

2020/01/14 (Machine Learning) - Andrew Ng

→ supervised (right answers were given)

CLASSIFICATION
Date
Page
Decision (discrete)
regression (continuous) continuous

misclassified → audio example (separating vowel from noisy environment)

cost function: measure of accuracy of hypothesis

$$\text{cost funct. } J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

$$= \frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2.$$

difference between the predicted value & actual value

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

$$j = m\theta_0 + c$$

minimize $J(\theta_0, \theta_1)$

θ_0, θ_1

sq. error function

* Gradient Descent: α : learning rate

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

repeat this until convergence
→ simultaneously update θ_0, θ_1
 $\text{temp } \theta = \theta_0 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}$

$\text{temp } \theta_1 = \theta_1 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1}$ partial derivative
(not derivative)
 $\theta_0 := \text{temp}_0$
 $\theta_1 := \text{temp}_1$ simultaneous update step

$$\theta_0 := \theta_0 + \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}, \quad \theta_1 := \theta_1 + \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1}$$

$$\theta_0 := \theta_0 + \alpha \theta_1 = 1 + \theta_1$$

$$\theta_1 := \theta_1 + \alpha \theta_0 = 2 + \theta_0$$

lit. $J(\theta)$ is convex $\theta_0 \geq 0$ & $\theta_1 \in \mathbb{R}$



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

$\frac{\partial J(\theta_1)}{\partial \theta_1}$ → slope of tangent at θ_1

positive no. / positive slope
 $= \theta_1 - \alpha$ (true no.) , α - always true
 $= \theta_1 - \text{value}$ (decrease θ_1)

Similarly, $J(\theta_0)$



$$\theta_0 := \theta_0 - \alpha \frac{\partial J(\theta_0)}{\partial \theta_0} \quad (\text{slope } -ve)$$

$= \theta_0 - \alpha (-ve \text{ no.})$
 $= \theta_0 + \text{some value}$ (increase θ_0)

$\alpha \rightarrow$ too small: gradient descent

becomes slow
→ too large: gradient descent can overshoot the minimum (fails to converge)



what happens?

$$\theta_1 := \theta_1 - \alpha (\text{slope at } \theta_1)$$

$$= \theta_1 - \alpha \times 0$$

$$= \theta_1$$

θ_1 is not changed

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

$$y = \theta_0 + \theta_1 x$$

$$1 = \theta_0 + \theta_1 x$$

$$\theta_0 + 2\theta_1 = 1$$

$$\theta_0 + 1\theta_1 = 0.5x$$

$$(1 - 2\theta_1)$$

$$1 - 2\theta_1 + \theta_1 = 0.5$$

$$1 - \theta_1 = 0.5$$

$$\theta_0 + 2\theta_1 = 1$$

$$\theta_0 = 0$$

$$y = mx$$

$$y = 0.5$$

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

$$= -1 + 2x$$

$$= -1 + 2 \times 6$$

$$= 12 - 1$$

$$= 11$$

$$m=0 = n \times$$

$$= \frac{0}{n}$$

$$= 0$$

* multivariate linear regression:

hypothesis:

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

for convenience, $n=1$

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

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad \text{and } X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad (\text{+1 dimension vector})$$

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

* Feature Scaling:

$$x_1 = \text{size}(0 - 5000)$$

$$x_2 = \text{no. of bedrooms}(1 - 5)$$

x_1, x_2 has very high difference when making contour plots difficult and making calculation of gradient descent difficult.

$$\text{So, } x_1 = \frac{\text{size}}{2000}, \quad x_2 = \frac{\text{no. of bedrooms}}{5}$$

$$0 \leq x_1, x_2 \leq 1 \text{ (easy)}$$

* Mean normalization: Replace x_1 with $(x_1 - \bar{x}_1)$ except x_0

$$x_1 = \frac{\text{size} - 1000}{2000}$$

$$x_2 = \frac{\text{no. of bedrooms} - 2}{5}$$

Multivariate linear regression fit cost function

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

And?

Gradient Descent,

repeat until convergence

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

}

for $j = 0 \dots n$

→ features can be improved by combining multiple features. (Ex: $x_3 = x_1 \times x_2$)

→ Polynomial Regression: → if hypothesis does not fit the data, we can use non-linear representations too.

→ We can change the behavior or curve of the hypothesis function by making it quadratic, cubic, sq. root or any other form.

Ex: $h_\theta(x) = \theta_0 + \theta_1 x_1$.

quadratic: $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$

cubic: $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$

sqrt: $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_1}$

(Note: Features should be considered (numerous))

* Normal Equation

→ gradient descent is one way of minimizing the cost function but this can also be converted by normal equation.

→ In Normal Equation method, we minimize the cost function (J) by explicitly taking its derivatives with respect to θ_j (partial derivative) and setting them to 0.

→ Normal eqn allows us to find the optimum θ without iterations.

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

Normal eqn.

→ No need to scale the features!

Gradient Descent

→ need to choose α
→ needs many iterations
→ works well when n is large.
($n \rightarrow$ features)

→ $O(Kn^2)$ terms
(mathematical complexity)

• Similar to finding minima of a curve using

$$\frac{\partial J(\theta)}{\partial \theta} = 0$$

Normal Equation

→ no need of α
→ no need to iterate
→ slow if $n \gg 1$
large
→ $O(n^3)$ terms
(mathematical complexity)

* Normal Equation Non-invertibility

($X^T X$) - term could be non-invertible because of:

→ redundant features (when two features are very closely related) (maybe linearly dependent).

→ too many features ($m \leq n$), where m is the no. of samples and n is the no. of features (in this case, use "regularization" or delete some features)

* Classification

→ attempt classification one method is to use linear regression and map all predictions greater than 0.5 as 1 and less than 0.5 as 0. However, this method does not work well because classification is not a linear function.

→ classification problem is like the regression problem except that the values we now want to predict take on only formally no. of discrete values. (Ex: 0 and 1) → bringing classification problem

→ 1 → 1 called positive class
→ 0 → 0 called negative class
→ $x^{(i)}$ → features of data
 $y^{(i)}$ → aka label of data. $y \in \{0, 1\}$

* Hypothesis Representation

→ for classification problems, there is no sense for hypothesis function to output values beyond 0 and 1. as we know that $y \in \{0, 1\}$

→ To fix this, we change the hypothesis function $h_\theta(x)$ so that: $0 \leq h_\theta(x) \leq 1$

→ How to achieve this?
Ans: by plugging θ_0
into the logistic function.
(sigmoid function)

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

$$h_\theta(x) = \theta_0 + \theta_1 x_1$$

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

$$g(z) \text{ maps any real number to } (0, 1)$$

$$h_\theta(x) = P(y=1|x; \theta)$$

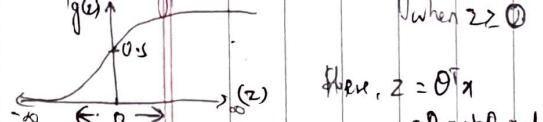
$$(h_\theta(x)) \text{ gives the probability that our output is 1.}$$

Decision Boundary

In order to get our discrete 0 or 1 classification we can think the output of the hypothesis function as follows:

$$\begin{cases} h_{\theta}(x) \geq 0.5 \rightarrow y = 1 \\ h_{\theta}(x) < 0.5 \rightarrow y = 0 \end{cases}$$

→ Logistic function (sigmoid) characterizes the output of hypothesis function: $g(z) \geq 0.5$



As, $z \rightarrow \infty$, $g(z) \approx 1$

As, $z \rightarrow -\infty$, $g(z) \approx 0$

As $z = 0$, $g(z) = 0.5$

→ Decision boundary i) the line that separates the area where $y=0$ and $y=1$.

$$\text{Let, } \theta^T = [5 \ -1 \ 0]$$

$$\text{Then } \theta^T x = 5x_0 + (-1)x_1 + 0x_2$$

$$\therefore \theta^T x = 5 - x_1 \quad (\text{as } x_0 = 1)$$

Now, for $y=1$, we know, $x=5 \text{ km}$

$$\theta^T x \geq 0$$

$$\therefore 5 - x_1 \geq 0$$

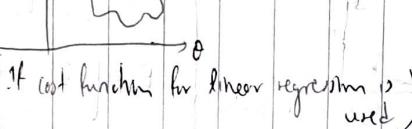
$$\therefore -x_1 \geq -5$$

$$\therefore x_1 \leq 5$$

$$\therefore n \leq 5$$

* Cost Function for Logistic Regression

→ if we use the same cost function i.e. $(h_{\theta}(x^{(i)}) - y^{(i)})^2$ for logistic regression, the logistic function will output a wavy curve (having many local minima (a non-convex function))



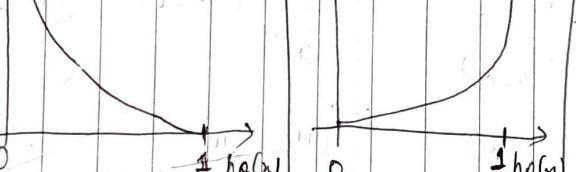
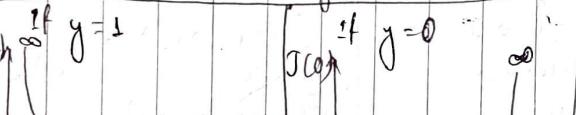
If cost function for linear regression is used

So, we use:

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

where,

$$\text{cost}(h_{\theta}(x^{(i)}), y^{(i)}) = \begin{cases} -\log(h_{\theta}(x^{(i)})) & \text{for } y=1 \\ -\log(1 - h_{\theta}(x^{(i)})) & \text{for } y=0 \end{cases}$$



$\text{cost}(h_{\theta}(x^{(i)}), y^{(i)}) \rightarrow 0$ if $h_{\theta}(x^{(i)}) = y^{(i)}$
 $\text{cost}(h_{\theta}(x^{(i)}), y^{(i)}) \rightarrow \infty$ if $y^{(i)} = 1$ and $h_{\theta}(x^{(i)}) \rightarrow 0$
 $\text{cost}(h_{\theta}(x^{(i)}), y^{(i)}) \rightarrow \infty$ if $y^{(i)} = 0$ and $h_{\theta}(x^{(i)}) \rightarrow 1$

→ writing the cost function in this form guarantees that $J(\theta)$ is convex for logistic regression

* Simplified Cost Function and Gradient Descent

→ Cost function can be compressed into a single equation such as:

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

When, $y=1$, $\text{cost}(h_{\theta}(x^{(i)}), y) = -\log(h_{\theta}(x^{(i)}))$

and when $y=0$, $\text{cost}(h_{\theta}(x^{(i)}), y) = -\log(1 - h_{\theta}(x^{(i)}))$

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

(entire cost function)

In vectorized implementation,

$$h = g(x\theta) \quad \text{where, } g(z) = \frac{1}{e^{-z} + 1} = \text{logistic function}$$

And

$$J(\theta) = \frac{1}{m} [-y^T \log(h) - (1 - y)^T \log(1 - h)]$$

* Gradient Descent:

$$\text{Repeat } \left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j} \end{array} \right\}$$

Also,

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

→ gradient descent for logistic regression

$$\theta := \theta - \alpha/m \sum_{i=1}^m x^T (g(x\theta) - y)$$

→ vector representation

(3)

* Advance Optimization

→ Conjugate Gradient", "BFGS" and "LBFGS" are more sophisticated and faster ways to optimize.

Θ Mat can be used instead of gradient descent.

→ Use library because they are tested & highly optimized

→ for a given value Θ, we need to calculate $J(\Theta)$ and $\frac{\partial J(\Theta)}{\partial \Theta_j}$

→ In octave, we can write a function to return these values as:

function [jval, gradient] = costFunction(theta)

jval = % code to compute J(theta)

gradient = % code to compute derivative of J(theta)

end

→ Octave has "fminunc()" optimization algorithm along with "optimset()" function that creates an object containing the options we want to send to "fminunc()".

```
Code:
options = optimset ('GradObj', 'on', 'MaxIter', 100);
initialTheta = zeros(2,1);
[theta, funVal, exitFlag] = fminunc(
    @costFunction, initialTheta, options);
```

* Multiclass Classification

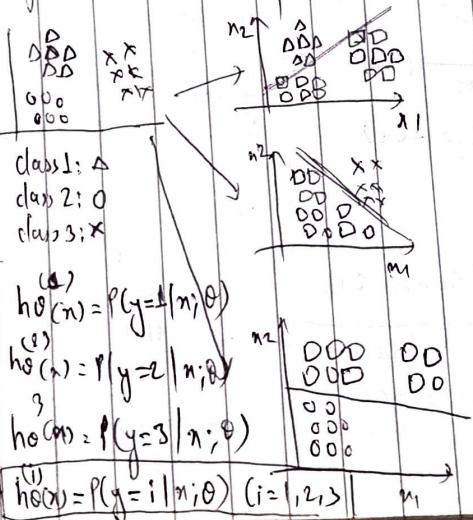
(One vs. all) (One vs rest)

→ for multiclass, $y \in \{0, 1, \dots, n\}$

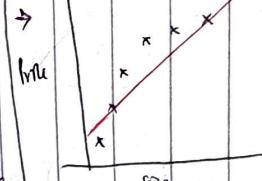
→ To achieve this, we divide our problem into $(n+1)$ binary classification problem (+1 because index starts at 0)

→ We are basically choosing one class and lumping all the others into a single second class.

→ We do this repeatedly applying binary logistic regression to each case and then use the hypothesis that returned the highest value as our prediction.



* Over-fitting



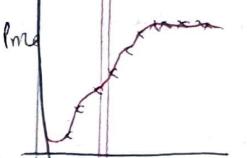
$$h_{\theta}(x) = \theta_0 + \theta_1 x_1$$

Under-fitting (high bias)*



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

"Slight fitting"



$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^3 + \theta_3 x_3^4$$

(Overfitting)*
(High-variance)

→ 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.

→ Underfitting: caused by a complicated function that creates a lot of unnecessary curves and angles unrelated to the data.

→ Underfitting: when hypothesis function maps poorly to the trend of data, usually caused by function that is too simple or uses too few features.

→ Address Overfitting: (i) Reduce the no. of features:

- manually selected the features to keep
- use a model selection algorithm

(ii) Regularization:

- keep all features but reduce the magnitude of parameters
- regularization works well when we have lot of slightly useful features

Date
page
classmate

* Regularization:

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

→ If $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$, to make $h_{\theta}(x)$ quadratic, we have to eliminate the influence of θ_3 & θ_4 .

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

→ The last two terms increases the cost of θ_3 and θ_4 significantly and thus to decrease the overall cost of the hypothesis, θ_3 and θ_4 should be greatly reduced. This results in a simple quadratic function where, $\theta_3 \approx 0$ and $\theta_4 \approx 0$

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

where λ is the regularization parameter

$\sum_{j=1}^n \theta_j^2$ is regularization component / term

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

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

} $\{ i \in \{1, 2, \dots, n\}$ (add to upgrade separate θ_j in θ)

→ Term $(1 - \alpha \frac{\lambda}{m})$ is always less than 1

→ Term $(\frac{\lambda}{m} \theta_j)$ performs the regularization

$$\theta = (X^T X + \lambda I)^{-1} X^T y \quad \text{modified normal equation}$$

$$\text{where, } L = \begin{bmatrix} 0 & & & \\ 1 & 1 & & \\ & \ddots & \ddots & \\ & & 1 & 1 \end{bmatrix} \quad (n \times (n+1))$$

(Normal eq is alternative to gradient descent, iterative way)

(Note) → If $m < n$, then $X^T X$ is non-invertible (having no inverse). However, if we add the term (λI) , then $(X^T X + \lambda I)$ becomes invertible

* Regularization of Logistic Regression:

→ Logistic regression can be regularized similarly as the linear regression

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

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

Gradient Descent

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

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

* Neural Networks

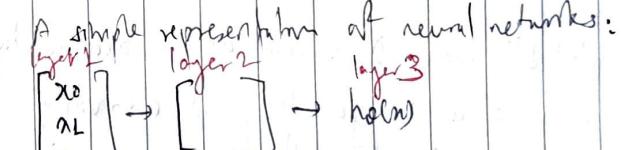
→ At a very simple level, neurons are basically computational units that take inputs (through dendrites) as electrical inputs (called "spikes") that are channeled to output axons.

→ In our model, our dendrites are like the input features x_1, x_2, \dots, x_n and the output is the result of our hypothesis function.

→ In neural networks, one input node is sometimes called the biased unit ($\theta_0 = 1$, always)

→ In neural networks, we use the same logistic function as in classification $(\frac{1}{1 + e^{-\theta^T x}})$, sometimes, we call it a sigmoid (logistic) activation function.

(These "theta" are sometimes called "weights".)



Input Node
Hidden layer
Output Node
Output layer

→ Intermediate / hidden layer nodes are also called activation units.

a_i = "activation" of unit i in layer j

$\theta_{ij}^{(l)}$ = matrix of weights controlling my function mapping from layer j to layer $j+1$

Suppose, a one hidden layer neural network represented

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ a_{1(1)} \\ a_{2(1)} \\ \vdots \\ a_{n(1)} \end{bmatrix} \rightarrow h\phi(x)$$

$a_1^{(2)}, a_2^{(2)}, a_3^{(2)}$ → activation units for layer 2.
Link unit

$$a_i^{(1)} = g\left(\pi_{10} \theta_{10}^{(1)} + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3\right)$$

$$a_3^{(2)} = g \left(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3 \right)$$

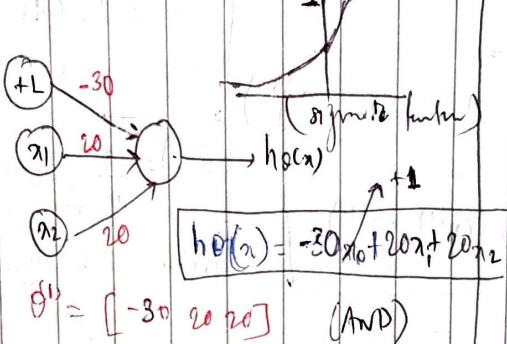
$$h_0(x) = q_1^{(3)} = g\left(\theta_{1,0}^{(2)}q_0^{(2)} + \theta_{1,1}^{(2)}q_1^{(2)} + \theta_{1,2}^{(2)}q_2^{(2)} + \theta_{3}^{(2)}q_3^{(2)}\right)$$

Here, $\theta^{(1)} = (3 \times 4)$ matrix $\theta^{(2)} = (1 \times u)$ matrix \rightarrow each layer gets its own matrix of weight b, $\theta^{(j)}$

Note: If a network has s_j units in layer j and s_{j+1} units in layer $(j+1)$, then

* (if 1 decision comes from the addition of binary unit) O(j) will be of dimension $S_{j+1} \times (S_j + 1)$

(*) logical AND of x_1, x_2 .



x_1	x_2	$h_0(n)$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(10) \approx 1$

$$\Theta^1 = [10 \ -20 \ -20] \text{ (NUR)}$$

We can combine these to get the XNOR (cycle operator (which gives 1 if all end up with 0 or both L))

→ for the transition between the first & the second layers
we'll use a matrix that combines the values for AWS & MRE

$$\theta^1 = \begin{bmatrix} -3 & 0 & 20 & 20 \\ 10 & -20 & -70 \end{bmatrix}$$

4. For 2nd ~~layer~~ the 3rd layer, we use,

$$= \begin{bmatrix} -10 & 20 & 20 \end{bmatrix}$$

$$K_{\theta, n} g^{(2)} = g \left(\theta^1 n \right)$$

$$a(s) = g(\theta, a^c)$$

The diagram shows three neural networks:

- AND Function:** Input nodes \$x_1\$ and \$x_2\$ both feed into a single output node labeled \$h_{\text{AND}}(x)\$. The connection from \$x_1\$ has weight 1.0, and the connection from \$x_2\$ has weight 0.5.
- NOT Function:** Input node \$x_1\$ feeds into a single output node labeled \$h_{\text{NOT}}(x)\$. The connection has weight -1.0.
- OR Function:** Input nodes \$x_1\$ and \$x_2\$ both feed into a single output node labeled \$h_{\text{OR}}(x)\$. The connection from \$x_1\$ has weight 1.0, and the connection from \$x_2\$ has weight 0.5.

x_1	x_2	$a_1^{(2)}$	$a_2^{(2)}$	$h_0(n)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
.	.	1	?	1

* Cost function for Neural Networks

Here,

L = total no. of layers in the network

s_L = no. of units (except bias unit) in layer L

K = no. of output units/classes

$h_{\theta}(x^{(i)})_k$ = hypothesis that results in the k th output

→ In Neural Networks, we may have multiple outputs.
(multiple h_{θ})

→ We know, for logistic regression (regularized), the cost function is as:

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

for Neural Networks;

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

→ a few nested summations have been added to account for multiple output nodes

→ In the first part of the equation, before the square brackets, we have an additional nested summation that loops through the number of output nodes.

→ In the regularization part, after the sq. brackets, we must account for multiple theta matrices, the number of columns in our current theta matrix is equal to the number of nodes in our current layer (including the bias unit)

→ The no. of rows in our current theta matrix is equal to the number of nodes in the next layer (excluding the bias unit).

Note

- The double sum simply adds up the logistic regression costs calculated for each cell in the output layer.
- The triple sum simply adds up the squares of all the individual θ 's in the entire network
- The i in the triple sum does not refer to training example i .

* Back-propagation Algorithm

→ Back propagation is the neural network terminology for minimizing our cost function, just like what we were doing with gradient descent in logistic and linear regression. Our goal is to compute: $\min J(\theta)$

→ Concretely, we want to minimize our cost function $J(\theta)$ using an optimal set of J parameters in Theta.

$$\frac{\partial J(\theta)}{\partial \theta_{j,k}}$$

• Steps:

- We are given: Training set of $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$
- set $\Delta_{ij}^{(l)} = 0$ for all (l, i, j) ($\Delta_{ij}^{(l)}$ becomes matrix full of zero)

for • training (example $t=1$ to m):

- ① Set $a^{(1)} := x^{(t)}$
- ② Perform feed-forward propagation to compute $a^{(l)}$ for $l=2, 3, \dots, L$
- ③ Using $y^{(t)}$, compute $\delta^{(L)} = a^{(L)} - y^{(t)}$
Now, L = total no. of layers

$a^{(l)}$ = vector of outputs of the activation units for the l th layer

$$y^{(t)} = \text{correct outputs in } y$$

→ an error vector is simply the difference between our actual result from the last layer and the correct outputs in y .

- (4) Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$
 using $\delta^{(l)} = ((\theta^{(l)})^T \delta^{(l+1)} \cdot a^{(l)}) \cdot (1 - a^{(l)})$
- $$g'(z) = a^{(l)} \cdot (1 - a^{(l)})$$

$$(5) \Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l+1)} \delta_i^{(l+1)}$$

In normalized form

$$\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + \delta^{(l+1)} (a^{(l)})^T$$

$$(6) D_{i,j}^{(l)} := \frac{1}{m} (\Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)}), \text{ if } (j \neq 0)$$

$$D_{i,j}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad \text{if } j=0$$

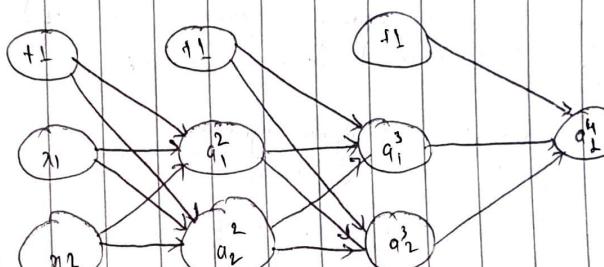
Next, capital-delta matrix D is used as an "accumulator" to add up values as we go along and eventually complete our partial derivative. Thus:

$$\boxed{\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}} = D_{i,j}^{(l)}}$$

→ Intuitively, $\delta_j^{(l)}$ is the "error" for $a_j^{(l)}$ (unit j in layer l). More formally, the delta values are actually the derivative of the cost function.

$$\delta_j^{(l)} = \frac{\partial \text{cost}(t)}{\partial z_j^{(l)}}, \text{ i.e., } \text{cost}(t) = y(t) \log(h(t)) + (1 - y(t)) \log(1 - h(t))$$

→ Recall that the derivative is the slope of a line tangent to the cost function, so the steeper the slope the more incorrect we are.



Note:

$\delta_j^{(l)}$ = "error" of cost for $a_j^{(l)}$ (unit j in layer l)

→ Back-propagation is similar to forward propagation but from right to left.

$$\delta_1^{(4)} = a_1^{(4)} - y$$

$$\text{Now, } \delta_1^3 = \theta_{11}^3 \cdot \delta_1^{(4)} \text{ and } \delta_2^3 = \theta_{21}^3 \cdot \delta_1^{(4)}$$

$$\delta_1^2 = \theta_{11}^2 \delta_1^3 + \theta_{21}^2 \delta_2^3 \quad \text{and} \quad \delta_2^2 = \theta_{12}^2 \delta_1^3 + \theta_{22}^2 \delta_2^3$$

* Gradient Checking:

→ Gradient Checking ensures that our back propagation works as intended.

→ We can approximate the derivative of our cost function with:

$$\frac{\partial J(\theta)}{\partial \theta_j} \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

Here $\epsilon \rightarrow$ Epsilon, a small real number $\approx 10^{-4}$ - guarantees that our math works

If ϵ is too small, we can end up with numerical problems.

→ with multiple theta parameters we can write:

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{J(\theta_1, \dots, \theta_j + \epsilon, \dots, \theta_n) - J(\theta_1, \dots, \theta_j - \epsilon, \dots, \theta_n)}{2\epsilon}$$

→ Run the gradient checking only till you have verified the back propagation, since the code to compute gradApprox. is very slow.

* Zero Initialization:

→ Initializing all parameters (theta) as 0 at first.

→ This causes an error called symmetry as theta's for a single activation unit likes to previous units/nodes input nodes, will end up taking the same values and this causes redundancy and not go good representation of a problem/model. (also all nodes will take up the same value repeatedly)

* Random Initialization (symmetry breaking)

→ Initialize each $\theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$

$$\text{re } [-\epsilon \leq \theta_{ij}^{(l)} \leq \epsilon]$$

$$\text{Ex: } \theta_{init} = \text{randn}(10, 11) * (2 * \text{INIT_EPSILON}) - \text{INIT_EPSILON}$$

* Putting it Together (Conclusion)

→ first, pick a network architecture, choose the layout of your neural network including how many hidden units in each layer and how many layers in total you want to have.

- No. of input units = dimension of features $x^{(i)}$
- No. of output units = no. of classes
- No. of hidden units per layer = usually more the better (must balance the cost of computation as it increases with more hidden units)
- Defaults: 1 hidden layer. (if you have more than 1 hidden layer, then it is recommended that you have the same no. of units in every hidden layer).

→ Training a Neural Network:

- randomly initialize the weights
- implement forward propagation to get $h_\theta(x^{(i)})$ for any $x^{(i)}$
- implement the cost function
- implement the back propagation to compute partial derivatives
- use gradient checking to confirm that your back propagation works. Then disable gradient checking.
- use gradient descent or a built-in optimization function to minimize the cost function with the weights in theta.

for $i = m$;
perform forward propagation & back propagation using example $(x^{(i)}, y^{(i)})$
(got activations all) and delta terms of $\delta^{(l)}$ for $l=2, \dots, L$

When we perform forward & backward propagation, loop on every training example.

Note: ideally we want, $h_\theta(x^{(i)}) \approx y^{(i)}$. This will minimize our cost function. However, we should keep in mind that $J(\theta)$ is not convex and thus we can end up in a local minimum instead.

* Evaluating the hypothesis:

- What should we do to solve the unacceptable large errors???
 - get more training examples - fixes high variance
 - try small sets of features - fixes high variance
 - try getting additional features - fixes high bias
 - try adding polynomial features - fixes high bias
 - try decreasing λ - fixes high bias
 - try increasing λ - fixes high ~~variance~~ variance
- 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.
- Diagnostic takes time to implement, but doing so saves and makes good use of your time.
- Once we have trouble-shoot the errors in our hypothesis, we now evaluate it. A hypothesis may be over-fitted and thus to evaluate a hypothesis, given a dataset of training examples, we can split the data into two sets: a training set and a test set. (Typically, training set: 70%, test set: 30%).

* Model Selection & Train / Validation / Test sets

→ The error of your hypothesis as measured on the data set with which you trained the ~~model~~ ~~your hypothesis~~ parameters will be lower than the error on any other data set.

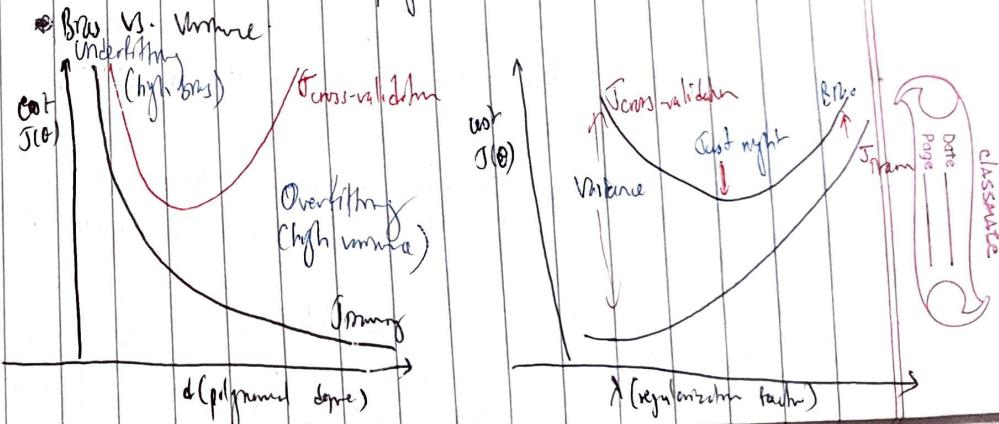
→ Given many models with different polynomial degrees, we can use a systematic approach to identify the best function. In order to choose the model of your hypothesis, you can test each degree of polynomial and look at their error result.

→ One way to break down is into the three sets:

- Training set 60%
- Cross-validation set 20%
- Test set 20%

→ Now → optimize parameters in a way the training set

- find the polynomial degree d with the least error using the cross-validation set
- estimate the generalization error using the test set with $J_{\text{test}}(\theta^{(d)})$,
 Here ($d = \text{degree of polynomial with best terms}$)

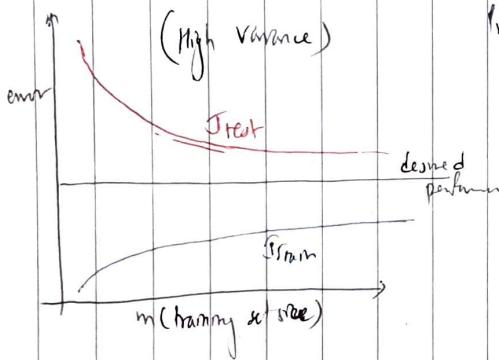
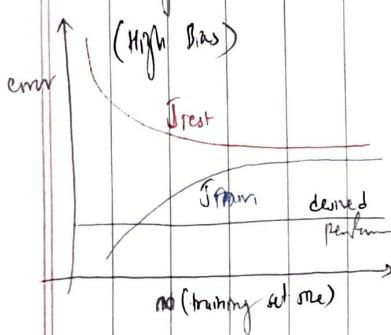


* Learning Curves

→ Training an algorithm on a very few no. of data points (2, 3, 4) will easily have 0 errors because we can always find a quadratic curve that touches exactly those number of points. Hence:

- as training sets get larger, the error for the a quadratic function ~~there~~ increases.

- the error value will plateau out after a certain m, or training set size.



* Error Analysis

→ 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.

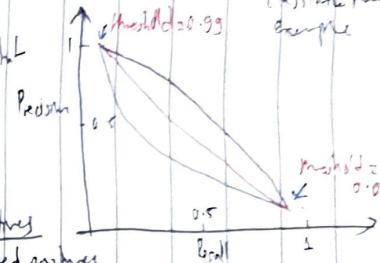
→ Error analysis: manually examine the examples (in cross-validation set) that your algorithm made errors on. See if you spot any systematic trend in what type of examples it is making errors on.

*(Reduction in recall)
↳ True positives / Total positives
↳ True negatives / Total negatives*

→ Skewed Classes: have very large no. of one type of data than other. Ex: Cancer (\rightarrow 1's) vs non-cancer (\rightarrow 0's)

function $y = \text{predict}(x)$
 $y = 0;$ // ignore $x!$
return

outputs 1 no matter what



* Precision: $\frac{\text{True positives}}{\# \text{predicted positives}}$

$$= \frac{\text{True positives}}{\text{True positives} + \text{False positives}}$$

* Recall: $\frac{\text{True positives}}{\# \text{actual positives}}$

$$= \frac{\text{True positives}}{\text{True positives} + \text{False negatives}}$$

* Trading off Precision & Recall:

→ Predict 1 if $h(x) \geq \text{threshold}$

logistic regression: $0 \leq h(x) \leq 1$

Ex: Suppose we want to predict $y=1$ (cancer) only if very confident. (high precision, low recall) $[h(x) \geq 0.9]$

→ Suppose we want to predict $y=0$ (no cancer) avoid missing too many cases of cancer (avoid false negatives). (high recall, low precision) $[h(x) \geq 0.3]$

→ Measuring and deciding on values regarding both precision & recall is non-trivial and difficult. So, F1 score is used

→ F1 Score (combines precision/recall into 1)

$$\text{F1 Score} = \frac{2 \cdot \text{P} \cdot \text{R}}{\text{P} + \text{R}}$$