

Machine Learning - Andrew Ng - Stanford

- Arthur Samuel: The field of study that gives computers the ability to learn without being explicitly programmed to do so.
- Tom Mitchell: "A computer program is said to learn from experience 'E' with respect of some class 'T' and performance measure 'P', if its performance at tasks in 'T', as measured by 'P', improves with experience 'E'"

Example playing checkers:

E = the experience of playing many times checkers 10K+

T = the task of playing checkers

P = the probability that the program will win the next game

In general, any machine learning can be assigned to one of two broad classifications:

- Supervised learning: we give answers for examples
- Unsupervised learning

Supervised learning

In supervised learning, we are given a dataset and already know what our correct output should look like, having the idea that there is a ~~non-linear~~ relationship between the input and the output.

Supervised learning problems are categorized into "regression" and "classification" problems. In a regression problem, we are trying to predict results with a continuous output, meaning that we are trying to map input values into a continuous function. In classification problem, we are trying to predict results in a discrete output.

Example 1: Given the size of houses on the real estate market, try to predict their price.
The price as a function of size is a continuous output (numerical value).

We can turn the example into a classification problem, if we try to predict if a house will sell for more or less than the asking price. We are getting two discrete (label) values, 'yes' or 'not'.

Example 2: Given a picture of a person, we have to predict

② Regression: their age - on the basis of the given picture

③ Classification: ④ Given a patient with a tumor - we have to predict whether the tumor is malignant or benign.

Unsupervised learning:

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

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

Example: Taking a group of 1,000,000 genes genes, and clustering. Find a way to automatically group these genes into groups that are somehow similar or related to different variables such as lifespan, location, role, etc.

Example: Non Clustering; The 'cocktail party algorithm' allows us to find structure in a chaotic environment, for example, identifying individual voices and music from a mesh of sounds at a cocktail party.

Linear Regression

Through the course superscripts will be used to denote the i^{th} input variable.

We'll use $X^{(i)}$ to denote the input variables and $y^{(i)}$ the output or target variable that we are trying to predict; a pair $(X^{(i)}, y^{(i)})$ is called a training example and 'm' is used to denote the training set.

h = hypothesis

Training set

$$h = \theta_0 + \theta_1 x$$

Learning algorithm

$\theta_{0,1}$ = Parameters

$X \rightarrow h \rightarrow$ predicted y

If the target variable is continuous, we call this learning problem a regression problem. Classification is when we want to predict discrete values.

We can measure the accuracy of our hypothesis function by using a cost function. This takes an average difference (actually, a fancier version of average) of all the results of the hypothesis with inputs from X 's and the actual output y 's.

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (y_i - h(x_i))^2$$

In cost function, not in gradient descent
to break it apart, it is $\frac{1}{2} \bar{x}$ where \bar{x} is the mean of the squares of $h(x_i) - y_i$, or the difference between the predicted value and the actual value.

This function is otherwise called the "Square error function" or "mean squared error". The mean is halved ($1/2$) as a convenience for the computation of the gradient ~~and~~ descent, as the derivative term of the square function will cancel out the $\frac{1}{2}$ term.

Cost Function - Intuition I

If we try to think of it in visual terms, our training data set is scattered on the x - t plane. So we are trying to make a straight line (defined by $h_\theta(x)$) which passes through those scattered data points.

The objective is to get the best possible line. The best line will be such that the average squared vertical distances of the scattered points from the line will be the least.

Gradient Descent - to minimize the cost function J
repeat until convergence {

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

}

learning rate

correct simultaneous update

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

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

$$\theta_0 := \text{temp1}$$

$$\theta_1 := \text{temp2}$$

The learning rate can be set big or small, having small 'baby steps' will take longer but will avoid overshooting, thus that happens if α is too big.

The intuition behind the convergence is that $\frac{d}{d\theta} J(\theta_n)$ will approach to 0 as we approach the bottom of our convex function. Once at the minimum, the derivative will be always 0 (zero), thus we get

$$\theta_i := \theta_i - \alpha \cdot 0 \quad \text{which is}$$

$$\theta_i := \theta_i - 0 \therefore \theta_i = \theta_i$$

Gradient Descent for Linear Regression

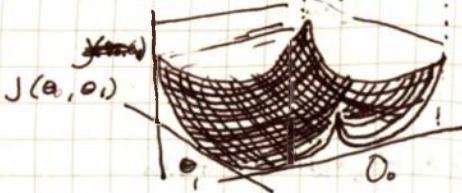
repeat until convergence ↗

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{y}_0(x^{(i)}) - y^{(i)}) \quad \checkmark$$

$$\theta_i := \theta_i - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{y}_0(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \quad \checkmark$$

↗ \Rightarrow Simultaneous

Cost function for linear regression will be always a convex function (bowl shape)



Multivariate Linear Regression

Notation:

n = number of features m = sample size

$x^{(i)}$ = input (features) of "ith" training example : $x^{(i)}$
- not exponent but superscript

$x_j^{(i)}$ = value of feature "j" in "ith" training example

Size (Foot ²)	Number of bedrooms	Number of floors	Age of Home (Years)	Price (k\$) "Y"
2104	5	1	45	460
1416	3	2	40	222
1524	3	2	30	715
852	2	1	36	178
...

$$x^{(i)} = \begin{bmatrix} 1 & x_1^{(i)} & x_2^{(i)} & \dots & x_n^{(i)} \end{bmatrix}$$

m
sample

The regression formula for hypothesis now becomes

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

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

For convenience of notation, we define
 x_0 as 1
 $\theta_0 = 1$

↓

$$\theta^T = \underbrace{\begin{bmatrix} \theta_0, \theta_1, \theta_2, \dots, \theta_n \end{bmatrix}}_{(n+1) \times 1 \text{ matrix}} \quad \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

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

For gradient descent with more than one feature, the formula changes.

repeat until convergence

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

(simultaneously update θ_j for $j=0 \dots j=n$)

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

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

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

⋮

Feature normalization / Scaling

Put all the variables between reasonable small limits

Mean normalization (don't apply it to the bias)

examples

$$x_i := \frac{\text{value} - \text{average}}{\text{max} - \text{min}}$$

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

$$x_2 = \frac{\text{bedrooms} - 2}{5}$$

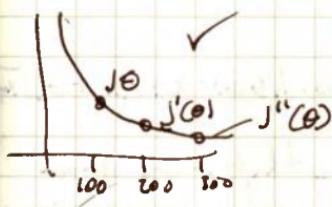
$$x_i := \frac{x_i - \mu_i}{\sigma_i}$$

~~mean~~
max - min or standard deviation.

Choose the right learning rate

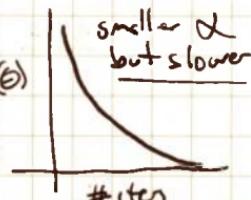
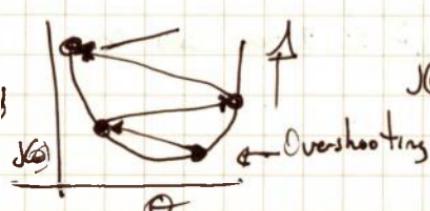
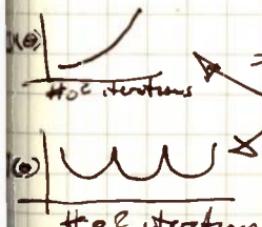
The best approach to test and choose the right learning rate, is to debug it.

By plotting the cost function across each iteration, if the cost function $J(\theta)$ ever increases, then the learning rate should be decreased.



A way to declare convergence of $J(\theta)$ decreases less than 'E' in one iteration, where 'E' is some small value such as 10^{-3} (0.001). However is difficult in practice to choose this threshold value.

It has been mathematically proved that if learning rate α is sufficiently small, then $J(\theta)$ will decrease on every iteration.



Features and Polynomial Regression

We can improve our features and the form of our hypothesis function in a couple different ways.

We can combine different features into one, for example, we can combine x_1 and x_2 into a new feature x_3 by taking $x_1 \cdot x_2$.

- Polynomial regression

~~Our hypothesis function needs to be linear (straight line)~~

Our hypothesis function need not to be linear (straight line), if that does not fit the data well.

We can change the behavior of the curve of our hypothesis function by making it a quadratic, cubic or square ~~function~~ root function (or any other form)

For example, if our hypothesis function is $h_\theta(x) = \theta_0 + \theta_1 x_1$, then we can create additional features based on x_1 to get the quadratic function $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$ or the cubic function $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$

In the cubic version we've made new features x_2 and x_3 where $x_2 = x_1^2$ and $x_3 = x_1^3$

To make a square root function, we could do: $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_1}$

One important thing to keep in mind is that if we choose our features this way, then feature scaling becomes crucial.

e.g. if x_1 has a range of $1 \leq x_1 \leq 1000$, then

$$1 \leq x_1^2 \leq 1000 \text{ 000 and}$$

$$1 \leq x_1^3 \leq 1000 \text{ 000 000}$$

Normal equation

Gradient descent gives a way of minimizing J . But there is a second way of doing so, this time performing the minimization explicitly and without resorting to an iterative algorithm. In the normal equation method, we will minimize J by explicitly taking its derivatives with respect to the θ_j 's, and setting them to zero. This allows us to find the optimum theta without iteration.

The normal equation formula is given below:

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

There is no need to do feature scaling with the normal equation.

Comparison between gradient descent and the normal equation.

Gradient Descent | Normal Equation

(-) Need to choose α
Needs many iterations

(+) No need to choose α
No need to iterate

(+) $O(n^2)$

Works well when 'n' is large

$O(n^3)$, need to calculate inverse of $X^T X$
(-) Slow if 'n' is very large
 \hookrightarrow cube of the dimension of the matrix

In practice, a larger set of features will imply problems in current computing; if $n < 10000$ still is viable, but if $n > 10000$, then iterative process with gradient descent.

This do not work on logistic regression or classification.

Normal Equation Non-invertibility

what if $X^T X$ is non-invertible? what to do if happens?

- Redundant features:

$$x_1 = \text{size in feet}^2$$

$$x_2 = \text{size in m}^2$$

- Too many features ($m \leq n$)

↳ not a good idea

↳ delete some features, or use regularization

Week 3 ~ Logistic Regression

Logistic regression is a method for classifying data into discrete outcomes. For example, to classify transactions as valid or fraud, email as spam or not, etc.

In this module we will introduce the notion of classification, the cost function for logistic regression and the application of logistic regression to multi-class classification.

We are also covering regularization. Machine learning models need to generalize well to new examples that the model has not seen in practice. This will be used to prevent overfitting the training data.

Classification

Linear regression is not a suggested algorithm, it fails to encompass all the features properly and ~~take~~ miss out, because it is not a problem for a linear function.

A classification problem is similar to a regression problem but the values we want now to predict are a small number of discrete values. For now we will focus on binary classification which can take only 2 values, 0 and 1.

$$y \in \{0, 1\}$$

0 is usually called the negative class and denotes absence of something.
1 is usually called the positive class and denotes presence of something.

They may be also denoted by symbols '+' and '-'. Given the $x^{(i)}$, the corresponding $y^{(i)}$ is also called the label for the training example.

Hypothesis representation

Logistic Regression model: Want $0 \leq h_{\theta}(x) \leq 1$

$$h_{\theta}(x) = g(\theta^T x) \quad \left. \begin{array}{l} \\ \downarrow \\ g(z) = \frac{1}{1+e^{-z}} \end{array} \right\} h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$$



$h_{\theta}(x)$ = estimated prob. of $y=1$ given input x $P(y=1|x)$

$$x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{feature size} \end{bmatrix} \quad h_{\theta}(x) = 0.7 \approx 70\% \text{ chance}$$

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

As $y \in \{0, 1\}$

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

$$\therefore P(y=0|x; \theta) = 1 - P(y=1|x; \theta) \quad \text{just a subtraction}$$

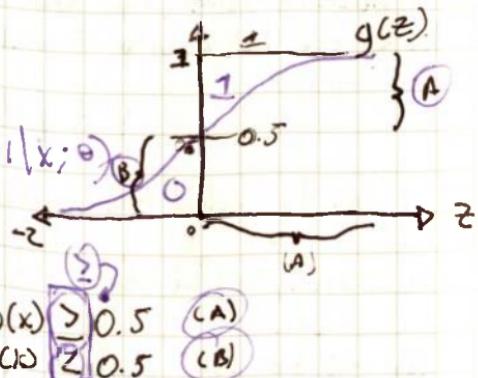
Decision Boundary

We've got the formulas

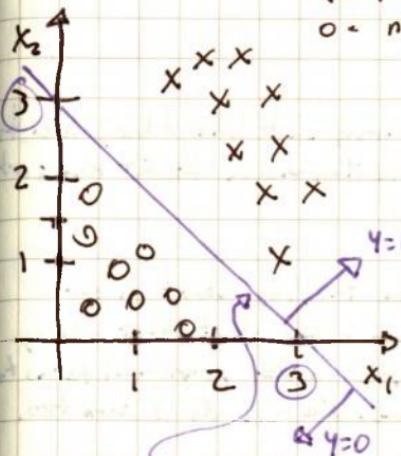
$$h_{\theta}(x) = g(\theta^T x) = P(y=1|x; \theta)$$

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

suppose predict $y=1$ if $h_{\theta}(x) \geq 0.5$
predict $y=0$ if $h_{\theta}(x) < 0.5$



x = positive example
 o = negative example



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

test " " "

$-3 \quad 1 \quad 1$

predict " $y=1$ ", if $\theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$$\Theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \quad \Theta^T = [-3 \ 1 \ 1]$$

$y=1$ if $\theta^T \geq 0$

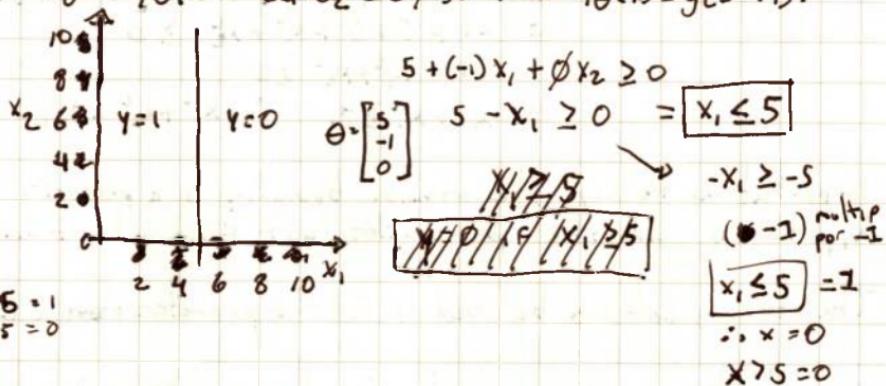
$$-3 + x_1 + x_2 \geq 0 \Leftrightarrow$$

$$x_1 + x_2 \geq 3$$

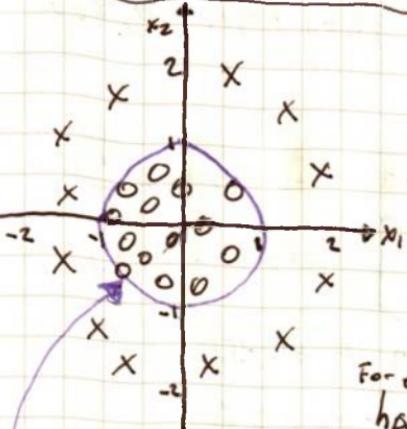
$$x_1 + x_2 = 3$$

Quiz:

Consider logistic regression with two features x_1 and x_2 . Suppose $\theta_0 = 5$, $\theta_1 = -1$ and $\theta_2 = 0$, so that $h_{\theta}(x) = g(5 - x_1)$.



Non linear decision boundaries



Given a training set like this, how to get Logistic Regression to fit the sort of data?

Earlier while talking about linear regression we talked about about adding extra higher order polynomial terms to the features. We can do it here too.

For example

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

For example, in the figure above, if we use the same test/example hypothesis: $h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$

sample values $\Rightarrow \begin{matrix} "1" \\ "1" \\ "1" \\ "1" \\ "0" \\ "0" \\ "0" \\ "0" \end{matrix}$

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

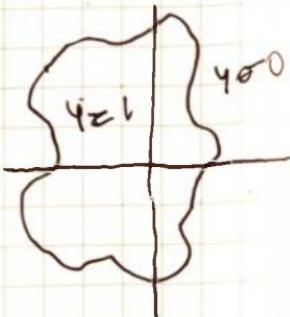
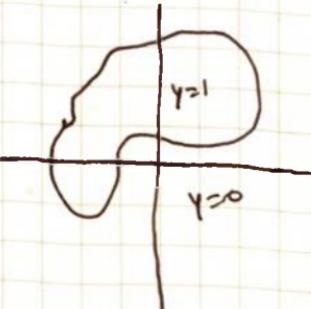
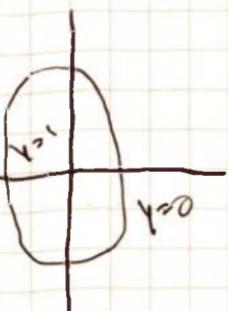
Predict "y=1" if $-1 + 0x_1 + 0x_2 + 1x_1^2 + 1x_2^2 \geq 0$
 $= -1 + 0 + 0 + x_1^2 + x_2^2 \geq 0$
 $= -1 + x_1^2 + x_2^2 \geq 0$
 $= x_1^2 + x_2^2 \geq 1 \quad \therefore \boxed{x_1^2 + x_2^2 = 1}$

Decision Boundary

Important to say, the decision boundary is a property of the hypothesis and the parameters, not from the training set.

The training set can be used to fit the parameters. i.e.,

given that, even more complex formulas may be used, which would have a representation as an ellipse or another shape on the hyperplane. i.e.:



Logistic Regression : Cost Function

Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(m)}, y^{(m)})\}$

(m) examples

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} \quad \boxed{x_0 = 1}, y \in \{0, 1\}$$

Vector
constant
 $x_0 = 1$ always

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

How to choose parameters for θ ?

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

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

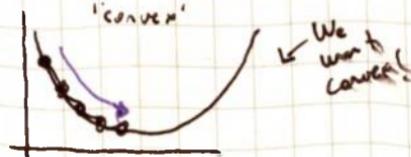
\downarrow
cost $(h_{\theta}(x^{(i)}), y^{(i)})$

$$\rightarrow \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad \downarrow \text{to make it simpler}$$

$$\text{Cost function}\bracket{\text{adopted but}} \text{Cost}(h_{\theta}(x), y) = \frac{1}{2} (h_{\theta}(x) - y)^2 \quad \frac{1}{1 + e^{-\theta^T x}}$$

NOT OPTIMAL

This cost function works fine for linear regression, but here we are trying to do logistic regression, to find labels, not predicted values. So if we could minimize the function ~~for~~ $J(\theta)$, it would work ok, but if we use this function, it would be a 'non' convex function $J(\theta)$.

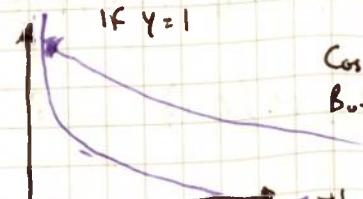


Non convex: it can fall in local minima but not guaranteed on global minima

Convex: the bowl shaped, it's guaranteed to find global minima

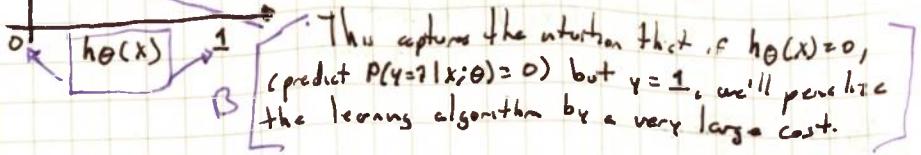
Logistic regression cost function

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



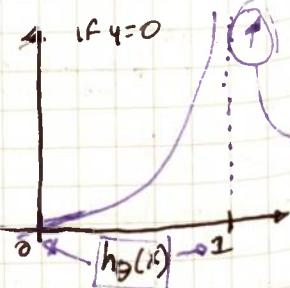
Cost is ϕ if $y=1$ and $h_\theta(x)=1$ (A)

But as $h_\theta(x) \rightarrow \phi$] because of the logarithmic nature
Cost $\rightarrow \infty$



(A) As $h_\theta(x) \geq 0.5 = 1$ (due to sigmoid), and the expected value matches exactly (1), then the cost is ϕ ; therefore the text in "B" explains why would be impossible to say that if $y=1$, that $h_\theta(x)$ would be ϕ or ≤ 0.5 on sigmoid

So that's why the cost approaches infinity, if $y=1$ but $h_\theta(x) \geq 0$,
because it's wrong, very very wrong.



What this curve does, it goes up, and goes to plus infinity as $x \rightarrow \infty$, and because as seen, that if $y=0$ but we predicted that $y=1$ with almost certainty, it would have a great cost.

And controversially, if $h_\theta(x)=\phi$ and $y=0$, at this point the cost function is going to be zero. The hypothesis melted.

Simplified cost function and Gradient Descent

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

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

Then it can be written as:

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

Demonstration:

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

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

Logistic Regression Cost Function :

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

This is a function or cost function coming from statistics, using the principle of maximum likelihood estimation. Which is an idea in statistics for how to efficiently find parameter's data for different models. Also it has a nice property, that it is convex.

Next step is to fit parameters to θ , so we need to minimize this cost function and then make a prediction given new x :

fit parameters θ : $\min_{\theta} J(\theta)$

To make a new prediction given new x :

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

Gradient Descent

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

We want $\min_{\theta} J(\theta)$:

repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad (\text{simultaneous update all } \theta_j)$$

}

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

∴

repeat {

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

Formula

$$\theta_j := \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (\text{simultaneous update all } \theta)$$

}

And the algorithm looks identical to linear regression!

- For linear regression, $h_{\theta}(x)$ was $\boxed{\theta^T x}$

- For logistic regression, $h_{\theta}(x)$ is $\boxed{\frac{1}{1+e^{-\theta^T x}}}$

Quiz: Suppose you are running gradient descent to fit a logistic regression model with parameter $\theta \in \mathbb{R}^{n+1}$. Which of the following is a reasonable way to make sure the learning rate α is set properly and the gradient descent is running correctly.

Answer: Plot $J(\theta) = -\frac{1}{n} \sum_{i=1}^n [y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log (1-h_{\theta}(x^{(i)}))]$

While running Logistic regression with gradient descent, we have all of these different parameter values $\theta_0, \dots, \theta_n$; we could use a for loop to update their values, but also it's possible to use vectorization to update all the parameters in one step: $\theta := \theta - \frac{\alpha}{n} X^T (g(X\theta) - \vec{y})$

Advanced Optimization

As we know, it's possible to use gradient descent.

There are other algorithms, more advanced and sophisticated ones, that if we provide them a way to compute $J(\theta)$ and $\frac{\partial J(\theta)}{\partial \theta}$, then these algorithms will compute the cost function for us.

- Gradient descent
- Conjugate gradient
- BFGS
- LBFGS

Advantages:

- no need to manually pick α
- often faster than gradient descent

Disadvantages:

- more complex

These are algorithms that we able to tune up automatically the learning rate on each iteration, and tend to converge faster.

How to use this algorithms?

Example

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

$$J(\theta) = (\theta_0 - 5)^2 + (\theta_1 - 5)^2$$

Function [jVal, gradient] =
costFunction(theta)

$$jVal = (\theta_0 - 5)^2 / 2 + \dots$$

$$+ (\theta_1 - 5)^2 / 2;$$

$$\text{gradient} = \text{zeros}(2, 1);$$

$$\text{gradient}(1) = 2 * (\theta_0 - 5);$$

$$\text{gradient}(2) = 2 * (\theta_1 - 5);$$

$$\frac{\partial}{\partial \theta_0} = 2(\theta_0 - 5)$$

$$\frac{\partial}{\partial \theta_1} = 2(\theta_1 - 5)$$

function [jVal, gradient] = costFunction(theta)

jVal = % Code to compute $J(\theta)$;

gradient = % Code to compute derivative of $J(\theta)$;

then:

```
options = optimset('GradObj', 'On', 'MaxIter', '100');
```

```
initialTheta = zeros(2, 1)
```

```
[optTheta, FunctionVal, exitFlag] = fminunc(@costFunction, initialTheta, options);
```

Multiclass Classification: One vs All classification

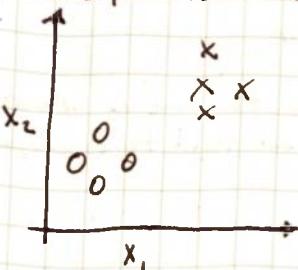
If we want to classify several options such as:

- Email folders/tags: Work, Friends, Family, Hobby

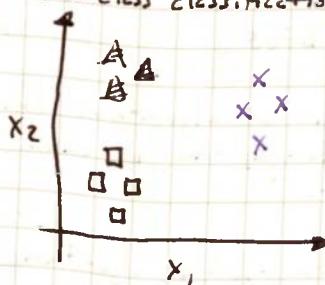
- Medical diagnosis: Not ill, Cold, Flu

- Weather: Sunny, Cloudy, Rain, Snow

Binary Classification

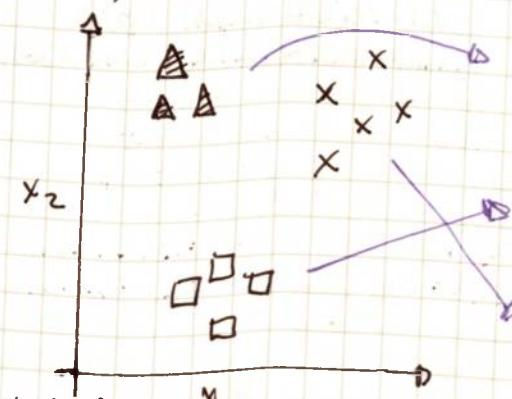


Multiclass classification



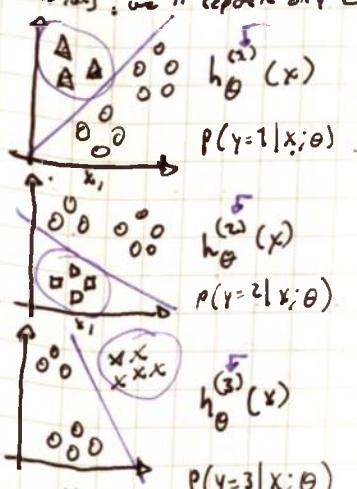
On this model, the challenge is to determine the class where an hypothesis belongs. For that, what we do is to do, in the example, is to do 3 different classifications, in other words, make 3 different binary classifications.

We create three different binary classifications, we'll separate only 2 classes, one at a time:



$$h_{\theta}^{(i)}(x) = P(y=i|x; \theta), i \in \{1, 2, 3\}$$

For a new sample, choose the class i that maximizes the one with higher probability

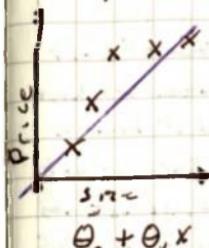


[Solving the problem of overfitting]

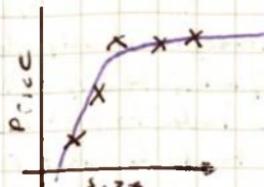
Overfitting and its problem

Overfitting describes a situation where the model learns to describe random error or noise & instead of 'real' meaningful relationships between features and the output. In other words, when it fits the training sample "too exactly, too accurate", and new samples, unseen to the model, can't and won't fit.

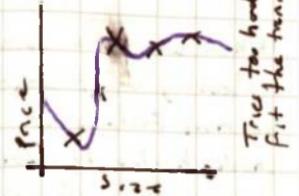
Example with **linear** regression (housing price)



"Underfitting" or
"high bias", doesn't fit
very well



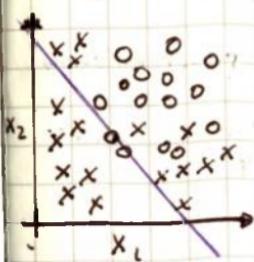
"Just right"



"Overfit", "high variance"
Trained to fit the training set

Thus, if we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \approx 0$), but fail to generalize to new examples.

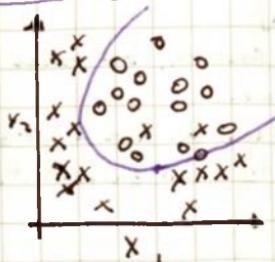
Example with **logistic** regression



$$h(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

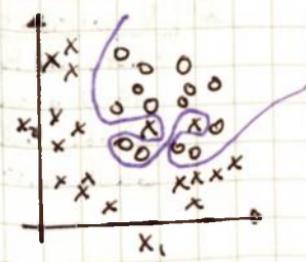
g = sigmoid

"Underfit"



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$

"Just right"



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_1^2 x_2^2 + \theta_4 x_1^2 x_2 + \dots)$$

"Overfit"

Addressing Overfitting

In practice, we will have sometimes several attributes/features

x_1 = size of house

x_2 = no. of bedrooms

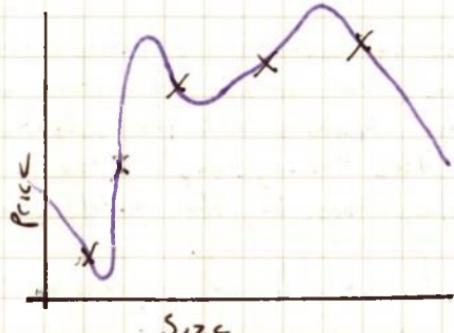
x_3 = no. of floors

x_4 = age of the house

x_5 = average income in neighborhood

x_6 = kitchen size

x_{100}



When we have several features, is not only ~~about~~ a matter of just selecting what degree polynomial to use, and in fact, when we have so many features it becomes harder to plot the data and it becomes much harder to visualize it, to decide what features to keep or not. So concretely, if we are trying to predict housing prices, sometimes we can have lots of different features, and they ~~are~~ seem kind of useful, but if we have a lot of features and very little training data, then overfitting can become a problem.

Options:

① Reduce number of features

- Manually select which features to keep
- Model selection algorithm (later in course)

feature, feature and feature contains information about the problem

② Regularization

- Keep all the features, but reduce the magnitude/values of θ parameters.
- Works well when we have a lot of features, each of which contributes a bit on predicting y

Cost function

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.

Say we wanted to make the following function more quadratic:

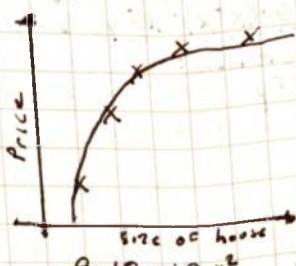
$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

We will want to reduce (almost eliminate) the influence of θ_3 and $\theta_4 x^4$ without actually getting rid of these features or changing the form of our hypothesis, then we can instead modify our cost function:

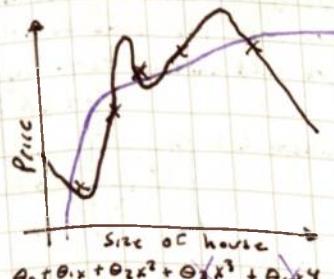
$$\min_{\theta} = \frac{1}{2m} \sum_{i=1}^m (\hat{h}_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \cdot \theta_3^2 + 1000 \cdot \theta_4^2$$

We've added two extra terms at the end to inflate the cost of θ_3 and θ_4 . Now in order to get close to zero, we will have to reduce the values of θ_3 and θ_4 near to zero. This will in turn greatly reduce the values of $\theta_3 x^3$ and $\theta_4 x^4$ in our hypothesis function. As a result, we see that the new hypothesis (depicted by the blue curve) looks like a quadratic function but fits the data better due to the extra small terms $\theta_3 x^3$ and $\theta_4 x^4$.

Intuition:



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



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

We also could regularize all of our ~~theta~~ parameters in a single summation:

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

The λ or lambda, is the regularization parameter, and determines how much the cost is inflated (with that, we can smooth the output of our hypothesis function to reduce overfitting). If λ is chosen to be too large, it may smooth too much and cause underfitting.

[Regularized Linear Regression]

We can apply regularization to both linear regression and logistic regression. On the case of linear regression:

Gradient descent

We will modify our gradient descent function to separate out θ_0 from the rest of the parameters because we don't want to penalize θ_0 .

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

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

The term $\frac{\lambda}{m} \theta_j$ performs our regularization. With some manipulation our update rule can also be represented as:

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

The first term in the above equation, $1 - \alpha \frac{\lambda}{m}$ will always be less than 1 (SI). Intuitively can be seen as reducing the value of θ_j on every update. Notice the second term is exactly the same as it was before.

Normal Equation

Using the alternate method of the non-intercepting normal equation, we use the same formula as before: $\theta = (X^T X)^{-1} X^T y$ but, now we add another term inside the parenthesis:

$$\theta = (X^T X + \lambda \cdot L)^{-1} X^T y$$

$$\text{where } L = \begin{bmatrix} 0 & & & \\ 1 & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}$$

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

$m \leq n$
samples # Features

$$(n+1) \times (n+1)$$

L is a matrix with 0 at the top left and 1's down the diagonal with 0's everywhere else. It should have the dimension $(n+1) \times (n+1)$

Intuitively, this is an identity matrix (though we are not including θ_0), multiplied with a single real number λ .

Important to recall that if $m < n$ (samples less than features), then $X^T X$ is non-invertible. However when we add the term $\lambda \cdot L$, then $X^T X + \lambda \cdot L$ becomes invertible.

Regularized Logistic Regression

Similar to regularized linear regression, now we will adapt our functions to optimize or regularize logistic regression.

$$h(\theta) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^3 x_2^3 + \dots)$$

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

..... Dotted line = unregularized
— Blue line = regularized

G_1, G_2, \dots, G_n

That's our cost function for logistic regression and we've modified it to use regularization, we add $\frac{\lambda}{2m} \sum_{j=2}^n \theta_j^2$. Which also has the effect of penalizing parameters θ_j for being too large.

Gradient descent:

repeat {

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

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

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

$$\approx \frac{d}{d\theta} J(\theta)$$

Advanced optimization

Function [JVal, gradient] = costFunction(theta)

JVal = [code to compute $J(\theta)$];

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

in diag 200, 20
and 400 = 2

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \\ \theta_5 \end{bmatrix}$$

gradient(1) = [code to compute $\frac{d}{d\theta_0} J(\theta)$];

$$\% \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \quad \text{term } \theta_0 \quad \text{intercept}$$

gradient(2) = [code to compute $\frac{d}{d\theta_1} J(\theta)$];

$$\% \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)} + \frac{\lambda}{m} \theta_1 \quad \text{theta 1}$$

gradient(3) = [code to compute $\frac{d}{d\theta_2} J(\theta)$];

$$\% \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)} + \frac{\lambda}{m} \theta_2 \quad \text{theta 2}$$

gradient(n+1) = [code to $\frac{d}{d\theta_n} J(\theta)$];

$$\% \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_n^{(i)} + \frac{\lambda}{m} \theta_n$$

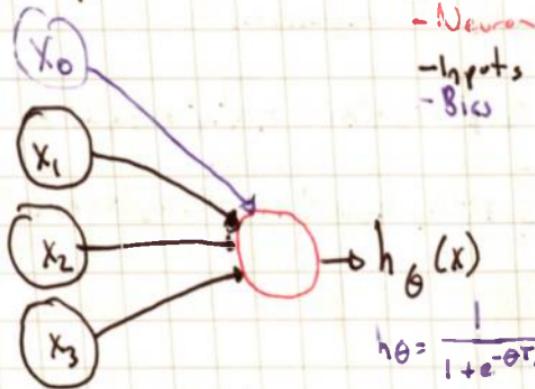
A vectorized implementation is:

$$h = g(X\theta)$$

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

(i) week 4 - Neural Networks

Model representation, Logistic unit, a neuron

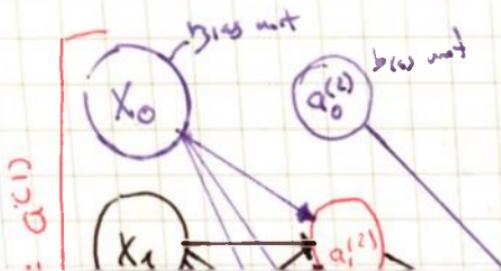


$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}, \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

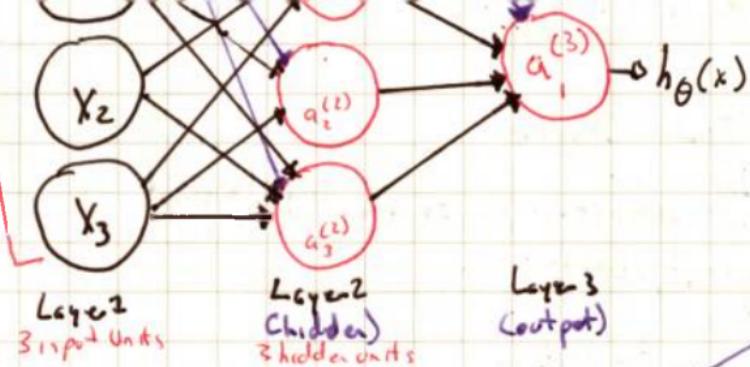
Weights or parameters

This is a neuron with a sigmoid or logistic activation function.

Neural Networks



$a_i^{(j)}$ = "activation" of unit i in layer j
 $\Theta^{(j)}$ = matrix of weights controlling function mapping from layer j to layer $j+1$



3×4 matrix

$$\begin{aligned}
 1. \quad a_1^{(1)} &= g(\theta_{1,0}^{(1)}x_0 + \theta_{1,1}^{(1)}x_1 + \theta_{1,2}^{(1)}x_2 + \theta_{1,3}^{(1)}x_3) \\
 2. \quad a_2^{(1)} &= g(\theta_{2,0}^{(1)}x_0 + \theta_{2,1}^{(1)}x_1 + \theta_{2,2}^{(1)}x_2 + \theta_{2,3}^{(1)}x_3) \\
 3. \quad a_3^{(1)} &= g(\theta_{3,0}^{(1)}x_0 + \theta_{3,1}^{(1)}x_1 + \theta_{3,2}^{(1)}x_2 + \theta_{3,3}^{(1)}x_3)
 \end{aligned}$$

$$\Theta^{(1)} = E \in \mathbb{R}^{3 \times (3+1)} = 3 \times 4$$

horizontal matrix

$$h_\theta = a_1^{(2)} = g(\theta_{1,0}^{(2)}a_0^{(1)} + \theta_{1,1}^{(2)}a_1^{(1)} + \theta_{1,2}^{(2)}a_2^{(1)} + \theta_{1,3}^{(2)}a_3^{(1)})$$

If a network has s_j units in layer j , s_{j+1} units in layer $j+1$, then $\Theta^{(j)}$ will be of dimensions:

$$S_{j+1} \times (S_j + 1)$$

units
of
output

The size of the matrix is outputs \times inputs, where "outputs"

is the number of units in the next layer and "inputs" is the number of units in previous layer including bias unit.

Model representation 2: Vectorized implementation

Previously we said that the sequence of steps we need to take in order to compute the output of a hypothesis are the equations given before on the previous page, compute the activation values of the three hidden units and use those to compute the final output of our hypothesis $h_{\theta}(x)$. To simplify the notation we'll introduce new terms:

$$z = \text{input} \quad z_i^{(2)} := a_i^{(1)} = g(z_i^{(1)})$$

$$a_0^{(2)} = g(\theta_{00}^{(1)} x_0 + \theta_{10}^{(1)} x_1 + \theta_{20}^{(1)} x_2 + \theta_{30}^{(1)} x_3)$$

$$a_1^{(2)} = g(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3) \quad z_1^{(2)}$$

$$a_2^{(2)} = g(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3) \quad z_2^{(2)}$$

$$a_3^{(2)} = g(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3) \quad z_3^{(2)} \therefore a_3^{(2)} = g(z_3^{(1)})$$

The dotted area resembles a matrix-vector operation, or $\Theta^{(1)} X$, so we can write a vectorized computation:

$\Theta \rightarrow$

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(2)} = \begin{bmatrix} z_0^{(2)} \\ z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \sim 3d \text{ vector } (3 \times 1)$$

3d vector
IR³

$$\boxed{z^{(2)} = \Theta^{(1)} X \xrightarrow{\text{activation: sigmoid}} \text{fraction - element wise} \quad \text{IR}^3 \sim 3d \text{ vector}}$$

$$a^{(2)} = g(z^{(2)}) \quad z^{(2)} = G^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)}) \quad a^{(2)} = g(z^{(2)})$$

And to account for the bias, $a_0^{(2)}$:

$$a_0^{(2)} = 1$$

$a^{(2)}$ (second layer) becomes $a^{(2)} \in \mathbb{R}^4$ (4×1), and to compute the first value, on the third layer, of our hyp.

$$z^{(3)} = \Theta^{(2)} a^{(2)} \therefore$$

$$h_{\theta}(x) = a^{(3)} = g(z^{(3)})$$

Alejandro: note that the superscripts are being assigned forward. Because of this, it's called "forward propagation"

Repeating:

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(1)} = \begin{bmatrix} z_0^{(1)} \\ z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \end{bmatrix} \rightarrow \quad z^{(2)} = \Theta^{(1)} a^{(1)} \quad z^{(2)} = \Theta^{(2)} a^{(2)}$$

$$a^{(1)} = g(z^{(1)}) \rightarrow \quad a^{(2)} = g(z^{(2)})$$

$$a^{(3)} = g(z^{(3)}) \Rightarrow h_{\theta}(x)$$

The neural network works as logistic regression, first layer (input) takes x to perform calculations, but subsequent layers uses already completed or "learned" parameters from the previous layers as input to the second one.

Examples and Intuitions (1)

Predicting the **AND** function, the graph of our function will look like

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow g(z^{(2)}) \rightarrow h_{\theta}(x)$$

$$0 \cdot 0 = 0$$

$$0 \cdot 1 = 0$$

$$1 \cdot 0 = 0$$

$$1 \cdot 1 = 1$$

x_0 is our bias value
 x_1, x_2 = 2

Assuming weights of -30 for x_0 , 20 for x_1 , and 20 for x_2 , we have:

$$\Theta^{(2)} = [-30 \ 20 \ 20]$$

This will cause the output of our hypothesis if both x_1 and x_2 are 1

$$h_{\theta}(x) = g(-30x_0 + 20x_1 + 20x_2)$$

$$h_{\theta}(x) = g(-30 + 20x_1 + 20x_2) \therefore$$

$$x_1 = -10 \quad x_2 = 0 \quad \text{then} \quad g(-30) \approx 0$$

$$x_1 = 0 \quad x_2 = 1 \quad \text{then} \quad g(-10) \approx 0$$

$$x_1 = 1 \quad x_2 = 0 \quad \text{then} \quad g(-10) \approx 0$$

$$x_1 = 1 \quad x_2 = 1 \quad \text{then} \quad g(40) \approx 1$$

OR function

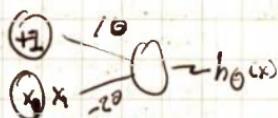
$$\Theta^{(2)} = [-10 \ 20 \ 20] = h_{\theta}(x) = (-10x_0 + 20x_1 + 20x_2)$$

OR

$-10x_0$	$(w)x_1$	$20x_2$	\therefore	
-10	$20(0)$	$20(0)$	$\therefore \approx$	$g(-10) \approx 0$
-10	$20(0)$	$20(1)$	$\therefore \approx$	$g(-10) \approx 1$
-10	$20(1)$	$20(0)$	$\therefore \approx$	$g(-10) \approx 1$
-10	$20(1)$	$20(1)$	$\therefore \approx$	$g(-10) \approx 1$

Intuitions 2

NOT

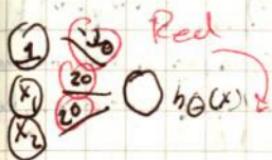


x_i	$h_{\theta}(x)$
0	$g(10-0) \approx 1$
1	$g(10-20) \approx 0$

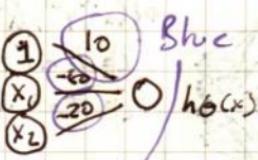
$$g(10) \rightarrow$$

$$h_{\theta}(x) = g(10 - 20x_1) \therefore$$

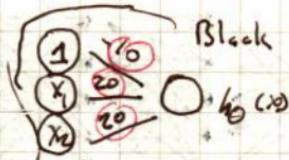
Putting them altogether: $X_1 \times \text{NOR} \times X_2$



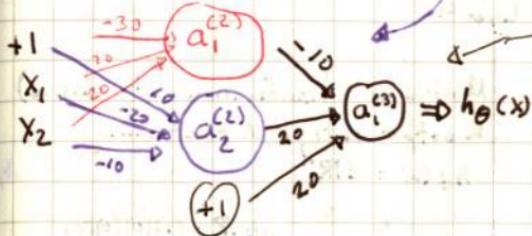
$$\rightarrow X_1 \text{ AND } X_2$$



$$\rightarrow (\text{NOT } X_1) \text{ AND } (\text{NOT } X_2)$$



$$\rightarrow X_1 \text{ OR } X_2$$



X_1	X_2	$a_1^{(l1)}$	$a_2^{(l1)}$	$h_\theta(x)$
0	0	0	0	1
0	1	0	0	0
1	0	0	0	0
1	1	1	1	1

Multiclass Classification

To classify data into multiple classes, we let our hypothesis function to return a vector of values, say, if we want to distinguish pictures from a pedestrian, a car, a motorcycle and a truck, what we would have is for example:

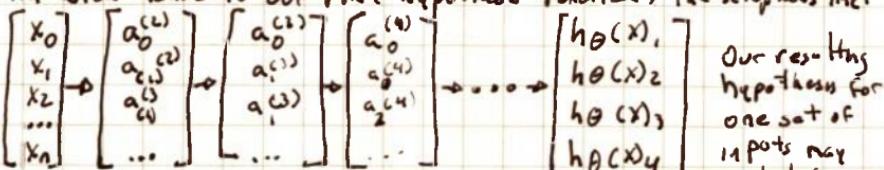
$$h_\theta(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad h_\theta(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad h_\theta(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \quad h_\theta(x) \approx \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

pedestrian car motor. truck

We can define our set of resulting classes as $\mathcal{Y} =$

$$\mathcal{Y}^{(i)} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Each $\in \mathcal{Y}^{(i)}$ represents a different image corresponding to either a car, pedestrian, truck or motorcycle. The inner layers, each provide us with some new information which leads to our final hypothesis function. The setup looks like:

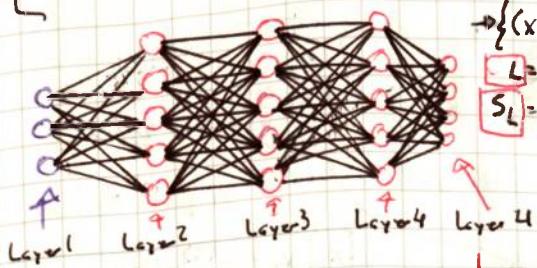


$$h_\theta(x) = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

which in our example would represent a motorcycle

Our resulting hypothesis for one set of inputs may look like

[Week 5 Neural Networks: Learning]



$$\rightarrow \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})\}$$

$L = \text{total \# of layers in the network}$

$S_L = \text{\# of units in layer } L \text{ (not counting bias)}$

$\therefore L = 4$

$S_1 = 3, S_2 = 5, S_3 = 5, S_4 = 4 \Rightarrow S_L = 4$

$S_L = \text{the final layer}$

- Binary Classification: only 2 labels, $y = 0 \text{ or } 1$
- One output unit
- $h(x) = \text{EIR}$ Area of number
- $K=1$
- $S_L = 1$
- One output but with two possible values {0,1} $K=2$

Multiclass classification (K classes)
K output units, K distinct classes
 $h_{\theta}(x) = \epsilon \mathbb{R}^K$ -> K dimensions! output.
 $S_k = K$
 $y \in \mathbb{R}^K = \begin{bmatrix} \text{---} \\ \text{---} \\ \vdots \\ \text{---} \\ \text{---} \end{bmatrix}^3$
 $K \geq 3$ because
 For 2 classes we would use
 a binary classification
 with one output unit


Cost Function:

The cost function for a neural network will be a generalization of the one used for logistic regression used to minimize the function of $J(\theta)$, with its regularization function (still without regularizing the bias).

Now instead of having one basically one logistic regression output unit, we may instead have K of them.

$$J(\theta) = -\frac{1}{n} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)})) + (1-y_k^{(i)}) \log(1-h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^m \sum_{j \neq 0} (\theta_j^{(l)})^2$$

Sum of the K output units on the final layer
for each sample

pattern θ
 $\theta \neq 0$

Looks complicated but what's the term $\theta_j^{(l)}$?

Notes:

- Double sum simply adds the logistic regression costs calculated for each cell in I_2
 - The triple sum simply adds up the scores of all the individuals in the entire network
 - The i in the triple sum does NOT refer to training examples!

Sum of the k^{th} output
units on the final layer

Looks complicated but what it's doing is summing over terms $\beta_{j1} x_j$ for all values of $j \neq 0$, except that we don't sum over the terms corresponding to this bias values like we had for logistic regression. Concretely, we don't sum over the terms corresponding to where $i = 0$, that.

because it would yield into regularizing the bisequent, and we should not.

Back propagation

We need to have code to compute:

$J(\theta)$ \rightarrow uses the formula on
the previous page

$\frac{d}{d\theta_{ij}^{(l)}} J(\theta)$ \rightarrow the partial derivative

Gradient computation

To simplify, imagine we have only one training example (x, y) , no subscripts, so (x, y) , only one per

Forward Propagation

$$\alpha^{(0)} = x$$

$$z^{(1)} = \Theta^{(0)} \alpha^{(0)}$$

$$\alpha^{(2)} = g(z^{(1)}) \text{ (add } \alpha_0^{(2)})$$

$$z^{(3)} = \Theta^{(2)} \alpha^{(2)}$$

$$\alpha^{(3)} = g(z^{(3)}) \text{ (add } \alpha_0^{(3)})$$

$$z^{(4)} = \Theta^{(3)} \alpha^{(3)}$$

$$\alpha^{(4)} = g(z^{(4)}) = h_{\Theta}(x)$$

Back Propagation

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

On the above and last page sample, $l=4$

$\Delta_j^{(4)} = \alpha_j^{(4)} - y_j$ = activation of unit - value of target

$$\Delta^{(4)} = (\Theta^{(3)})^T \Delta^{(4)} \cdot g'(z^{(3)})$$

$$\Delta^{(3)} = (\Theta^{(2)})^T \Delta^{(4)} \cdot g'(z^{(2)})$$

- No $\Delta^{(1)}$, because it's the input layer

$$\alpha^{(3)} \cdot (1 - \alpha^{(3)})$$

$$\alpha^{(4)} \cdot (1 - \alpha^{(4)})$$

Back propagation is a neural network terminology for minimizing our cost function, just as we were doing with gradient descent in logistic and linear regression. Our goal is to compute $\min_{\theta} J(\theta)$.

That is, we want to minimize our cost function J using an optimal set of parameters in there. In this section we'll look at the equations we use to compute the partial derivative of $J(\theta)$:

$$\frac{d}{d\theta_{ij}^{(l)}} J(\theta)$$

Given the training set $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

• Set $\Delta_{ij}^{(1)} := \emptyset$ for all i, j, l ; have hence a matrix full of zeroes.



- for training example $t=1$ to m : $(x^{(t)}, y^{(t)})$ training example
- 1 Set $a^{(2)} = x^{(t)}$
 - 2 Perform forward propagation to compute $a^{(l)}$ for $t=1, 2, \dots, L$ ($a^{(L)}$ is our target)
 - 3 Using $y^{(t)}$ compute $\delta^{(L)} = a^{(L)} - y^{(t)}$ with target value $y^{(t)}$
 - 4 Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$. Denote $\delta^{(L)}$ to $\delta^{(2)}$
 - 5 $\Delta_{ij}^{(t)} := \Delta_{ij}^{(t)} + a_j^{(t)} \delta_i^{(t+1)}$ ~~$\Delta^{(t)} := \Delta^{(t)} + \delta^{(t+1)} (a^{(t)})^T$~~ ~~$\Delta^{(t)} := \Delta^{(t)} + \delta^{(t+1)} (a^{(t)})^T$~~ with vectorization
 - 6 $D_{ij}^{(t)} := \frac{1}{m} \Delta_{ij}^{(t)} + \lambda \Theta_{ij}^{(t)}$ if $j \neq 0$ $\left| \frac{d}{d \Theta_{ij}^{(t)}} J(\theta) = D_{ij}^{(t)} \right.$
 - 7 $D_{ij}^{(t)} := \frac{1}{m} \Delta_{ij}^{(t)}$ if $j = 0$

③ Using $y^{(t)}$, compute $\delta^{(L)} = a^{(L)} - y^{(t)}$

(Where L is our total number of layers and $a^{(L)}$ is the vector of outputs of the activation units for the last layer. So our "error values" for the last layer are simply the differences of our actual results in the last layer and the correct results in y . To get the delta values of the layers before the last layer, we can use an equation that steps us back from right to left: ④)

④ Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$ using $\delta^{(L)} = ((\theta^{(L)})^T \delta^{(L+1)}). a^{(L-1)}$

$$\delta^{(L)} = ((\theta^{(L)})^T \delta^{(L+1)}). * a^{(L)} * (1 - a^{(L)})$$

The delta values of layer l are calculated by multiplying the delta values in the next layer with the ~~the~~ matrix of layer l .

Then we multiply 'element wise' that with a function called g' , or g prime, which is the derivative of the activation function g evaluated with the input values given by $z^{(L)}$. The g -prime derivative terms, can also be written out as:

$$g'(z^{(L)}) = a^{(L)} * (1 - a^{(L)})$$

⑤ $\Delta_{ij}^{(t)} := \Delta_{ij}^{(t)} + a_j^{(t)} \delta_i^{(t+1)}$, or with vectorization; $\Delta^{(t)} := \Delta^{(t)} + \delta^{(t+1)} (a^{(t)})^T$

Hence we update our new matrix Δ :

$$D_{ij}^{(t)} := \frac{1}{m} (\Delta_{ij}^{(t)} + \lambda \Theta_{ij}^{(t)}), \text{ if } j \neq 0$$

the capital Delta matrix
is used as a "calculator"
to add up our values as we go
along and eventually compute our
partial derivative $\frac{\partial J(\theta)}{\partial \theta_{ij}} = \delta_{ij}^{(t)}$

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

Implementation note: unrolling parameters

With neural networks, we are working with sets of matrices:

$$\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}, \dots$$

$$D^{(1)}, D^{(2)}, D^{(3)}, \dots$$

Optimization function tend to need vectors as an input, that's why we will want to 'unroll' all the elements and put them into a long vector:

$$\text{ThetaVector} = [\text{theta}_1(:); \text{theta}_2(:); \text{theta}_3(:)]$$

We use $(:)$

$$\Delta\text{thetaVector} = [D1(:); D2(:); D3(:)]$$

- If the dimension of theta1 is 10×11 , theta2 is 10×11 , and theta3 is ~~1×11~~ 1×11 , we will get a vector with length of 232 elements.
- To get back a set of matrices from a vector, we use the reshape command, so we get the unrolled version of it:

reshape (vector, newRows, newCols)

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

Theta2 = reshape(thetaVector(111:210), 10, 11)

Theta3 = reshape(thetaVector(221:232), 1, 11)

Learning algorithm:

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

function [Jval, gradientVec] = costFunction(thetaVec)

From thetaVec, get $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$

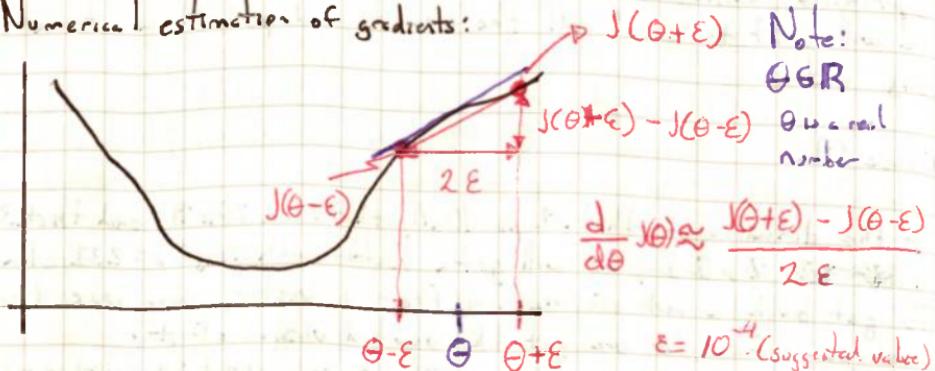
Use forward propagation and backward propagation to get $D^{(1)}, D^{(2)}, D^{(3)}$ and $J(\Theta)$

Unroll $D^{(1)}, D^{(2)}, D^{(3)}$ to get gradient Vec.

Gradient checking \sim "Debug forward and back prop"

While implementing the algorithm, it might look that the cost function indeed gets reduced, but with fairly complex algorithms, it might happen that by a small bug in the implementation, we end with a neural network that has a big error rate, and to avoid that, it's recommended to test it, like with the following procedure.

Numerical estimation of gradients:



Implement: $\text{gradApprox} = (J(\theta_{\text{true}} + \text{EPSILON}) - J(\theta_{\text{true}} - \text{EPSILON})) / 2 * \text{EPSILON}$

A way to estimate the gradient, is to choose a value of θ_{true} and place two $\theta_{\text{true}} \pm \text{EPSILON}$. So we can numerically approximate.

Then we connect those points with a straight line $J(\theta + \epsilon)$ and $J(\theta - \epsilon)$, also then draw a line perpendicular line down from $J(\theta + \epsilon)$ and horizontal line from $J(\theta - \epsilon)$, forming an isosceles triangle; then, similar to the calculation of the slope of the line, we do so, with $\frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$; Usually, a small value works fine like $\epsilon = 10^{-4}$, but choosing smaller values can cause numerical problems.

Quiz: Let $J(\theta) = \theta^3$. Furthermore, let $\theta = 1$ and $\epsilon = 0.01$, what's the approximate derivative

$$\theta = 1 \quad \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} = \frac{(1+0.01)^3 - (1-0.01)^3}{2 \times 0.01} = \frac{(1.01)^3 - (0.99)^3}{0.02}$$

$$= \frac{1.0303 - 0.9703}{0.02} = \frac{0.0600}{0.02} = 3.001$$

$J(\theta) = \theta^3$

Now, on the previous page we did consider Θ as \mathbb{R} , now we will look at a more general case where $\Theta \in \mathbb{R}^n$ (vector) (Θ is an unrolled version of $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$, the parameters of four layers), so Θ is vector with n elements, from $[\Theta_1 \dots \Theta_n] = \Theta$; we can then use a similar idea to compute the approximate partial derivatives:

$$\frac{d}{d\Theta_1} J(\Theta) \approx \frac{J(\Theta_1 + \epsilon, \Theta_2, \Theta_3, \dots, \Theta_n) - J(\Theta_1 - \epsilon, \Theta_2, \Theta_3, \dots, \Theta_n)}{2\epsilon}$$

$$\frac{d}{d\Theta_2} J(\Theta) \approx \frac{J(\Theta_1, \Theta_2 + \epsilon, \Theta_3, \dots, \Theta_n) - J(\Theta_1, \Theta_2 - \epsilon, \Theta_3, \dots, \Theta_n)}{2\epsilon}$$

⋮

$$\frac{d}{d\Theta_n} J(\Theta) \approx \frac{J(\Theta_1, \Theta_2, \Theta_3, \dots, \Theta_{n-1}, \Theta_n + \epsilon) - J(\Theta_1, \Theta_2, \Theta_3, \dots, \Theta_{n-1}, \Theta_n - \epsilon)}{2\epsilon}$$

In code: for $i=1:n$

thetaPlus = theta;

thetaPlus(i) = thetaPlus(i) + EPSILON;

thetaMinus = theta;

thetaMinus(i) = thetaMinus(i) - EPSILON;

gradApprox(i) = $(J(\text{thetaPlus}) - J(\text{thetaMinus})) / 2 * \text{EPSILON}$;

end;

$$\begin{bmatrix} \Theta_1 \\ \Theta_2 \\ \Theta_3 \\ \Theta_i + \epsilon \\ \Theta_n \end{bmatrix}$$

$\Theta - \epsilon$

And then, check that gradApprox ≈ DVec; should be pretty close.

↑ derivatives from back prop. (vector)

Implementation note:

- Implement backprop to compute DVec (unrolled $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$).
- Implement numerical gradient check to compute gradApprox
- Make sure they give similar values
- Turn off gradient checking and use back propagation for learning.

Important: Be sure to double gradient checking code before training the classifier. If numerical gradient computation is run on every iteration of gradient descent (or in inner loop of costFunction), the code will be very slow.

[Random Initialization]

Initiating all the weights to zero does not work with neural networks. When we back propagate, all nodes will update to the same value ~~except~~ repeatedly. Instead we can randomly initialize our weights for our Θ matrices using random values.

The random values should be symmetrically distributed, meaning that the range of possible values have the same 'distribution'; to achieve that we expect to have initialized:

$$\Theta_{ij}^{(l)} \text{ to a random value in } [-\epsilon, \epsilon] \\ \text{ s.t. } -\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$$

Note that ϵ here
doesn't relate to the
previous value used on
numerical reduction for
gradient checking.

In code, for example:

$$\text{theta1} = \text{rand}(10, 11) * (2 * \text{EPSILON}) - \text{EPSILON}$$

$$\text{theta2} = \text{rand}(1, 11) * (2 * \text{EPSILON}) - \text{EPSILON}$$

Putting it together: Pick the network architecture

- Number of input units: # of dimensions of features $\times^{(1)}$
- Number of output units: # of classes
- Number of hidden units per layer: Usually the more the better, must balance with computation costs.

→ Defaults: 1 hidden layer; if we have more than 3, then it's recommended to have the same number of units in every hidden layer.

Training a NN:

- 1: Randomly initialize the weights
- 2: Implement forward propagation to get $h_\Theta(x^{(i)})$ for any $x^{(i)}$
- 3: Implement the cost function
- 4: Implement back propagation to compute partial derivatives
- 5: Use gradient checking to confirm that your backprop works, then disable it
- 6: Use gradient descent or a built-in optimization function to minimize the cost function with the weights in theta.

For $i=1:m$

Perform forward propagation and back propagation using example $(x^{(i)}, y^{(i)})$
get activations $a^{(l)}$ and delta terms $d^{(l)}$ for $l=2, 3, \dots, L$

[Evaluating the hypothesis]

As usual we should choose the parameters that minimize the errors. It can happen that we overfit.

Learn parameter Θ from training data (minimizing training error $J(\Theta)$). At 70% of the samples

Once we've done some troubleshooting for errors in our predictions, we can try to:

- get more training examples
- try a smaller set of features
- try additional features
- try polynomial features
- increase or decrease λ

Ideally don't chose them randomly, we'll explore diagnostic techniques & the best of above's solutions following on.

Evaluating an hypothesis

An hypothesis may have a low error for the training examples but still be incorrect (because of overfitting). Thus to evaluate a hypothesis, given a dataset of training examples, we can split up the data into two sets: training set and test set. Typically it's chosen at a ratio 70-30%.

The new procedure consists now on:

- 1: Learning Θ and minimize $J_{\text{train}}(\Theta)$ using the training set
- 2: Compute the test error $J_{\text{test}}(\Theta)$

The test set error

1 For linear regression: $J_{\text{test}}(\Theta) = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\Theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$

2 For classification: misclassification error (0/1 misclassification error)

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

this gives us a binary 0 or 1 result based on a misclassification.

The average test error for the test set is:

$$\text{TestError} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_{\Theta}(x_{\text{test}}^{(i)}, y_{\text{test}}^{(i)}))$$

this gives us the proportion of the test data that was misclassified.

[Diagnosing Bias vs Variance]

CV = cross validation

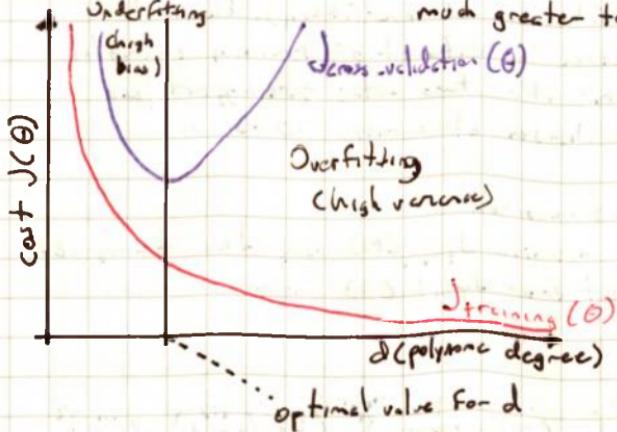
60-20-20%

In this section we examine the relationship between the degree of the polynomial ' d ' and the underfitting or overfitting of our hypothesis.

- We need to distinguish whether bias or variance is the problem contributing to bad decisions.
- High bias is underfitting and high variance is overfitting. Ideally we need to find a 'golden mean' between them.

The train error will tend to decrease as we increase ' d ' (the polynomial degree) up. At the same time, the cross validation error will tend to decrease as we increase ' d ' up to a point, and then it will increase as ' d ' is increased, forming a curvy curve.

- High Bias (Underfitting): both $J_{\text{train}}(\theta)$ and $J_{\text{cv}}(\theta)$ will be high.
Also $J_{\text{cv}}(\theta) \approx J_{\text{train}}(\theta)$
- High Variance (Overfitting): $J_{\text{train}}(\theta)$ will be low and $J_{\text{cv}}(\theta)$ will be much greater than $J_{\text{train}}(\theta)$



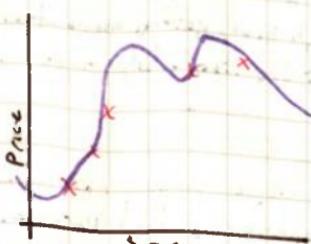
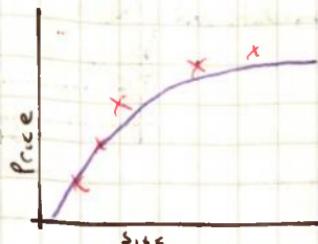
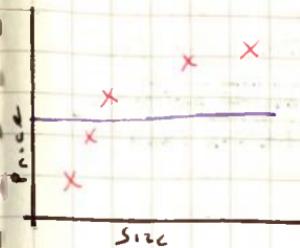
[Regularization and Bias / Variance]

Linear regression with regularization

For the cost function $J(\theta)$ we add the regularization term to adapt/soften/rounded regularize the cost. Now with $J_{cv}(\theta)$ or $J_{test}(\theta)$, we still use it but only for $J(\theta)$; the regularization is not used on $J_{cv}(\theta)$ neither on $J_{test}(\theta)$.

Model: $h_{\theta}(x) = \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

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



Large λ

High bias (underfit)

$\lambda > 10000$, $\theta_0 \approx 0$, $\theta_1 \approx 0 \dots$

$h_{\theta}(x) \approx \theta_0$

Intermediate λ

"Just right"

Small λ

High Variance (overfit)

$\lambda_{cv}, \lambda=0$

In the figures above, we see that as λ increases, our fit becomes more rigid. On the other hand, as λ approaches to 0, we tend to overfit the data. So, how do we choose our parameter λ to get it "just right"? In order to choose the model and λ :

1: Create a list of λ 's, i.e. $\lambda \in \{0, 0.01, 0.02, 0.04, 0.09 \dots\}$

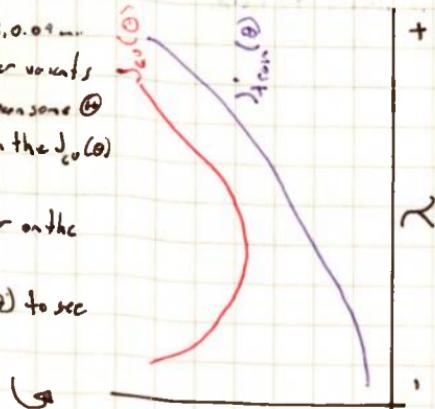
2: Create a set of models with different degrees or any other variants

3: Iterate through λ 's and for each λ , go thru all models to learn some θ

4: Compute CV error using the learned θ (computed with λ) on the $J_{cv}(\theta)$ without regularization or $\lambda=0$

5: Select the best combo that produces the lowest error on the CV set

6: Using the best combo θ and λ , apply it on $J_{test}(\theta)$ to see if it has a good generalization of the problem



Learning Curves

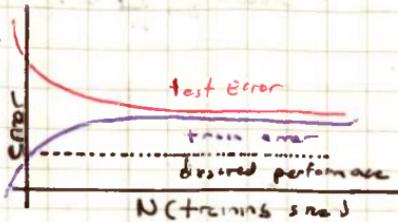
Training an algorithm on a very few number of datapoints (such as 1, 2 or 3) will easily have errors because we can always find a quadratic curve that touches exactly those points. Therefore:

- As the training set gets larger, the error for a quadratic function increases.
- The error will 'plateau' out after a certain size, or training set size.

Experiencing high bias:

- Low training set size: causes $J_{\text{train}}(\theta)$ to be low and $J_{\text{CV}}(\theta)$ to be high.
- Large training set size: causes both $J_{\text{train}}(\theta)$ and $J_{\text{CV}}(\theta)$ to be high, having $J_{\text{train}}(\theta) \approx J_{\text{CV}}(\theta)$.

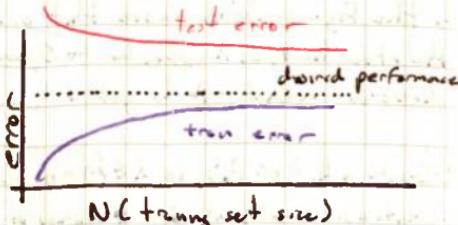
If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much.



Experiencing high variance:

- Low training set size: $J_{\text{train}}(\theta)$ will be low and $J_{\text{CV}}(\theta)$ will be high.
- Large training size: $J_{\text{train}}(\theta)$ increases with training set size and $J_{\text{CV}}(\theta)$ continues to decrease without leveling off.
Also, $J_{\text{train}}(\theta) < J_{\text{CV}}(\theta)$, but the difference between them remains significant.

If a learning algorithm is suffering from high variance, getting more training data is likely to help.



Deciding what to do next (continued)

Our decision process can be broken down as follows:

- Getting more training examples: Fixes high variance
- Trying smaller set of features: Fixes high variance
- Adding features: Fixes high bias
- Adding polynomial features: Fixes high bias
- Decreasing λ : Fixes high bias
- Increasing λ : Fixes high variance

Diagnosing Neural Networks:

- A neural network with fewer parameters is prone to underfitting, it's also computationally cheaper.
- A large neural network with more parameters is prone to overfitting. It's also computationally expensive. In this case, using regularization (increase λ) will address the overfitting.

Using a single layer ($L=3$) is a good starting default. You can train your neural network on a number of hidden layers using your cross-validation set. Then just select the one that performs the best.

Model Complexity Effects

- Lower order polynomials (low model complexity) have big bias and low variance. In this case the model fits poorly consistently.
- Higher-order polynomials (high model complexity) fit the training data extremely well and the test data extremely poorly. These have low bias on the training data, but very high variance.
- In reality, we would want to choose a model in between, that generalizes well, but also fits the data reasonably well.

[Error metrics for skewed classes]

When we have two samples which consist mainly on one specific classification (1 or 0), and our logistic regression model is highly oriented towards predicting ~~one~~ that class type. So it's influenced.

For example in the cancer detection problem, $h(x)$. ($y=1$ if cancer, $y=0$ otherwise) we find that we've got 1% error on the test set, but only 0.5% of patients have cancer.

When we get this kind of situation, we need a different error metric or error evaluation metric called Precision / Recall.

Precision: Of all patients that were predicted $y=1$, what fraction actually had cancer.

Actual Class		
	1	0
1	true positive	false positive
0	false negative	true negative

$$\frac{\text{true positives}}{\# \text{predicted } +} = \frac{\text{true positives}}{\text{true pos.} + \text{false pos.}}$$

R Recall: of all patients that actually have cancer, what fraction of them did we correctly detect as having cancer?

$$\frac{\text{true positives}}{\# \text{actual positives}} = \frac{\text{true positives}}{\text{true pos.} + \text{false neg.}}$$

Fraudulence

For some applications we want to control the tradeoff between precision and recall, to do it, we increment the threshold (or reduce it), as we want to provide confident information, and we do this in our sigmoid or logistic function, by predicting 1 if $h(x) \geq 0.7$ and 0 if $h(x) < 0.7$. This will have higher precision but lower recall. Likewise if we would like to be more 'flexible' into classifying more samples as positive, then the threshold should be lowered. This depends on the purpose, i.e. need to correct flag to avoid problems. This will have less precision and higher recall.

[F-score]: How to compare precision and recall numbers?

	Precision (P)	Recall (R)	Average	F-Score
Algorithm 1	0.5	0.4	0.45	0.444
Algorithm 2	0.7	0.1	0.4	0.175
Algorithm 3	0.02	1.0	0.51	0.0392

$$\text{Average} = \frac{P+R}{2}$$

$$P=0 \text{ or } R=0 \Rightarrow F\text{-score} = 0$$

$$F_1 \text{ score} = 2 \frac{PR}{P+R}$$

$$P=1 \text{ and } R=1 \Rightarrow F\text{-score} = 1$$

perfect F-score

Support Vector Machine Optimization Objective

(Week 7)

The Support Vector Machine (SVM) is yet another type of supervised machine learning algorithm. It is sometimes cleaner and more powerful.

Recalling that in logistic regression, we used the following rules:

If $y=1$, then $h_\theta(x) \approx 1$ and $\Theta^T x \gg 0$

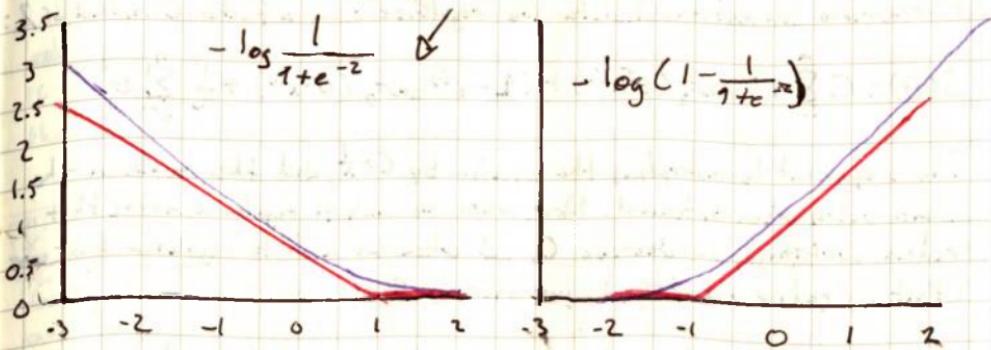
If $y=0$, then $h_\theta(x) \approx 0$ and $\Theta^T x \ll 0$

And the cost function for [unregularized] logistic regression:

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

$$= \frac{1}{m} \sum_{i=1}^m -y^{(i)} \log\left(\frac{1}{1+e^{-\Theta^T x^{(i)}}}\right) - (1-y^{(i)}) \log\left(1 - \frac{1}{1+e^{-\Theta^T x^{(i)}}}\right)$$

To make a support vector machine, we will modify the first term of the cost function: $\log(h_\theta(x)) = -\log\left(\frac{1}{1+e^{-\Theta^T x}}\right)$, so that when $\Theta^T x$ (from now on, we shall refer to this as z) is greater than 1, it outputs 0. Furthermore, for values of z less than one, we shall use a straight decreasing line instead of the sigmoid curve. In the literature this is called a "Hinge Loss Function"



Similarly, we modify the second term of the cost function $-\log(1-h_\theta(x)) = -\log\left(1 - \frac{1}{1+e^{-\Theta^T x}}\right)$, so that when Z is less than ~~-1~~ ≈ -1 , it outputs 0. We also modify it so that for values of z greater than -1 , we use a straight increasing line instead of the sigmoid curve.

(1) we shall denote these as $\text{cost}_1(z)$ and $\text{cost}_0(z)$ respectively, note that $\text{cost}_1(z)$ is the cost for classifying when $y=1$ and $\text{cost}_0(z)$ is the cost for classifying when $y=0$, and we define them as follows (where K is an arbitrary constant defining the magnitude of the slope of the line).

$$z = \theta^T x \quad \text{cost}_0(z) = \max(0, K(1+z)) \\ \text{cost}_1(z) = \max(0, K(1-z))$$

Recalling the cost function from (unregularized) logistic regression:

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

Note that the negative sign has been distributed into the sum of the above equation. We may transform this into the cost function for SVM by substituting $\text{cost}_0(z)$ and $\text{cost}_1(z)$:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

This can be optimized a bit by multiplying all by " -1 " (thus removing it from the denominators). Note that this does not affect our optimization since it's multiplied by a positive constant.

$$J(\theta) = \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) + \frac{\lambda}{2} \sum_{j=1}^n \theta_j^2$$

Furthermore convention dictates we regularize using a factor C instead of λ :

$$J(\theta) = C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

This is equivalent to multiplying the equation by $C=\frac{1}{\lambda}$, and thus results in the same values when optimized. Now, when we wish to regularize more (that is, reduce overfitting), we decrease C , and when we wish to regularize less, more (that is, reduce underfitting), we increase C .

One final note is that the hypothesis of the Support Vector Machine is not interpreted as the probability of y being 1 or 0 (as is the hypothesis of logistic regression). Instead, it outputs either 1 or 0 straight (its a discriminant function).

$$h_\theta(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

[Large Margin Intuition]

A useful way to think about Support Vector Machines, is to think about them as Large Margin Classifiers..

If $y=1$, we want $\theta^T x \geq 1$ (not just ≥ 0)

If $y=0$, we want $\theta^T x \leq -1$ (not just < 0)

Now when we set our constant to a very large value (eg 100,000), our optimizing function will constrain θ such as the first part of the equation equals zero. We impose the following constraints on θ :

$$\theta^T x \geq 1 \text{ if } y=1 \text{ and } \theta^T x \leq -1 \text{ if } y=0$$

If C is very large, we must choose θ parameters such that:

$$\sum_{i=1}^n y^{(i)} \text{cost}_0(\theta^T x) + (1-y) \text{cost}_1(\theta^T x) = 0 \quad \text{and this reduces our cost function to:}$$

$$J(\theta) = C \cdot 0 + \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Recalling the decision boundary from logistic regression (the line separating positive and negative examples). In SVM, the decision boundary has the special property that is as far away as possible from both, the positive and negative values.

The distance of the decision boundary to the nearest example is called the margin and thus, large margin is only achieved when C is very large. Data is linearly separable when a straight line can separate the positive and negative examples. If we have outlier examples that we don't want to affect the decision boundary, then we can reduce C (Increasing and decreasing C is similar to respectively decreasing and increasing λ), and can simplify our decision boundary.

Kernels 1 Kernels allow us to make complex, non-linear classifiers. Given X , compute new feature ϕ depending on proximity to landmarks $\ell^{(1)}, \ell^{(2)}, \dots$ and to do so, we should find the "similarity" of x and some landmark $\ell^{(i)}$.

$$f_i = \text{similarity}(x; \ell^{(i)}) = \exp\left(-\frac{\|x - \ell^{(i)}\|^2}{2\sigma^2}\right), \quad \text{using SVMs, these are the similarity func.}$$

$$f_i = \text{similarity}(x, \ell^{(i)}) = \exp\left(-\frac{\sum_{j=1}^n (x_j - \ell_j^{(i)})^2}{2\sigma^2}\right)$$

is called a "gaussian kernel". It is a specific example of a kernel. The function can also be written as follows:

There are a couple properties of the similarity function:

$$\text{if } x \approx l^{(i)}, \text{ then } f_i = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right) \approx 1 \quad \text{and if } x \text{ is far from } l^{(i)}$$
$$f_i = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0$$

In other words, if x and the landmark are close, the similarity will be close to 1, and if x and the landmark are far away from each other, the similarity will be close to 0.

Each landmark gives us the features in our hypothesis: $f_1^{(i)} = f_1$, $f_2^{(i)} = f_2$, \dots

$$f_1^{(i)} = f_1$$

$$f_2^{(i)} = f_2$$

$$f_3^{(i)} = f_3$$

$$\dots$$

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots$$

σ^2 is a parameter of the Gaussian kernel, and it can be modified to increase or decrease the drop-off for our feature f_i . Combined with looking at the values inside Θ , we can choose these landmarks to get the general shape of the ^{decision} boundary.

Kernels II: One way to get the landmarks is to put them in the exact same location as all the training examples. This gives us "m" landmarks, with one landmark per training example.

Given example x : $f_1 = \text{similarity}(x, l^{(1)})$, $f_2 = \text{similarity}(x, l^{(2)})$, $f_3 = \text{similarity}(x, l^{(3)})$, and so on.

This gives us a 'feature vector', $f^{(i)}$ of all our features for the example $x^{(i)}$. We may also set $f_0 = 1$ to correspond with θ_0 . Thus given training example $x^{(i)}$:

$$x^{(i)} = \begin{bmatrix} f_1^{(i)} = \text{similarity}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} = \text{similarity}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_m^{(i)} = \text{similarity}(x^{(i)}, l^{(m)}) \end{bmatrix}$$

Now to get the parameters Θ we can use the SVM minimization algorithm but with $f^{(i)}$ substituted in for $x^{(i)}$:

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_+ (\theta^T f^{(i)}) + (1 - y^{(i)}) \text{cost}_- (\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Using Kernels to generate $f^{(i)}$, is not exclusive for SVMs and may also be applied to logistic regression. However, because of computational optimizations on SVMs, kernels combined with SVMs is much faster than with other algorithms, so kernels are almost always found combined only with SVMs.

Choosing SVM parameters

Choosing C ($C \approx \frac{1}{\lambda}$)

If C is large, then we get higher variance and lower bias

If C is small, then we get lower variance and higher bias

Choosing σ^2

Large σ^2 : Features $f^{(i)}$ very non-smoothly, causing higher bias and lower variance

Small σ^2 : Features very smooth, causing lower bias and higher variance.

Using a SVM

There are lots of good SVM libraries out there, such as 'liblinear' or 'libsvm'. And in real life is much better to use them rather than reinventing the wheel. In real life there are choices to be made, such as:

- Choice of parameter C
- No Kernel ("linear" kernel) gives a flat decision boundary.
- Choice of Kernel (the similarity function)
- Choose when α is large and m is small
- Gaussian Kernel - need to choose $\sigma^2 \rightarrow$ choose when α is small and m is large
- The library may not provide the kernel function.

Notes:

- ~~DO~~ perform feature scaling before using Gaussian Kernel
- Not all similarity functions are valid kernels, they must satisfy the 'Mercer's' theorem. Which worries SVM package optimizers can't solve and won't converge.

Multi Class Classification:

Many SVM's have multiclass classification built in. It is possible to use as well one-vs-all method just as we did for logistic regression, where $y \in \{1, 2, 3, \dots, k\}$ with $\Theta^{(1)}, \Theta^{(2)}, \dots, \Theta^{(k)}$. We pick the class i with the largest $(\Theta^{(i)})^T x$.

Logistic Regression vs. SVMs

- If n is large (relative to m), then use logistic regression, or SVM without kernel (linear).
- If n is small and m is intermediate, then use SVM with Gaussian Kernel.
- If n is small and m is large, then manually create/add more features, then use logistic regression or SVMs without Kernel

In the first case, we don't have enough examples to need a complicated polynomial hypothesis

In the second example, we have enough examples, that we might need a complex non-linear hypothesis

In the third example, we want to increase our features so that logistic regression becomes applicable.

NOTE: A neural network is likely to work well for any of those situations but may be slower to train.

[Week 8 - Clustering]

[Unsupervised learning: Introduction]

(Unsupervised) learning is contrasted from supervised learning because it uses an **unlabeled training** set rather than a labeled one. In other words, we don't have the vector Y of expected results, we only have a dataset of features where we can find structure.

Clustering is good for:

- Market segmentation
- Social networks analysis
- Organizing computer clusters
- Astronomical data analysis

[K-means algorithm]

The K-Means algorithm is the most popular and widely used algorithm for automatically grouping data into coherent subsets. It works as follows:

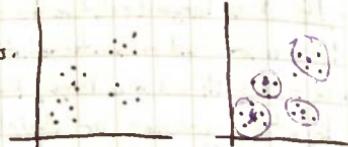
- 1: Randomly initialize $2(K)$ points in the dataset, called the **cluster centroids**.
- 2: Cluster assignment: assign all examples into one of $2(K)$ groups based on which cluster centroid the example is closest to.
- 3: Move centroid: Compute the averages of all the points in each of the two cluster centroid groups, then move the centroids to those averages.
- 4: Repeat #2: Until we've found our clusters.

Our main variables are:

- K (number of clusters)
- Training set $x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(n)}$ } Note we want use the $x_0 = 1$ constant.
- Where $x^{(i)} \in \mathbb{R}^n$ - n dimensional

The algorithm:

- 1 Randomly initialise K cluster centroids $m_0(1), m_0(2), \dots, m_0(K)$
- 2 Repeat:
 - 3 For $i = 1$ to K :
 - 4 $c(i) :=$ index (from 3 to K) of cluster centroids closest to $x^{(i)}$
 - 5 for $k = 1$ to K :
 - 6 $m_k(k) :=$ average (mean) of points assigned to cluster k



The first for-loop (lines 3 and 4) is the 'cluster assignment' step. We make a vector c where $c(i)$ represents the centroid assigned to example $x^{(i)}$. We can write the operation of cluster assignment more mathematically as follows: $c^{(i)} = \arg \min_k \|x^{(i)} - \mu_k\|^2$, that is each $c^{(i)}$ contains the index of the centroid that has minimal distance to $x^{(i)}$.

By convention, we square the right-hand side, which makes the function we are trying to minimize "more sharply increasing". It is mostly just a convention. But a convention that helps reduce the computational load because the Euclidean distance requires a square root but it's cancelled.

- Without the square: $\|x^{(i)} - \mu_k\| = \sqrt{(x_1^{(i)} - \mu_{k(1)})^2 + (x_2^{(i)} - \mu_{k(2)})^2 + (x_3^{(i)} - \mu_{k(3)})^2 + \dots}$
 - With the square: $\|x^{(i)} - \mu_k\|^2 = (x_1^{(i)} - \mu_{k(1)})^2 + (x_2^{(i)} - \mu_{k(2)})^2 + (x_3^{(i)} - \mu_{k(3)})^2 + \dots$
- so squaring convention serves two purposes: minimize sharply and less computationally*

The second for-loop is the "move centroid" step where we move each centroid to the average of its group. More formally, the equation is as follows:

$$\mu_k = \frac{1}{n} [x^{(k_1)} + x^{(k_2)} + \dots + x^{(k_n)}] \in \mathbb{R}^n$$

Where each of $x^{(k_1)}, x^{(k_2)}, x^{(k_3)}, \dots, x^{(k_n)}$ are the training examples assigned to group ~~mpk~~ μ_k . If we happen to have a cluster centroid with 0 points assigned to it, we can randomly-reintroduce that centroid to a new point, or simply eliminate that cluster group.

After a number of iterations, the algorithm will converge where new iterations won't affect the clusters.

[Optimization Objective]

Recalling some of the quantities we've used in our algorithm:

- $c^{(i)}$: index of cluster ($1, 2, \dots, K$) to which example $x^{(i)}$ is assigned currently
- μ_k : cluster centroid K ($\mu_k \in \mathbb{R}^n$)
- $\mu_{c(i)}$: cluster centroid of cluster to which example $x^{(i)}$ has been assigned
- Using these variables, we can define our cost function:

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c(i)}\|^2$$

Our optimization objective is to minimize all of our parameters using the cost function. That is, we are finding all values in sets C , representing all of our clusters, and N , representing all of our centroids, that will minimize the average of the distances of every training example to its corresponding centroid cluster centroid. This ~~cost~~ function is also called sometimes "the 'distortion' of the training examples".

With K-means, it is not possible for the cost function to increase, it should always decrease.

[Random Initialization]

There is one particular recommended method for randomly initializing C if $K < n \sim N$: Make sure there are less clusters than training samples.

2: Randomly pick K training examples \sim be sure they are unique.

3: Set N_1, \dots, N_K equal to these K examples.

K-means can get stuck on local optima. To decrease the chance of this happening you can run the algorithm on many different random initializations. In cases where $K > 10$ it's strongly recommended to run a loop of random initializations.

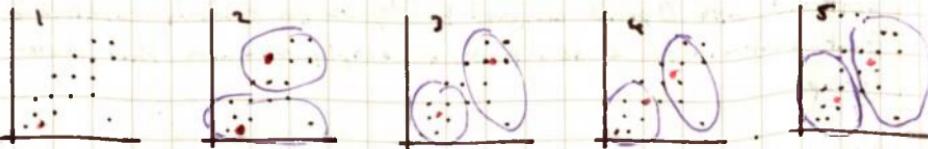
2 for $i=1$ to 100 :

2 randomly initialize K-means

3 For K-means to get ' c ' and ' m '

4 compute the cost function (distortion) $J(C, m)$

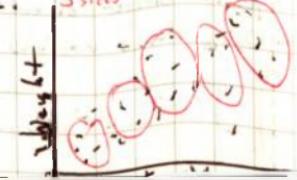
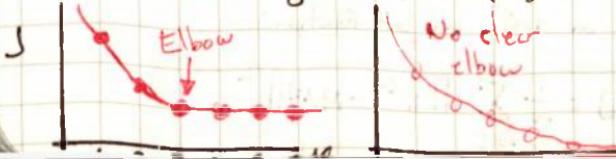
5 pick the clustering that gave us the lowest cost



[Choosing the Number of Clusters]

Choosing K can be quite arbitrary and ambiguous.

- The elbow method: Plot J and # of clusters K ; J should ~~reduce~~ as K increases; However fairly often, the curve is very gradual, so there is no clear elbow, and also J will always decrease as K is increased. Only exception is if K gets stuck in a local optimum.
- The downstream purpose: Assign it arbitrarily, choose K that ~~proves~~ to be more useful to the goal we are trying to achieve.

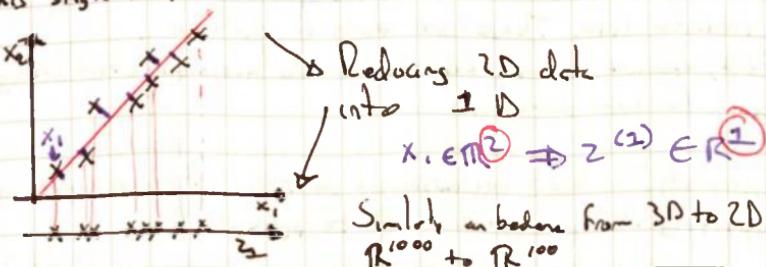


[Dimensionality Reduction]

There are a couple of reasons why we would like to do Dimensionality Reduction, whether is to do data compression, which will use less computer memory but also speed up our learning algorithms.

[Motivation 1]: Data Compression:

- We may want to reduce the dimension of our features if we have redundant data
- To do this we find highly correlated features, plot them, and make a new line that seems to describe both features accurately. We place new features on this single line.



this will reduce the total data we have to store in computer memory and will speedup our learning algorithm.

Note: In dimensionality reduction, we are reducing the number of features rather than our number of examples. Our variable m will stay exactly the same, however the number of n will be reduced.

[Motivation 2]: Visualization

It is not easy to visualize data that is more than 3D, we can then reduce the dimensions of data to 3 or less in order to plot it. One example is: Having hundreds of features related to a country's economic system, which may be combined into one feature that you can call "Economic Activity"

Principal Component Analysis - Problem Formulation

The most popular dimensionality reduction algorithm is Principal Component Analysis (PCA)

Problem Formulation: Given two features X_1 and X_2 , we want to find a single line that effectively describes both features at once. We then map our old features onto this new line to get a new single feature.

The goal of PCA is to reduce the average of all the distances of every feature to the projection line, thus is: the projection error.

Reduce from n -d to 2d to 1d: Find a direction (a vector $u^{(1)} \in \mathbb{R}^n$) onto which project the data so as to minimize the projection error.

At more general case:

- Reduce from n -dimensional to k -dimension: Find k vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project the data so as to minimize the projection error. If we are converting from 3d into 2d, we will project our data onto two directions (a plane), so k will be 2.

[PCA is not linear regression]:

- In linear regression, we are minimizing the squared error from every point to our predictor line. Those are vertical distances
- In PCA, we are minimizing the shortest distance, or shortest orthogonal distances, to our datapoints

More generally, in linear regression we are taking all our examples in X and applying the parameters in Θ to predict Y .

In PCA, we are taking a number of features x_1, x_2, \dots, x_n , and finding a closest common dataset among them. We aren't trying to predict any result and we aren't applying any theta weights to the features.

[Principal Component Analysis Algorithm]

Before we can apply PCA, there is a data-preprocessing we must perform

• Data Preprocessing:

- Given a training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$
- Preprocess (feature scale / mean normalization):

$$\mu_j = \frac{1}{n} \sum_{i=1}^n x_j^{(i)}$$

- Replace each $x_j^{(i)}$ with $x_j^{(i)} - \mu_j$

- If there are different features in different scales (e.g. x_1 = size of the house and x_2 = number of bedrooms), scale features to have comparable range of values.

3 Compute Covariance Matrix &

Above we first subtract the mean of each feature from the original feature and then we can scale all the features $x_j^{(i)} = \frac{x_j^{(i)} - \mu_j}{\sigma_j}$

We can define what it means specifically to reduce from 2D to 1D as:

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T \quad \begin{array}{l} \text{Vectorized} \\ \text{Implementation} \end{array} = \text{Sigma} = 1/m * X^T * X$$

We denote the covariance matrix with a capital sigma (Σ) which happens to be the same symbol for summation, confusingly - they represent entirely different things.

Note that $x^{(i)}$ is a $n \times 1$ vector, $x^{(i)T}$ is a $1 \times n$ vector and X is a $m \times n$ matrix (row-wise stored examples). The product of those will be an $n \times n$ matrix, which are the dimensions of Σ .

② Compute 'eigenvectors' of covariance matrix Σ

$$[U, S, V] = svd(\text{Sigma});$$

"SVD" is the "singular value decomposition", a built-in function.

What we actually went out of $svd()$ is the " U " matrix of the Sigma covariance matrix: $U \in \mathbb{R}^{n \times n}$. U contains $U^{(1)}, \dots, U^{(n)}$, which is what we want.

③ Take the first K columns of the U matrix to compute Z

We'll assign the first K columns of U to a variable called "UReduce", this will be a $n \times K$ matrix. We compute Z with:

$$Z^{(i)} = \text{UReduce}^T \cdot x^{(i)}$$

$\text{UReduce}^T \cdot x^{(i)}$ will have dimensions $K \times 1$, while $x^{(i)}$ will have dimensions $n \times 1$. The product $\text{UReduce}^T \cdot x^{(i)}$ will have dimension $K \times 1$.

To summarize, the whole algorithm is:

$$\boxed{\begin{aligned} & \text{Sigma} = \frac{1}{m} X^T \cdot X; \\ & [U, S, V] = svd(\text{Sigma}); \\ & \text{UReduce} = U(:, 1:K); \\ & Z = X \cdot \text{UReduce}; \end{aligned}}$$

Reconstruction from Compressed Representation

If PCA can be used to compress, also can be used to get an approximation of the original data and same original amount of features. (e.g. $Z \in \mathbb{R}^{n \times D} \approx \hat{X} \in \mathbb{R}^{n \times n}$)

This can be done with the equation $X_{\text{approx}}^{(i)} = \text{UReduce} \cdot Z^{(i)}$

Turns out that the U matrix has the special property that is an Unitary Matrix and one of the matrix properties is that $U^{-1} = U^*$ (*=Conjugate transpose). Since we are dealing here with real numbers, this is equivalent to: $U^{-1} = U^T$ so we could compute the inverse and use that, but it would be a waste of energy and cycles.

[Choosing the Number of Principal Components]

How do we choose k , also called the 'number of principal components'? Recall that k is the dimension we are reducing to.

One way to choose k s is by using the following formula:

Given the squared projection error $\frac{1}{m} \sum_{i=1}^m \|(\mathbf{x}^{(i)} - \mathbf{x}_{\text{approx}}^{(i)})\|^2$

Also given the total variation of data $\frac{1}{m} \sum_{i=1}^m \|\mathbf{x}^{(i)}\|^2$

(choose k to be the smallest value such that: $\frac{\frac{1}{m} \sum_{i=1}^m \|(\mathbf{x}^{(i)} - \mathbf{x}_{\text{approx}}^{(i)})\|^2}{\frac{1}{m} \sum_{i=1}^m \|\mathbf{x}^{(i)}\|^2} \leq 0.01$)

In other words: the squared projection error divided by the total variation should be less than one percent, so that 99% of the variance is retained.

Algorithm for choosing K

1 Try PCA with $k=1, 2, \dots$

2 Compute $\text{Var}(z_i)$, $z_i \in \mathbb{R}^n$

3 Check the formula compliance, if not repeat for 99% with the S matrix as:

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99$$

[Advice for Applying PCA]

- The most common use of PCA is to speed up supervised learning

Given a large training set with a large number of features (e.g. $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$) for \mathbb{R}^{1000} we can use PCA to reduce the number of features in each example of the training set, so $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(n)} \in \mathbb{R}^{1000} \rightarrow \mathbf{z}^{(1)}, \dots, \mathbf{z}^{(n)} \in \mathbb{R}^{100}$

- Note that we should define the PCA reduction from $\mathbf{x}^{(1)} \dots \mathbf{x}^{(n)}$ only on training set, and not on Cross Validation or Test sets. We can then apply the mapping $\mathbf{z}^{(i)}$ to the CV or Test sets, after it is defined on the training set.

- Applications:

• Compression: Reduce space of data
Speed up algorithms

• Visualization: Choose 2 or 3 dimensions

Bad use of PCA: Trying to prevent overfitting. We may think that

reducing the number of features with PCA would be a good way to address overfitting. It might work but using regularization is preferred because PCA doesn't take into account \mathbf{Y} 's. Try first the full algorithm and then attempt PCA

Week 9: Anomaly Detection

As in any other machine learning problem, we are given a dataset $x^1, \dots, x^{m'}$ and given an example x_{test} we want to know if this example is normal or anomalous.

We define a "model" $P(x)$ that tells us the probability that the example is not anomalous. We also use a threshold ϵ (Capita) as a dividing line so we can say which examples are anomalous and which aren't.

A very common application of anomaly detection is detecting a fraud:

- $x^{(i)}$ = features of user's i 's activities
- model $P(x)$ from the data
- identify unusual users by checking which have $P(x) < \epsilon$

If our anomaly detector is flagging too many anomalous examples, then we need to decrease our threshold ϵ .

Gaussian Distribution

The Gaussian Distribution is a familiar bell-shaped curve that can be described by a function $N(\mu, \sigma^2)$.

Let $X \in \mathbb{R}$, if the probability distribution of x is Gaussian with mean μ , variance σ^2 , then: $x \sim N(\mu, \sigma^2)$ | the field \sim means "distributed as"

The Gaussian Distribution is parametrized by a mean and a variance, μ or N describes the center of the curve, called the mean. The width of the curve is described by sigma or σ , called the standard deviation.

The full function is as follows: $P(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$

We can estimate μ from a given dataset by taking the average of samples \bar{x}

$\mu = \frac{1}{n} \sum_{i=1}^n x^{(i)}$ and for σ^2 , with our familiar squared error formula: $\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu)^2$

Algorithm

Given a training set of examples, $\{x^{(1)}, \dots, x^{(n)}\}$ where each example is a vector $x \in \mathbb{R}^n$:

$$p(x) = p(x_1; \mu_1, \sigma_1^2) \cdot p(x_2; \mu_2, \sigma_2^2) \cdot \dots \cdot p(x_n; \mu_n, \sigma_n^2)$$

In statistics this is called "independence assumption" on the values of the features inside training examples. More compactly, the above expression can be written as follows:

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

Steps

- Choose features x_j that I think might be indicative of anomalous examples
- fit parameters $\mu_1, \dots, \mu_n; \sigma_1^2, \dots, \sigma_n^2 \rightarrow \boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$
- Calculate $\mu_j = (\sum_{i=1}^n x_i^{(j)}) / n$
- Calculate $\sigma_j^2 = 1/n \cdot \sum_{i=1}^n (x_i^{(j)} - \mu_j)^2$
- Given a new example x , compute $p(x)$:

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

- anomaly if $p(x) < \epsilon$

Developing and Evaluating an Anomaly Detection System

To evaluate our learning algorithm, we take some labeled data, categorized into anomalous and non-anomalous examples ($y=0$ if normal, $y=1$ if anomalous)

Among that data, take a large proportion of 6000, non-anomalous data for the training set on which to train $P(x)$. Then take a smaller proportion of mixed anomalous and non-anomalous data examples (you will usually have many more non-anomalous examples) for your cross-validation and test sets.

For example, we may have a set where 0.2% of the data is anomalous. We take 60% of those examples, all of which are good ($y=0$) for the training set. We then take 20% of the examples for the cross-validation set (with 0.1% of the anomalous examples) and another 20% from the test set (with another 0.1% of anomalies).

In other words, we split the data in 60/20/20, training/CV/test, and then split the anomalous examples 50/50 between the CV and test sets.

Algorithm evaluation:

- Fit model $p(x)$ on training set $\{x^{(1)}, \dots, x^{(m)}\}$ } Possible evaluation metrics
- On a cross validation / test example x , predict:
 - if $p(x) < \epsilon$, then $y=1$
 - if $p(x) \geq \epsilon$, then $y=0$

- True positive / false positive
- False negative / true negative
- precision / recall
- F₁ score

* Note that we use the CV set to choose ϵ

[Anomaly Detection vs Supervised Learning]

When do we use anomaly detection and when do we use supervised learning?

User Anomaly Detection when... | Use Supervised Learning When...

We have very small number of positive examples ($y=1$ ~ 0-20 examples is correct) with a large number of negative ($y=0$) examples

We have many different types of anomalies and it's hard for any algo. to learn from positive examples what the anomalies look like; future anomalies may look nothing like any of the anomalies examples we've looked at.

We have a large number of both positive and negative examples; In other words, the training set isn't skewed, it is more evenly divided into classes.

We have enough positive examples for the algorithm to get a sense of what new positive examples look like. The future positive examples are likely similar to the ones in the training set.

Choosing what features to use

I.E.: create new features like

CPU
network usage or CPU²/Network

The features will greatly affect how well our anomaly detection algorithm works; we can check that our features are gaussian by plotting a histogram of our data and checking for the bell shaped curve. Some transforms we can try on an example feature x that does not have a bell shaped curve:

- $\log(x)$
- $\log(x+1)$
- $\log(x+c)$ for some constant
- $\sqrt{x} = x^{1/2}$
- $x^{-1/3}$

We can play with each of these to try and achieve the gaussian shape in our data

There is an error analysis procedure for anomaly detection that is very similar to the one in supervised learning. Our goal is for $p(x)$ to be large for normal examples and small for anomalies. One way is to examine the anomalies and come up with new features that will better distinguish the data. In gen., choose features that will be unusually large or small in the event of anomaly

[Multivariate Gaussian Distribution]

Multivariate Gaussian Distribution is an extension of anomaly detection and may or may not catch more anomalies.

Instead of modeling $p(x_1), p(x_2), \dots$ separately, we will model $p(x)$ all in one go. Our parameters will be $\mu \in \mathbb{R}^n$ and $\Sigma \in \mathbb{R}^{n \times n}$.

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

The important effect is that we can model along gaussian contours, allowing us to better fit data that might not fit into the normal circular contours. Varying Σ changes the shape, width and orientation of the contours, changing μ will move the center of the distribution.

[Anomaly Detection using the Multivariate Gaussian Distribution]

When doing anomaly detection with multivariate gaussian distribution, we compute μ and Σ normally. We then compute $p(x)$ using the new formula in the previous section and flag an anomaly if $\text{log}(p(x)) < \epsilon$.

The original model for $p(x)$ corresponds to a multivariate gaussian where the contours of $P(x; \mu, \Sigma)$ are axis-aligned.

The original model for maintains some advantages: it is computationally cheaper as there is no matrix to invert, which is costly for a large number of features, ... and it performs well even with small training set size. (In multivariate gaussian models, it should be greater than the number of features for Σ to be invertible).

The multivariate gaussian model can automatically capture correlations between different features of x .