting up emails in different catogary, -->work on catogary that have more errors.

e.g: sell things (20pc), steal password(70pc).-->work on steal password deteting.

> Cues/ feauters of these mail, that may have helped the classifier class them correctly.

counting up emails in different features, -->work on features that have more errors.

e.g: deliberate misspellings (5pc), unusual email routing (16 pc), unusual punctuation (32 pc)-->work on email routing and punctuation.

Error Analysis

$m_{cv} = 500$ examples in cross validation set

Algorithm misclassifies 100 emails.

Manually examine the 100 errors, and categorize them based on:

- (i) What type of email it is → Pharma, replica, steal passwords, ...
- (ii) What cues (features) you think would have helped the algorithm classify them correctly.

Pharma: 12

Replica/fake: 4

Steal passwords: 53

Other: 31

→ Deliberate misspellings: 5

(m0rtgage, med1cine, etc.)

→ Unusual email routing: 16

→ Unusual (spamming) punctuation: 32

3. Principle

Implement simple algorithm quick and dirty-->find the most difficult examples to identify and focus effert on those.

Note: Different algorithm almost always find similar categories of examples difficult.

4. Importance of numerical evaluation

Incredibly helpful to have a way of evaluating learning algorithm that gives back a single rule number (maybe accuracy, maybe error) tell how well learning algorithm is doing.

Summary:

when developing new algorithm, need to try lots of new ideas and lots of new versions of algorithm. if every time try new idea, if end up manually examining a bunch of examples again to see better or worse, makes harder to make decision for this idea effect on algorithm performance.

single rule number could let you just look and see if the error go up or down, use it more rapidly to try out new ideas, and almost right away tell if new idea improved algorithm or not.

Note: do error analysis on dev set, instead of test set.

(note:

1. train/dev/test error: algorithm performance on training/dev/test set

2. performance/error representation could use: cost, accuracy, recall, precision, F1, etc.

3. error usage:

>> check bias/variance

>> as signal evaluation to choose new ideas/features

>> choose hyperparameters based on cost on dev

4. cost function vs iteration : usage:

--> debug

(only for linear regression, logistic regression, neural network use gradient check for debug as cost vs parameter is not convex)

--> if converged or not

The importance of numerical evaluation

Should discount/discounts/discounted/discounting be treated as the same word?

Can use "stemming" software (E.g. "Porter stemmer")

universe/university

Error analysis may not be helpful for deciding if this is likely to improve performance. Only solution is to try it and see if it works.

Need numerical evaluation (e.g., cross validation error) of algorithm's performance with and without stemming.

Without stemming: 5% error With stemming: 3% error

Distinguish upper vs. lower case (Mom/mom): 3.2%

12_87. Error metrics for skewed classes_Machine learning system design

Error metrics: a single real number evaluation metric for learning algorithm to tell how well it's doing.

Skewed classes: particularly tricky to come up with an appropriate error metric.

1. Skewed class example:

e.g. trained classifier have 1% test error (tested on test set)

while only 0.5% of test set actually have cancer. --> algorithm error 1% is not good.

If take unlearning algorithm: just let $y=0$, given whatever input x , --> test error then is only 0.5%, better than above learning algorithm

(note: for skewed class-->use unsupervised learning: anomaly detection, positive $y=1$ is small volume, save for dev and test set; training set is $y=0$)

Cancer classification example

Train logistic regression model $h_\theta(x)$. ($y=1$ if cancer, $y=0$ otherwise)

Find that you got 1% error on test set.

(99% correct diagnoses)

Only 0.50% of patients have cancer.

skew class
function $y = \text{predictCancer}(x)$
 $\rightarrow y = 0$; ignore x
return

0.5% error

→ 99.5% error (0.5% error)

→ 99.5% accuracy (0.5% error)

Andrea Ng

2. Using accuracy/error evaluating skewed classifier

Skewed classes: ratio of positive to negative example (in data: including training, dev, test) is very close to one of two extremes (positive is much more or negative is much more than the other class)

in this case unlearning algorithm (let y always =1/0) could do pretty well on prediction.

using classification accuracy / classification error (cost on dev) as evaluation metric for skewed class:
Even result in high accuracy / low errors, do not always clear if doing so is really improving classifier quality, like right e.g.

Precision/Recall

$y=1$ in presence of rare class that we want to detect

Actual class		Predicted 1 class	
1	0	True positive	False positive
0	False negative	True negative	False negative

$y=0$
Recall = $\frac{\text{True negative}}{\text{True negative} + \text{False negative}}$

$\frac{\text{True positives}}{\text{True positives} + \text{False positives}}$

$\frac{\text{True positives}}{\text{True positives} + \text{False negatives}}$

e.g.:
accuracy improve from 99.2% (0.8% error) to 99.5% (0.5% error):

may just simplify algorithm just let $y=0$ or 1 (accuracy higher, yet algorithm itself not improve /particularly good classifier)

3. Precision/ Recall_for skewed class

Definition:

True positive: test data class is positive (1), and algorithm prediction is also positive (1);

True negative: test data class is negative (0), and algorithm prediction is also negative (0);

False Positive: test data class is negative (0), yet algorithm prediction is positive (1);

False negative: test data class is positive (1), yet algorithm prediction is negative (0);

Precision: true positive / predicted positive(true positive + false positive)
(of all test set that predict positive, what fraction actually is positive)

Recall: True positive /actual positive (True positive + false negative)

(of all test set that is actually positive, what fraction did algorithm correctly detect as positive)

for unlearning algorithm y always =0 (above e.g.) - true positive=0->precision and recall=0

Note:

1. when defining precision and recall, usually use the convention that $y=1$ in the presence of the more rare class.

2. Precision/Recall metric is a good way to tell if algorithm even for skewed class is doing well or not.

3. for skewed classes, precision/recall is a better way than accuracy/error for evaluation algorithm performance.

can prevent unlearning algorithm 'cheating' like y always =0 which have high accuracy/low error.

for other un-skewed data, which evaluation value used often: accuracy / cost

11_88. Trading off precision and recall_Machine learning system design

2. Large data rationale

Conditons massive training set that will be able to help.

2.1 Feature x have sufficient info. for specific y-prediction.

training set feature could cover all test set feature;
does it also mean even small training set cover all the data set features?
even 1/2 training sample have show all the data features, need more/ enough training data to constrain hypothesis parameter
1. have enough feature cover all data set feature;
2. have enough constraints to hypothesis parameter/degree/feature;

- >
1. for linear regression/logistic classification: $y/\text{decision boundary} = \Phi_0 + \Phi_1 x_1 + \Phi_2 x_2 + \dots + \Phi_{n \times n}$
have enough feature info. x_1, \dots, x_n (for each example), to capture training example info. to predict feature/polynomial degree is just fit or overfit training & dev & test set
 2. for neural network: enough final 'new created feature' feeding to last layer-->
input feature vector should be enough, and neural network architecture: hidden layer number, hidden units should be big enough.

Hypothesis function need:

1. have enough parameter to capture all training set feature;
2. Capture less fake feature from training set, while not applicable to test set.

And use a **learning algorithm with many parameters (high feature number)-low bias algorithm**

e.g: logistic regression / linear regression with many features; neural network with many hidden units;
> J_{train} will be small (overfitting)
> using very large training set , will be unlikely to overfit (when training example number large than hypothesis parameter number, feature space bigger, better cover dev/test set feature vector)
--> and $J_{\text{train}} - J_{\text{cv}} / J_{\text{test}}$ f is small

Adding more training set/more features in training set, add more constrain to hypothesis, make less overfit, to just right point .

J_{train} increase to flatten, but still small as a little bit overfit in last;
 J_{cv} decrease to flatten, close to J_{train} (as training set cover all test set feature, only hypothesis may have more 'fake features' captured in training set , that not applicable for test set--> J_{cv} f a bit higher than J_{train} f)

Algorithm that no high bias, and no high variance:

No high bias: using algorithm with many parameters (high degree/many features)

No high variance: using huge large training set: (algorithm could generalize well from training set to dev set,--> training set is large enough to cover all dev/test set features)

1. Key point for getting low bias on test set:-->a little bit overfit

- >1. hypothesis function high degree/data set features ;
>2. training set big enough -->have enough features that cover all test set feature;

2. Get low variance on test set: huge training set.-->a little bit overfit

(training set have all features could cover test set , and these features captured by hypothesis function. while also try to make the 'fake features' captured by hypothesis, which could not apply to test set, as small as possible)

huge data to train algorithm with many parameters might help get a high performance learning algorithm.

Large data rationale

→ Assume feature $x \in \mathbb{R}^{n+1}$ has sufficient information to predict y accurately.

Example: For breakfast I ate two eggs.

Counterexample: Predict housing price from only size (feet²) and no other features.

Useful test: Given the input x , can a human expert confidently predict y ?

Large data rationale

→ Use a learning algorithm with many parameters (e.g. logistic regression/linear regression with many features; neural network with many hidden units). low bias algorithms.

→ $J_{\text{train}}(\theta)$ will be small.

Use a very large training set (unlikely to overfit) low variance

→ $J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$

→ $J_{\text{test}}(\theta)$ will be small

Summary:

Pre-judge,

> first , if human expert could predict confidently based on these training set/feature.

This is a sort of certification if y can be predicted accurately from input feature: input feature is enough, algorithm also have enough parameters.

> 2nd: if can get large training set, to train algorithm with many parameters. : to make training example feature space (x_1, \dots, x_m), nearly cover all test example.

if can do both, then could get a good performance algorithm