# Machine Learning

**Nagubandi Krishna Sai**
**MS20BTECH11014**

**June 2021**

# Contents

# Chapter 1

# Fundamentals of Machine Learning

## 1.1 Supervised Learning

Supervised Learning gives "correct answers", the output values are same as real life values.

In Supervised Learning, we are given a set of data and we know what our correct output should look like, having an idea that there is relationship between the input and the output.

Supervised Learning problems has two types of problems,

1. Regression.

2. Classification.

### 1.1.1 Linear Regression

In regression type of problems, we are trying to predict results within a continuous output, means that we are trying to map input variables to some continuous function.

Linear regression has real-valued output, but the output will be same or near valued to the actual output.

1. m = Number of training examples.

2. x's = "input" variable (or) feature.

3. y's = "output (or) target" variable.

4. (x,y) = one training example.

5. $(x^{(i)}, y^{(i)}) = i^{th}$ training example.

| Size of feet$^2(x)$ | Price(\$) in 1000's (y) |
|:---:|:---:|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| $\vdots$ | $\vdots$ |

Table 1.1: Training set of housing prices.

Example :

1. m = 47.

2. $x^{(1)} = 2104$ and $y^{(1)} = 460$.

3. $x^{(2)} = 1416$ and $y^{(2)} = 232$.

Training set

Training set is feeded to learning algorithm

Learning algorithm

h    (hypothesis, This is a function that takes input x and estimate value of y.)

### 1.1.2 Hypothesis for Linear Regression

$$h_\theta(x) = \theta_0 + \theta_1.x \tag{1.1}$$

$\theta_i$'s = parameters.

This type of hypothesis mode is called "Linear regression with one variable" (or) "Univariate linear regression."

### 1.1.3 Cost function for Linear Regression

Cost function helps us know that how well to fit the best possible straight line over the given data.

Q⟩ *How to choose $\theta_i$'s ?*

1. Choose $\theta_i$ 's so that $h_\theta(x)$ is close to y for our training example (x,y).

2. minimise $\theta_0,\theta_1$ so, that $[h_\theta(x)$ - y] is small.

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} [h_\theta(x^{(i)}) - y^{(i)}]^2 \qquad (1.2)$$

J($\theta_0,\theta_1$) = Cost function (or) Squared error function.

### 1.1.4 Gradient descent for Linear Regression

Gradient descent is used to minimise cost function(J) in linear regression.
Gradient descent is used in many areas to minimise many functions in ML/AI.
**Gradient descent algorithm,**

Repeat until convergence (minimum) $\left\{ \theta_j := \theta_j - \alpha \dfrac{\partial}{\partial \theta_j} J(\theta_0, \theta_1), \; for \; j = 0, j = 1. \right.$

$$\qquad (1.3)$$

1. := is Assignment operator.

2. $\alpha$ is learning rate.

3. $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ is derivative.

Gradient descent is nothing but the derivative of the Cost function.

$$Slope \; of \; cost \; function \; curve = \frac{\partial J(\theta_1)}{\partial \theta_1}, \; when \; \theta_0 = 0. \qquad (1.4)$$

**Learning rate,**

1. If $\alpha$ is too small, gradient descent can be slow. After many such opera-tions(can be infinite times), the '$\theta_1$' could reach "global minimum".

2. If $\alpha$ is too large, gradient descent can overshoot the minimum. It may "fail to converge (or) even diverge".

3. If $\theta_1$ is at the local optima itself when we started or taken $\theta_1$, then there is no use of "$\alpha$ (or) gradient descent".

### 1.1.5 Linear Regression for multivariables

**Hypothesis,**

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + ... + \theta_n x_n. \tag{1.5}$$

In the total context of Supervised learning, hypothesis is just predicting the output.

For convenience of notation, declare $x_0 = 1$ ($x_0^{(i)} = 1$).

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \epsilon \, \Re^{n+1} \qquad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \epsilon \, \Re^{n+1}$$

The above matrix is 0 - indexed.

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + ... + \theta_n x_n. h_\theta(x) = \theta^\top x. \tag{1.6}$$

**Cost function,**

$$J(\theta) = J(\theta_0, \theta_1, \theta_2, ..., \theta_n) = \frac{1}{2m} \sum_{i=1}^{m} [h_\theta(x^{(i)}) - y^{(i)}]^2 \tag{1.7}$$

**Gradient descent,**

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1, \theta_2, ..., \theta_n) \tag{1.8}$$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \tag{1.9}$$

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^{m} [h_\theta(x^{(i)}) - y^{(i)}] x_j^{(i)} \tag{1.10}$$

**Feature scaling,**
Get every feature into approximately $-1 \leq x_i \leq 1$ *range.*
**Mean normalization,**
*Replace* $x_i$ *with* $x_i - \mu_i$ *to make features have approximately zero mean.*

### 1.1.6 Polynomial regression

**Hypothesis,**

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3. \tag{1.11}$$

| Housing price prediction |
|:---:|
| $x_1 = size$ |
| $x_2 = (size)^2$ |
| $x_3 = (size)^3$ |

Table 1.2: Features be-like in Polynomial regression.

## 1.1.7 Normal Equation

**Intuition,**

$$\frac{d}{d\theta} J(\theta) = 0.$$

**Cost function,**

$$\theta \; \epsilon \; \Re^{n+1}, \;\; J(\theta_0, \theta_1, \theta_2, ..., \theta_n) = \frac{1}{2m} \sum_{i=1}^{m} [h_\theta(x^{(i)}) - y(i)]^2 \qquad (1.12)$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = 0, \;\; (solve \; for \; \theta_0 \; , \theta_1 \; , ... \; , \theta_n) \qquad (1.13)$$

**Example,**

| | Size (feet$^2$) | No.of Bed rooms | No.of floors | Age of home (years) | Price($1000) |
|:---:|:---:|:---:|:---:|:---:|:---:|
| $x_0$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | y |
| 1 | 2104 | 5 | 1 | 45 | 460 |
| 1 | 1416 | 3 | 2 | 40 | 232 |
| 1 | 1534 | 3 | 2 | 30 | 315 |
| 1 | 1852 | 2 | 1 | 36 | 178 |

Table 1.3: Sample Training set for Multi-Variate Linear regression.

1. n = number of Features.

2. $x_j^{(i)}$ = value of j in the $i^{th}$ training example.

3. $x_{(i)}$ = the input(features) of the $i^{th}$ training example.

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \qquad Y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

X is m×(n+1)-dimensional matrix and Y is a m-dimensional vector.

$$\theta = (X^\top X)^{-1} X^\top Y \qquad\qquad (1.14)$$

The above $\theta$ value is optimal $\theta$ value.

For Normal equation method, then no need to use **feature scaling.**

We use 'Gradient descent' and 'Normal equation' methods to minimise cost function.

| Gradient descent | Normal equation |
|---|---|
| 1⟩ Need to choose '$\alpha$'. | 1⟩ No need to choose '$\alpha$'. |
| 2⟩ Need many iterations. | 2⟩ Don't need to iterate. |
| 3⟩ Works well even, when 'n' is large. (n¿10000) | 3⟩ Need to compute n× n matrix inverse $(X\top X)^{-1}$ |
| | 4⟩ Works Now if n is very large. |

Table 1.4: Why should we use the particular method? Advantages and Disadvantages of two methods.

$$\theta = (X^\top X)^{-1} X^\top Y \qquad\qquad (1.15)$$

Q⟩ What is $X^\top X$ is non-invertible(singular/degenerate) ?

Reasons,

1. Redundant features (linearly dependent)

    - $x_1 =$ Size in $feet^2$
    - $x_2 =$ Size in $m^2$
    - $x_2 = (3.28)^2 \ x_1$ , 1m = 3.28feet.

2. Too many features. (m≤n)

    - m = 10
    - n = 100, $\theta \ \epsilon \ \Re^{101}$, Delete some features (or) Regularization.

### 1.1.8   Classification

The output value 'y' is **discrete value.**

The algorithm used is **logistic regression.**

### 1.1.9   Logistic Regression Model

We want $0 \le h_\theta(x) \le 1$.

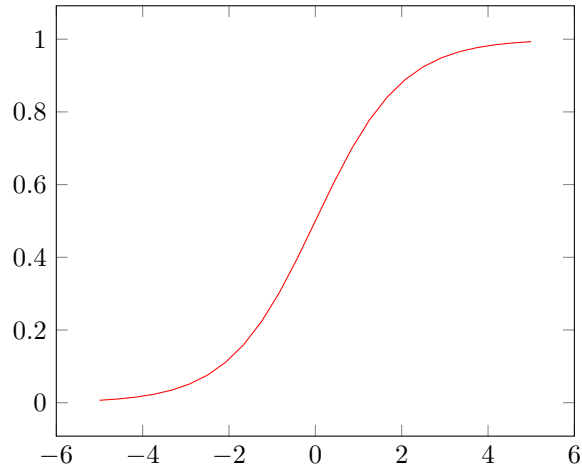$$h_\theta(x) = g(\theta^\top x) \tag{1.16}$$

$$g(z) = \frac{1}{1 + e^{-z}} \tag{1.17}$$

$$h_\theta(x) = g(\theta^\top x) \tag{1.18}$$

$$= \frac{1}{1 + e^{-\theta^\top x}} \tag{1.19}$$

The above g(z) is called sigmoid function (or) logistic function.
**Graph,**



$$g(z) \ge 0.5, \ when \ z \ge 0. \tag{1.20}$$

$$h_\theta(x) = g(\theta^\top x) \ge 0.5, \ when \ \theta^\top x \ge 0. \tag{1.21}$$

### 1.1.10   Interpretation of hypothesis output for Logistic Regression

$$h_\theta(x) = P(y = 1|x; \theta) \tag{1.22}$$
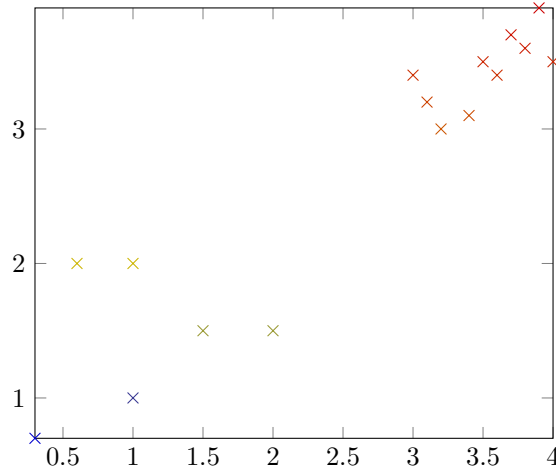
$$P(y = 0|x; \theta) + P(y = 1|x; \theta) = 1 \tag{1.23}$$

$$P(y = 0|x; \theta) = 1 - P(y = 1|x; \theta) \tag{1.24}$$

$$= 1 - h_\theta(x) \tag{1.25}$$

$$\tag{1.26}$$

$$y = 0 \text{ (or) } 1.$$

**Decision boundary,**



## Decision Boundary



For, the above diagram decision boundary will be a line separating the two output values of y(y=0 (or) y=1).

$$h_\theta(x) \geq 0.5 \rightarrow y = 1. \tag{1.27}$$
$$h_\theta(x) < 0.5 \rightarrow y = 0. \tag{1.28}$$
$$h_\theta(x) = g(\theta^\top x) \tag{1.29}$$
$$= g(\theta_0 + \theta_1 x_1 + \theta_2 x_2) \tag{1.30}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \qquad x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

**Example,**

Let,

$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

$$y = 1, if\ \theta^\top x \geq 0 \tag{1.31}$$
$$-3 + x_1 + x_2 \geq 0 \tag{1.32}$$
$$x_1 + x_2 \geq 3,\ y = 1 \tag{1.33}$$
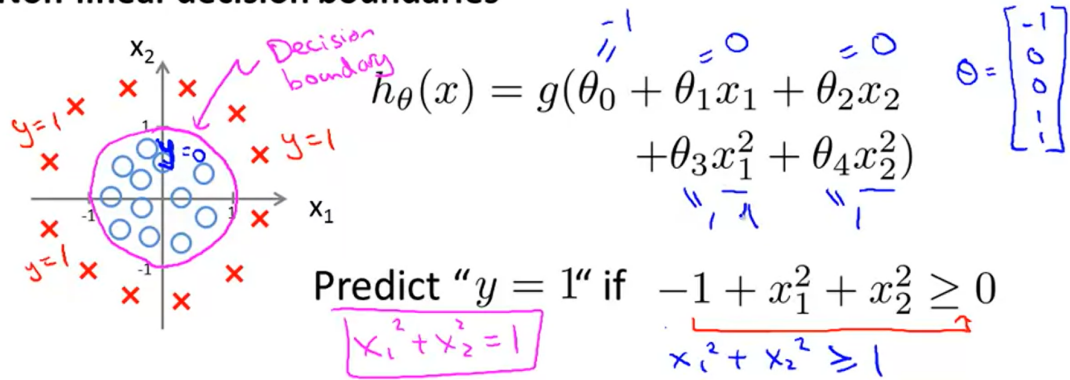$$x_1 + x_2 < 3,\ y = 0. \tag{1.34}$$

**Non-Linear Decision Boundary,**



So, in the above non-linear classification, the **decision boundary is a circle** of radius of 1unit.

$$Inside\ circle,\ y = 0 \tag{1.35}$$
$$Outside\ circle,\ y = 1. \tag{1.36}$$

## 1.1.11 Cost function for Logistic Regression

**Training set,**

| $x^{(1)}$ | $y^{(1)}$ |
|-----------|-----------|
| $x^{(2)}$ | $y^{(2)}$ |
| $x^{(3)}$ | $y^{(3)}$ |
| $\vdots$ | $\vdots$ |
| $x^{(m)}$ | $y^{(m)}$ |

Table 1.5: Training set of m-examples for Logistic Regression

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \epsilon \, \Re^{n+1} \; - \; n \; features. \; x_0 = 1, \; y \, \epsilon \, 0, 1.$$

For linear regression,

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} [h_\theta(x^{(i)} - y^{(i)}]^2 \qquad (1.37)$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} Cost(h_\theta(x^{(i)}, y^{(i)}) \qquad (1.38)$$

$$Cost(h_\theta(x^{(i)}), y^{(i)}) = \frac{1}{2} [h_\theta(x^{(i)} - y^{(i)}]^2 \qquad (1.39)$$

Convex function

The above graph is a convex.



The below graph is a non-convex.
If we use the cost function of linear regression in logistic regression, the the we would get non-convex cost function, because the **hypothesis is** $\frac{1}{1+e^{-\theta^\top x}}$.
For Logistic regression,

$$Cost(h_\theta(x^{(i)}), y^{(i)}) = \begin{cases} -\log(1 - h_\theta(x)), & if \ y = 0 \\ -\log(h_\theta(x)), & if \ y = 1. \end{cases} \qquad (1.40)$$

If y=1,



If y=0,

**Simplified Cost function,**

$$Cost(h_\theta(x^{(i)}), y^{(i)}) = -(1 - y)\log(1 - h_\theta(x)) - y\log(h_\theta(x)). \ \forall \ y \ \epsilon \ \{0, 1\}. \tag{1.41}$$

$$If \ y = 1 : Cost(h_\theta(x), y) = -\log(h_\theta(x)). \tag{1.42}$$

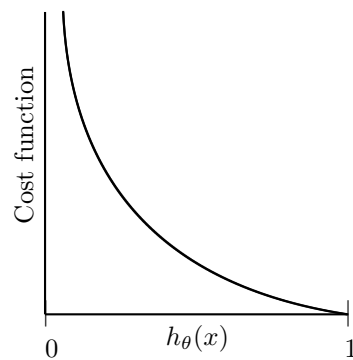$$If \ y = 0 : Cost(h_\theta(x), y) = -\log(1 - h_\theta(x)). \tag{1.43}$$

$$Cost function = J(\theta) = \frac{1}{m}\sum_{i=1}^{m} Cost(h_\theta(x^{(i)}, y^{(i)}) \tag{1.44}$$

$$= \frac{1}{m}[-\sum_{i=1}^{m}(1 - y^{(i)})\log(1 - h_\theta(x^{(i)})) + y^{(i)}\log(h_\theta(x^{(i)}))] \tag{1.45}$$

### 1.1.12 Gradient Descent for Logistic Regression

$$\text{Repeat until convergence (minimum)}\Big\{\theta_j := \theta_j - \alpha\frac{\partial}{\partial\theta_j}J(\theta) \tag{1.46}$$

$$\frac{\partial}{\partial\theta_j}J(\theta) = \frac{1}{m}\sum_{i=1}^{m}[h_\theta(x^{(i)}) - y^{(i)}]x_j^{(i)} \tag{1.47}$$

### 1.1.13 Optimization algorithm

1. Gradient descent.

2. Conjugate gradient.

3. BFGS.

4. L - BFGS.

These are the 4 algorithms to minimise **cost function.**
Advantages of the **last three advanced optimization algorithm.**

- No need to manually pick $\alpha$.

- Often faster than Gradient descent.

- They themselves choose $\alpha$, for faster convergence.

### 1.1.14 Multiclass Classification : One-vs-All

$$y \in \{0, 1...n\}$$
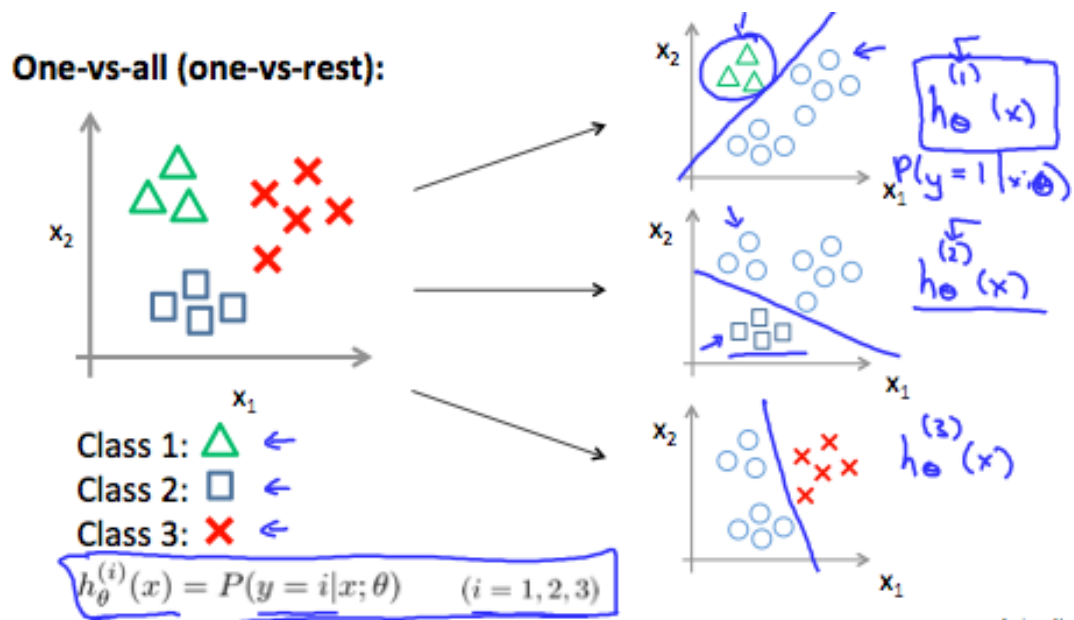$$h_\theta^{(0)}(x) = P(y = 0|x; \theta)$$
$$h_\theta^{(1)}(x) = P(y = 1|x; \theta)$$
$$\vdots$$
$$h_\theta^{(n)}(x) = P(y = n|x; \theta)$$
$$\text{prediction} = \max_i(h_\theta^{(i)}(x))$$

To summarize,



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

2. To make a prediction on a new x, pick the class that maximizes $h_\theta(x)$

15

### 1.1.15 Problem of Overfitting,

## Example: Linear regression (housing prices)



$$\rightarrow \theta_0 + \theta_1 x$$
"Underfit"   "High bias"

$$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$$
"Just right"

$$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$
"Overfit"   "High variance"

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

Similar for Logistic Regression.

**Addressing Overfitting,**

1. Reduce number of features.

   - Manually select which features to keep.
   - Model selection algorithm.

2. Regularization

   - Keep all features, but reduce magnitude/values of parameters $_j$.
   - Works well when we have a lot of features, each of which contributes a bit to predict $y^{(i)}$.

### 1.1.16 Regularization

Small values for parameters $\theta_0, \theta_1, \theta_2, ..., \theta_n$.

1. Simpler hypothesis.

2. Less prone to overfitting.

**Regularized Cost function,**

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} \left[ [h_\theta(x^{(i)} - y^{(i)}]^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right] \qquad (1.48)$$

If '$\lambda$' is extremely large, then the cost function will become underfitting (doesn't fit to our training data).

Repeat {

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

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

}

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

$$\theta_j := \theta_j(1 - \alpha\frac{\lambda}{m}) - \alpha\frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)} \qquad\qquad (1.49)$$

The first term in the above equation, $1 - \alpha\frac{\lambda}{m}$ will always be less than 1. Intuitively you can see it as reducing the value of $\theta_j$ by some amount on every update. Notice that the second term is now exactly the same as it was before.
**Normal equation after regularization,**

$$\theta = (X^\top X) + \lambda.L^{-1}X^\top Y \qquad\qquad (1.50)$$
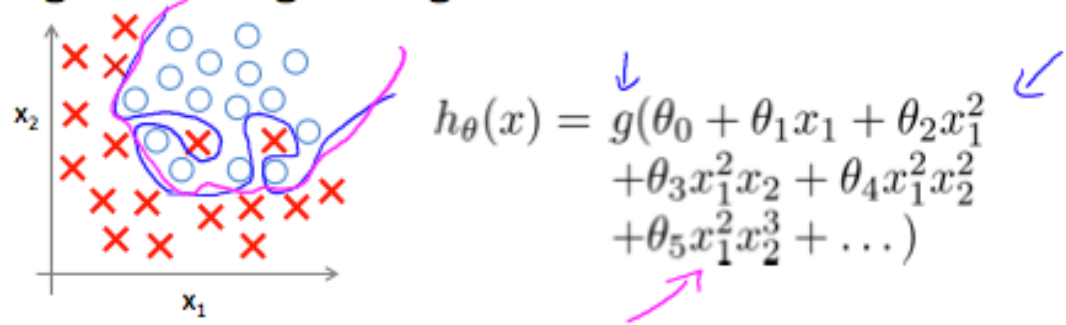
$$where \; L = \begin{bmatrix} 0 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & \ddots & \\ & & & & 1 \end{bmatrix}$$

### 1.1.17 Regularized Logistic Regression

**Cost function,**

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

$$(1.51)$$

**Regularized logistic regression.**



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

The second sum, $\sum_{j=1}^{n} \theta_j^2$ means to explicitly exclude the bias term, $\theta_0$. I.e. the vector is indexed from 0 to n (holding n+1 values, $\theta_0$ through $\theta_n$), and this sum explicitly skips $\theta_0$, by running from 1 to n, skipping 0. Thus, when computing the equation, we should continuously update the two following equations:

Repeat {

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

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

}

## 1.2 Neural Networks

### 1.2.1 Model representation I

Let's examine how we will represent a hypothesis function using neural networks. At a very simple level, neurons are basically computational units that take inputs (dendrites) as electrical inputs (called **"spikes"**) that are channeled to outputs (axons). In our model, our dendrites are like the input features $x_1 \cdots x_n$, and the output is the result of our hypothesis function. In this model our $x_0$ input node is sometimes called the **"bias unit"**. It is always equal to 1. In neural networks, we use the same logistic function as in classification, $\frac{1}{1+e^{-\theta^T x}}$, yet we sometimes call it a sigmoid (logistic) activation function. In this situation, our **"theta" parameters** are sometimes called **"weights".**
A simple representation looks like :

$$\begin{bmatrix} x_0 x_1 x_2 \end{bmatrix} \rightarrow [\,] \rightarrow h_\theta(x)$$

Our input nodes (layer 1), also known as the **"input layer"**, go into another node (layer 2), which finally outputs the hypothesis function, known as the **"output layer".**

We can have intermediate layers of nodes between the input and output layers called the **"hidden layers."**

In this example, we label these intermediate or **"hidden"** layer nodes $a_0^2, \cdots, a_n^2 a$ and call them **"activation units."**

1. $a_i^{(j)}$ = "activation" of unit i in layer j

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

If we had one hidden layer, it would look like :

$$\begin{bmatrix} x_0 x_1 x_2 x_3 \end{bmatrix} \rightarrow \begin{bmatrix} a_1^{(2)} a_2^{(2)} a_3^{(2)} \end{bmatrix} \rightarrow h_\theta(x)$$

The values for each of the **"activation"** nodes is obtained as follows :

$$a_1^{(2)} = g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3) \qquad (1.52)$$

$$a_2^{(2)} = g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3) \qquad (1.53)$$

$$a_3^{(2)} = g(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3) \qquad (1.54)$$

$$h_\Theta(x) = a_1^{(3)} = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)}) \qquad (1.55)$$

This is saying that we compute our activation nodes by using a 3×4 matrix of parameters. We apply each row of the parameters to our inputs to obtain the value for one activation node. Our hypothesis output is the logistic function applied to the sum of the values of our activation nodes, which have been multiplied by yet another parameter matrix $\Theta^{(2)}$ containing the weights for our second layer of nodes.

Each layer gets its own matrix of weights, $\Theta^{(j)}$.

The dimensions of these matrices of weights is determined as follows :

If network has $s_j$ units in layer $j$ and $s_{j+1}$ units in layer $j + 1$, then $\Theta^{(j)}$ will be of dimension $s_{j+1} \times (s_j + 1)$.

The +1 comes from the addition in $\Theta^{(j)}$ of the **"bias nodes,"** $x_0$ and $\Theta_0^{(j)}$. In other words the output nodes will not include the bias nodes while the inputs will.

## 1.2.2   Model representation II

we'll do a vectorized implementation of the above functions. We're going to define a new variable $z_k^{(j)}$ that encompasses the parameters inside our g function.

$$a_1^{(2)} = g(z_1^{(2)}) \tag{1.56}$$

$$a_2^{(2)} = g(z_2^{(2)}) \tag{1.57}$$

$$a_3^{(2)} = g(z_3^{(2)}) \tag{1.58}$$

$$x = a^{(1)} \tag{1.59}$$

$$a^{(2)} = g(\theta^{(1)}x) \tag{1.60}$$

$$= g(\theta^{(1)}a^{(1)}) = g(z^{(2)}) \tag{1.61}$$

$$h_\theta(x) = g(\theta^{(2)}a^{(2)}) \tag{1.62}$$

$$= g(z^{(3)}) \tag{1.63}$$

In other words, for layer j=2 and node k, the variable z will be :

$$z_k^{(2)} = \Theta_{k,0}^{(1)}x_0 + \Theta_{k,1}^{(1)}x_1 + \cdots + \Theta_{k,n}^{(1)}x_n \tag{1.64}$$

The vector representation of x and $z^j$ is :

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \qquad\qquad z^{(j)} = \begin{bmatrix} z_1^{(j)} \\ z_2^{(j)} \\ \vdots \\ z_n^{(j)} \end{bmatrix}$$

Setting x $= a^{(1)}$, we can rewrite the equation as :

$$z^{(j)} = \theta^{(j-1)}a^{(j-1)} \tag{1.65}$$

We are multiplying our matrix $\Theta^{(j-1)}$ with dimensions $s_j\times$ (n+1)(where $s_j$ is the number of our activation nodes) by our vector $a^{(j-1)}$ with height (n+1). This gives us our vector $z^{(j)}$ with height $s_j$. Now we can get a vector of our activation nodes for layer j as follows :

$$a^{(j)} = g(z^{(j)})a \tag{1.66}$$

Where our function g can be applied element-wise to our vector $z^{(j)}$
We can then add a bias unit (equal to 1) to layer j after we have computed $a^{(j)}$. This will be element $a_0^{(j)}$ and will be equal to 1. To compute our final hypothesis, let's first compute another z vector :

$$z^{(j+1)} = \Theta^{(j)}a^{(j)} \tag{1.67}$$

We get this final z vector by multiplying the next theta matrix after $\Theta^{(j-1)}$ with the values of all the activation nodes we just got. This last theta matrix $\Theta^{(j)}$

will have only one row which is multiplied by one column $a^{(j)}$ so that our result is a single number. We then get our final result with :

$$h_\Theta(x) = a^{(j+1)} = g(z^{(j+1)}) \tag{1.68}$$

Notice that in this last step, between layer j and layer j+1, we are doing exactly the same thing as we did in logistic regression. Adding all these intermediate layers in neural networks allows us to more elegantly produce interesting and more complex non-linear hypothesis.

**Examples,**

A simple example of applying neural networks is by predicting $x_1$ AND $x_2$, which is the logical 'and' operator and is only true if both $x_1$ and $x_2$ are 1.

$$h_\Theta(x) = g(-30 + 20x_1 + 20x_2)$$
$$x_1 = 0 \quad and \quad x_2 = 0 \quad then \quad g(-30) \approx 0$$
$$x_1 = 0 \quad and \quad x_2 = 1 \quad then \quad g(-10) \approx 0$$
$$x_1 = 1 \quad and \quad x_2 = 0 \quad then \quad g(-10) \approx 0$$
$$x_1 = 1 \quad and \quad x_2 = 1 \quad then \quad g(10) \approx 1$$

So we have constructed one of the fundamental operations in computers by using a small neural network rather than using an actual AND gate. Neural networks can also be used to simulate all the other logical gates. The following is an example of the logical operator 'OR', meaning either $x_1$ is true or $x_2$ is true, or both :

## Example: OR function



| $x_1$ | $x_2$ | $h_\Theta(x)$ |
|-------|-------|---------------|
| 0 | 0 | $g(-10) \approx 0$ |
| 0 | 1 | $g(10) \approx 1$ |
| 1 | 0 | $\approx 1$ |
| 1 | 1 | $\approx 1$ |

$g(-10 + 20 x_1 + 20 x_2)$

The $^{(1)}$ matrices for AND, NOR, and OR are :

$$AND :$$
$$\Theta^{(1)} = \begin{bmatrix} -30 & 20 & 20 \end{bmatrix}$$
$$NOR :$$
$$\Theta^{(1)} = \begin{bmatrix} 10 & -20 & -20 \end{bmatrix}$$
$$OR :$$
$$\Theta^{(1)} = \begin{bmatrix} -10 & 20 & 20 \end{bmatrix}$$

We can combine these to get the XNOR logical operator (which gives 1 if $x_1$ and $x_2$ are both 0 or both 1).

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \end{bmatrix} \rightarrow \begin{bmatrix} a^{(3)} \end{bmatrix} \rightarrow h_\Theta(x)$$

Let's write out the values for all our nodes :

$$a^{(2)} = g(\Theta^{(1)} \cdot x)$$
$$a^{(3)} = g(\Theta^{(2)} \cdot a^{(2)})$$
$$h_\Theta(x) = a^{(3)}$$

### 1.2.3  Multiclass Classification

To classify data into multiple classes, we let our hypothesis function return a vector of values. Say we wanted to classify our data into one of four categories. We will use the following example to see how this classification is done. This algorithm takes as input an image and classifies it accordingly :



Andrew Ng

We can define our set of resulting classes as y :

$$
y^{(i)} =
\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix},
\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix},
\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix},
\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix},
$$

Each $y^{(i)}$ 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 :
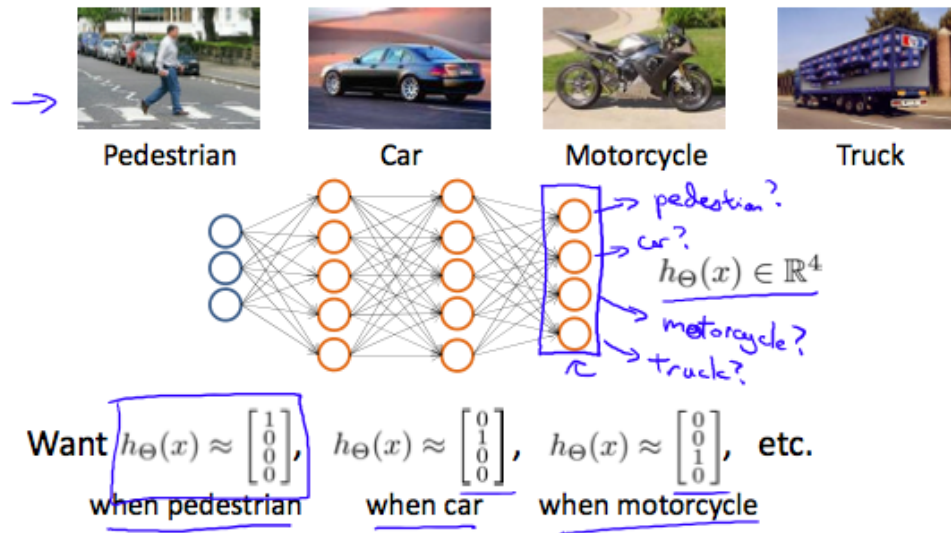
$$
\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \cdots \\ x_n \end{bmatrix} \rightarrow
\begin{bmatrix} a_0^{(2)} \\ a_1^{(2)} \\ a_2^{(2)} \\ \cdots \end{bmatrix} \rightarrow
\begin{bmatrix} a_0^{(3)} \\ a_1^{(3)} \\ a_2^{(3)} \\ \cdots \end{bmatrix} \rightarrow \cdots \rightarrow
\begin{bmatrix} h_\Theta(x)_1 \\ h_\Theta(x)_2 \\ h_\Theta(x)_3 \\ h_\Theta(x)_4 \end{bmatrix}
$$

Our resulting hypothesis for one set of inputs may look like :

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

In which case our resulting class is the third one down, or $h_\Theta(x)_3$, which represents the motorcycle.

### 1.2.4 Cost function for Neural networks

Let's declare some variables.

1. L = total number of layers in the network.

2. $s_l$ = number of units (not counting bias unit) in layer l.

3. K = number of output units/classes.

**The cost function for regularized logistic regression,**

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

$$(1.69)$$

**The cost function for Neural networks,**

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

$$(1.70)$$

In the regularization part, after the square brackets, we must account for multiple theta matrices. The number of columns in our current theta matrix is equal to the number of nodes in our current layer (including the bias unit). The number of rows in our current theta matrix is equal to the number of nodes in the next layer (excluding the bias unit). As before with logistic regression, we square every term.

1. The double sum simply adds up the logistic regression costs calculated for each cell in the output layer.

2. The triple sum simply adds up the squares of all the individual s in the entire network.

3. The i in the triple sum does not refer to training example i.

### 1.2.5 Backpropagation Algorithm

**"Backpropagation"** is neural-network terminology for minimizing our cost function, just like what we were doing with gradient descent in logistic and linear regression. Our goal is to compute :

$$\min_{\Theta} J(\Theta) \tag{1.71}$$

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

$$\frac{\partial}{\partial \Theta_{j,i}^{(l)}} J(\Theta) \tag{1.72}$$

**Backpropagation algorithm,**

**Backpropagation algorithm**

Training set $\{(x^{(1)}, y^{(1)}), \ldots, (x^{(m)}, y^{(m)})\}$

Set $\triangle_{ij}^{(l)} = 0$ (for all $l, i, j$).  $\quad$ (used to compute $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$)

For $i = 1$ to $m \leftarrow \quad (x^{(i)}, y^{(i)})$.

$\quad$ Set $a^{(1)} = x^{(i)}$

$\quad$ Perform forward propagation to compute $a^{(l)}$ for $l = 2, 3, \ldots, L$

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

$\quad$ Compute $\delta^{(L-1)}, \delta^{(L-2)}, \ldots, \delta^{(2)} \quad \cancel{\delta^{(1)}}$

$\quad$ $\triangle_{ij}^{(l)} := \triangle_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)} \leftarrow \qquad \triangle^{(l)} := \triangle^{(l)} + \delta^{(l+1)} (a^{(l)})^T.$

$D_{ij}^{(l)} := \frac{1}{m} \triangle_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)}$ if $j \neq 0 \qquad\qquad \frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$

$D_{ij}^{(l)} := \frac{1}{m} \triangle_{ij}^{(l)} \qquad\qquad$ if $j = 0$

1. Given training set $\{(x^{(1)}, y^{(1)}) \cdots (x^{(m)}, y^{(m)})\}$.

   - Set $\Delta_{i,j}^{(l)} := 0$ for all (l,i,j), (hence you end up having a matrix full of zeros)

2. For training example t = 1 to m :

   - Set $a^{(1)} := x^{(t)}$.
   - Perform forward propagation to compute $a^{(l)}$ for l=2,3,...,L.
   - 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
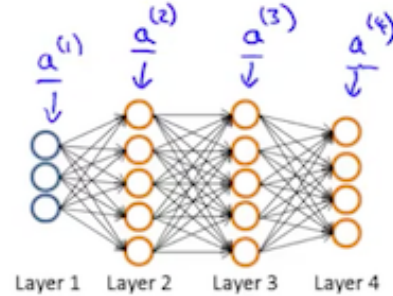
the last layer. So our **"error values"** for the last layer are simply the differences of our actual results in the last layer and the correct outputs in y. To get the delta values of the layers before the last layer, we can use an equation that steps us back from right to left :

## Gradient computation

Given one training example $(x, y)$:

Forward propagation:

$$a^{(1)} = x$$
$$z^{(2)} = \Theta^{(1)}a^{(1)}$$
$$a^{(2)} = g(z^{(2)}) \quad (\text{add } a_0^{(2)})$$
$$z^{(3)} = \Theta^{(2)}a^{(2)}$$
$$a^{(3)} = g(z^{(3)}) \quad (\text{add } a_0^{(3)})$$
$$z^{(4)} = \Theta^{(3)}a^{(3)}$$
$$a^{(4)} = h_\Theta(x) = g(z^{(4)})$$



Layer 1   Layer 2   Layer 3   Layer 4

- Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$ using $\delta^{(l)} = ((\Theta^{(l)})^T \delta^{(l+1)}) \mathbin{.*} a^{(l)} \mathbin{.*} (1 - a^{(l)})$.
  The delta values of layer l are calculated by multiplying the delta values in the next layer with the theta matrix of layer l. We then element-wise multiply 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 :

$$g'(z^{(l)}) = a^{(l)} \mathbin{.*} (1 - a^{(l)}) \tag{1.73}$$

- $\Delta_{i,j}^{(l)} := \Delta_{i,j}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$ (or) with vectorization, $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)}(a^{(l)})^T$.
  Hence we update our new $\Delta$ matrix.

  (a) $D_{i,j}^{(l)} := \frac{1}{m}(\Delta_{i,j}^{(l)} + \lambda\Theta_{i,j}^{(l)})$, if j$\neq$ 0.
  $\phantom{(a)} D_{i,j}^{(l)} := \frac{1}{m}(\Delta_{i,j}^{(l)}$, if j = 0.

  The capital-delta matrix D is used as an "accumulator" to add up our values as we go along and eventually compute our partial derivative. Thus we get $\frac{\partial}{\partial\Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$.
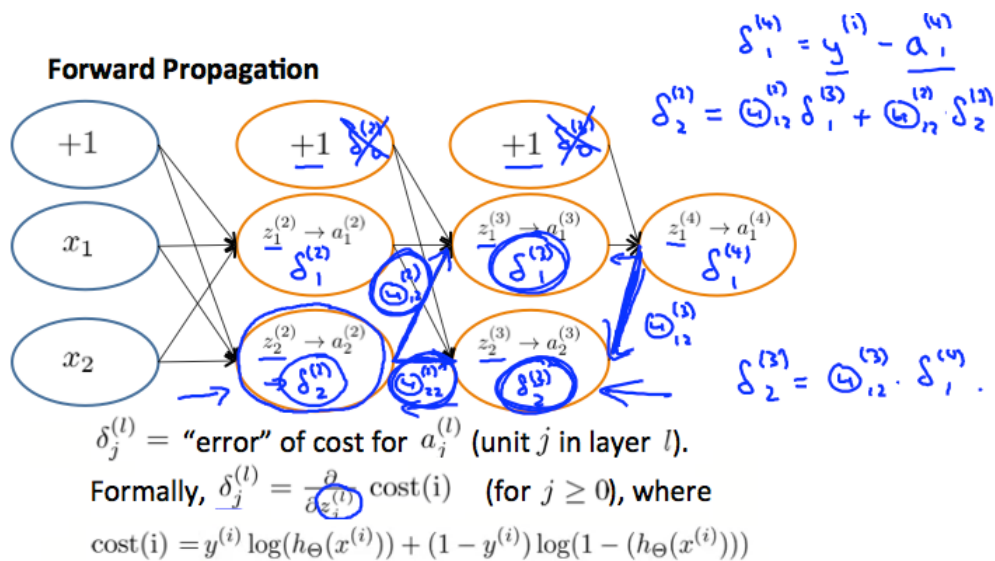
If we consider simple non-multiclass classification (k = 1) and disregard regularization, the cost is computed with :

$$Cost(t) = (1 - y^{(t)}) \log(1 - h_\theta(x^{(t)})) + y^{(t)} \log(h_\theta(x^{(t)})) \tag{1.74}$$

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

$$\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} cost(t) \tag{1.75}$$

Our derivative is the slope of a line tangent to the cost function, so the steeper the slope the more incorrect we are. Let us consider the following neural network below and see how we could calculate some $\delta_j^{(l)}$ :



**Forward Propagation**

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

$\delta_2^{(1)} = \Theta_{12}^{(i)} \delta_1^{(3)} + \Theta_{n}^{(i)} \delta_2^{(1)}$

$\delta_2^{(3)} = \Theta_{12}^{(3)} \cdot \delta_1^{(4)}$

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

Formally, $\delta_j^{(l)} = \frac{\partial}{\partial z_i^{(l)}} cost(i)$ (for $j \geq 0$), where

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

Andrew Ng

28

In the image above, to calculate $\delta_2^{(2)}$, we multiply the weights $\Theta_{12}^{(2)}$ and $\Theta_{22}^{(2)}$ by their respective $\delta$ values found to the right of each edge. So we get $\delta_2^{(2)} = \Theta_{12}^{(2)} * \delta_1^{(3)} + \Theta_{22}^{(2)} * \delta_2^{(3)}$. To calculate every single possible $\delta_j^{(l)}$, we could start from the right of our diagram. We can think of our edges as our $\Theta_{ij}$. Going from right to left, to calculate the value of $\delta_j^{(l)}$, you can just take the over all sum of each weight times the $\delta$ it is coming from. Hence, another example would be $\delta_2^{(3)} = \Theta_{12}^{(3)} * \delta_1^{(4)}$.

### 1.2.6   Gradient Checking

Gradient checking will assure that our backpropagation works as intended. We can approximate the derivative of our cost function with :

$$\frac{\partial}{\partial\theta}J(\Theta) \approx \frac{J(\Theta + \epsilon) - J(\Theta - \epsilon)}{2\epsilon} \tag{1.76}$$

With multiple theta matrices, we can approximate the derivative with respect to $\Theta_j$ as follows :

$$\frac{\partial}{\partial\theta_j}J(\Theta) \approx \frac{J(\Theta_1, ..., \Theta_j + \epsilon, ..., \Theta_n) - J(\Theta_1, ..., \Theta_j - \epsilon, ..., \Theta_n)}{2\epsilon} \tag{1.77}$$

A small value for $\epsilon$ (epsilon) such as $\epsilon = 10^{-4}$, guarantees that the math works out properly. If the value for $\epsilon$ is too small, we can end up with numerical problems.
We previously saw how to calculate the deltaVector. So once we compute our gradApprox vector, we can check that gradApprox $\approx$ deltaVector.
Once you have verified once that your backpropagation algorithm is correct, you don't need to compute gradApprox again. The code to compute gradApprox can be very slow.

### 1.2.7   Random Intialization

Initializing all theta weights to zero does not work with neural networks. When we backpropagate, all nodes will update to the same value repeatedly. Instead we can randomly initialize our weights for our $\Theta$ matrices using the following method :

**Random initialization: Symmetry breaking**

→ Initialize each $\Theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$
(i.e. $-\epsilon \le \Theta_{ij}^{(l)} \le \epsilon$ )

E.g.  _random_ $10 \times 11$ _matrix_ _(betw. 0_
_and 1)_

→ `Theta1 = rand(10,11)*(2*INIT_EPSILON)`
            `- INIT_EPSILON;`          $[-\epsilon, \epsilon]$

→ `Theta2 = rand(1,11)*(2*INIT_EPSILON)`
            `- INIT_EPSILON;`

Hence, we initialize each $\Theta_{ij}^{(l)}$ to a random value between $[-\epsilon, \epsilon]$. Using the above formula guarantees that we get the desired bound.

> **The epsilon used above is unrelated to the epsilon from Gradient Checking.**

### 1.2.8   Choosing Neural network

1. Pick a network architecture.

2. Choose the layout of your neural network.

3. Including how many hidden units in each layer and how many layers in total you want to have.

   - Number of input units = dimension of features $x^{(i)}$.

   - Number of output units = number of classes.

   - Number of hidden units per layer = usually more the better (must balance with cost of computation as it increases with more hidden units).

   - **Defaults** : 1 hidden layer. If you have more than 1 hidden layer, then it is recommended that you have the same number of units in every hidden layer.

### 1.2.9   Training a Neural Network

1. Randomly initialize the weights.

2. Implement forward propagation to get $h_\Theta(x^{(i)})$.

3. Implement the cost function.

4. Implement backpropagation to compute partial derivatives.

5. Use gradient checking to confirm that your backpropagation works. Then disable gradient checking.

6. Use gradient descent or a built-in optimization function to minimize the cost function with the weights in theta.

**Ideally, you want $h_\Theta(x^{(i)}) \approx y^{(i)}$. This will minimize our cost function. However,**

**keep in mind that $J(\Theta)$ is not convex and thus we can end up in a local minimum instead.**

# Chapter 2

# Deciding whether the algorithm is perfect or not

## 2.1  Evaluating the algorithm

### 2.1.1  Evaluating a Hypothesis

**Fails to generalize to new examples not in training set.**
**Once we have done some trouble shooting for errors in our predictions by :**

1. Getting more training examples.

2. Trying smaller sets of features.

3. Trying additional features.

4. Trying polynomial features.

5. Increasing or decreasing $\lambda$.

We can move on to evaluate our new hypothesis.
A hypothesis may have a low error for the training examples but still be inaccurate (because of overfitting). Thus, to evaluate a hypothesis, given a dataset of training examples, we can split up the data into two sets: a training set and a test set. Typically, the training set consists of 70% of your data and the test set is the remaining 30%.
The new procedure using these two sets is then :

**Dataset,**
**Model selection,**

| Size | Price | |
|---|---|---|
| 2104 | 400 | $(\mathrm{x}^{(1)}, y^{(1)})$ |
| 1600 | 330 | $(\mathrm{x}^{(2)}, y^{(1)})$ |
| 2400 | 369 | $(\mathrm{x}^{(3)}, y^{(1)})$ |
| 1416 | 232 | $\rightarrow$ $\vdots$ |
| 3000 | 540 | $\vdots$ |
| 1985 | 300 | $\vdots$ |
| 1534 | 315 | $(\mathrm{x}^{(m)}, y^{(m)})$ |
| 1427 | 199 | $(\mathrm{x}^{(1)}_{test}, y^{(1)}_{test})$ |
| 1380 | 212 | $\longmapsto$ $(\mathrm{x}^{(2)}_{test}, y^{(2)}_{test})$ |
| 1494 | 243 | $(\mathrm{x}^{(m_{test})}_{test}, y^{(m_{test})}_{test})$ |

Table 2.1: Evaluating your hypothesis

**Model selection**

$\rightarrow \boxed{d = \text{degree of polynomial}}$ $\downarrow$

$d=1$ 1. $\rightarrow h_\theta(x) = \theta_0 + \theta_1 x \longrightarrow \Theta^{(1)} \longrightarrow J_{test}(\Theta^{(1)})$

$d=2$ 2. $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \longrightarrow \Theta^{(2)} \longrightarrow J_{test}(\Theta^{(2)})$

$d=3$ 3. $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_3 x^3 \longrightarrow \Theta^{(3)} \longrightarrow J_{test}(\Theta^{(3)})$

$\vdots$ $\vdots$ $\vdots$

$d=10$ 10. $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_{10} x^{10} \longrightarrow \Theta^{(10)} \longrightarrow J_{test}(\Theta^{(10)})$

Choose $\boxed{\theta_0 + \ldots \theta_5 x^5} \leftarrow$

How well does the model generalize? Report test set $\boxed{\Theta_0, \Theta_1 \ldots}$
error $J_{test}(\theta^{(5)})$. $\qquad \Theta^{(5)}$

Problem: $J_{test}(\theta^{(5)})$ is likely to be an optimistic estimate of generalization error. I.e. our extra parameter ($\underline{d}$ = degree of polynomial) is fit to test set.

1. $m_{test}$ = number of test example.

2. Learn $\Theta$ and minimize $J_{train}(\Theta)$ using the training set.

3. Compute the test set error $J_{test}(\Theta)$.

### 2.1.2 The test set error

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

2. For classification   Misclassification error (aka 0/1 misclassification error):

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

$$(2.1)$$

This gives us a binary 0 or 1 error result based on a misclassification. The average test error for the test set is :

$$\text{Test Error} = \frac{1}{m_{test}} \sum_{i=1}^{m_{test}} err(h_\Theta(x_{test}^{(i)}), y_{test}^{(i)}) \tag{2.2}$$

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

## Train/validation/test error

### Training error:

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

### Cross Validation error:

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

### Test error:

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

### 2.1.3 Model Selection and Train/ Validation/ Test Sets

Just because a learning algorithm fits a training set well, that does not mean it is a good hypothesis. It could over fit and as a result your predictions on the test set would be poor. The error of your hypothesis as measured on the data set with which you trained the parameters will be lower than the error on any

other data set.

Given many models with different polynomial degrees, we can use a systematic approach to identify the 'best' function. In order to choose the model of your hypothesis, you can test each degree of polynomial and look at the error result. One way to break down our dataset into the three sets is :

1. Training set: 60%.

2. Cross validation set: 20%.

3. Test set: 20%.

**Evaluating your hypothesis**
Dataset:

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

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

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

$m_{cv}$ = no. of cv examples $(x_{cv}^{(i)}, y_{cv}^{(i)})$

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

$m_{test}$

Training set

Cross validation set (CV)

test set

We can now calculate three separate error values for the three different sets using the following method:

1. Optimize the parameters in   using the training set for each polynomial degree.

2. Find the polynomial degree d with the least error using the cross validation set.

3. Estimate the generalization error using the test set with $J_{test}(\Theta^{(d)})$, (d = theta from polynomial with lower error).

**This way, the degree of the polynomial d has not been trained using the test set.**
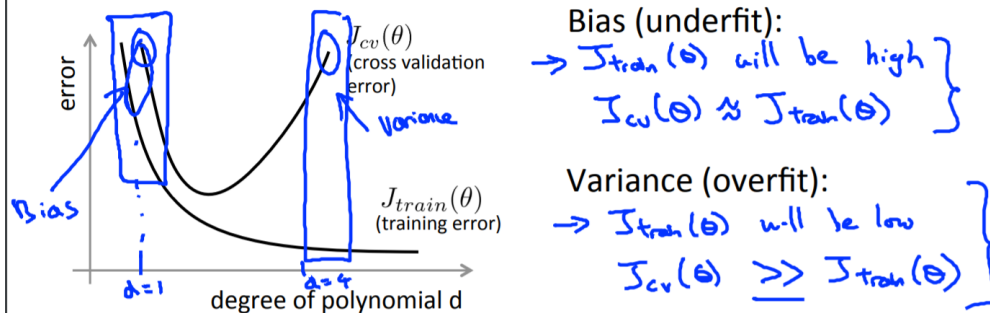
### 2.1.4   Diagnosing Bias vs Variance

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 predictions.

- High bias is underfitting and high variance is overfitting. Ideally, we need to find a golden mean between these two.



The training error will tend to **decrease** as we increase the degree d of the polynomial.
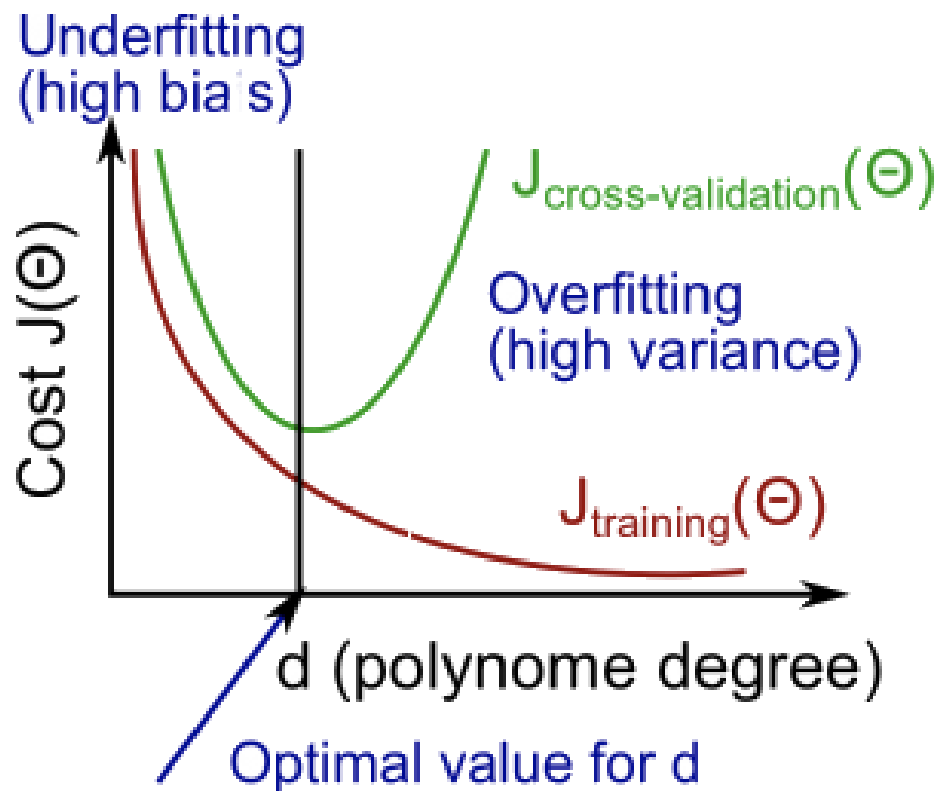
At the same time, the cross validation error will tend to **decrease** as we increase d up to a point, and then it will **increase** as **d** is increased, forming a convex curve.

**High bias (underfitting)** : both $J_{train}(\Theta)$ and $J_{CV}(\Theta)$ will be high. Also, $J_{CV}(\Theta) \approx J_{train}(\Theta)$.

**High variance (overfitting)** : $J_{train}(\Theta)$ will be low and $J_{CV}(\Theta)$ will be much greater than $J_{train}(\Theta)$.
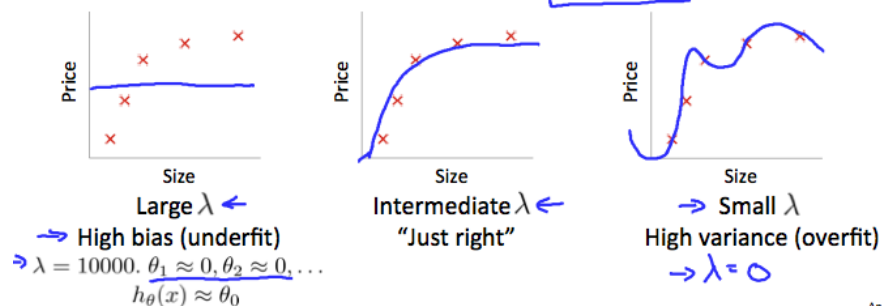
## 2.1.5   Regularization

In the figure above, we see that as $\lambda$ increases, our fit becomes more rigid. On the other hand, as $\lambda$ approaches 0, we tend to over overfit the data. So how do we choose our parameter $\lambda$ to get it 'just right' ? In order to choose the model and the regularization term $\lambda$, we need to:

Underfitting
(high bias)

Cost J(Θ)

$J_{\text{cross-validation}}(\Theta)$

Overfitting
(high variance)

$J_{\text{training}}(\Theta)$

d (polynome degree)

Optimal value for d

**Linear regression with regularization**

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

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



Price — Size
Large $\lambda$
High bias (underfit)
$\lambda = 10000.\ \theta_1 \approx 0, \theta_2 \approx 0, \ldots$
$h_\theta(x) \approx \theta_0$

Price — Size
Intermediate $\lambda$
"Just right"

Price — Size
Small $\lambda$
High variance (overfit)
$\lambda = 0$

1. Create a list of lambdas (i.e. $\lambda$ 0,0.01,0.02,0.04,0.08,0.16,0.32,0.64,1.28,2.56,5.12,10.24).

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

37

3. Iterate through the $\lambda$s and for each $\lambda$ go through all the models to learn some $\Theta$.

4. Compute the cross validation error using the learned $\Theta$ (computed with $\lambda$) on the $J_{CV}(\Theta)$ without regularization or $\lambda = 0$.

5. Select the best combo that produces the lowest error on the cross validation set.

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

## Choosing the regularization parameter $\lambda$

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \quad \leftarrow$$

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

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

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

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

### 2.1.6 Learning Curves

Training an algorithm on a very few number of data points (such as 1, 2 or 3) will easily have 0 errors because we can always find a quadratic curve that touches exactly those number of points. Hence :

1. As the training set gets larger, the error for a quadratic function increases.

2. The error value will plateau out after a certain m, or training set size.
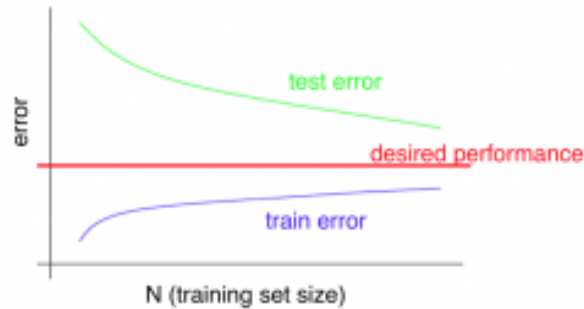
**Experiencing high bias :**
**Low training set size:** causes $J_{train}(\Theta)$ to be low and $J_{CV}(\Theta)$ to be high.
**Large training set size:** causes both $J_{train}(\Theta)$ and $J_{CV}(\Theta)$ to be high with $J_{train}(\Theta) \approx J_{CV}(\Theta)$.
If a learning algorithm is suffering from **high bias**, getting more training data

## More on Bias vs. Variance

Typical learning curve for high variance(at fixed model complexity):



will not (by itself) help much.

**Experiencing high variance:**

**Low training set size:** $J_{train}(\Theta)$ will be low and $J_{CV}(\Theta)$ will be high.

**Large training set size:** $J_{train}(\Theta)$ increases with training set size and $J_{CV}(\Theta)$ continues to decrease without leveling off. Also, $J_{train}(\Theta) < J_{CV}(\Theta)$ but the difference between them remains significant.

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

## More on Bias vs. Variance

Typical learning curve for high bias(at fixed model complexity):



### 2.1.7   Deciding What to Do Next Revisited

Our decision process can be broken down as follows:

1. **Getting more training examples:** Fixes high variance.

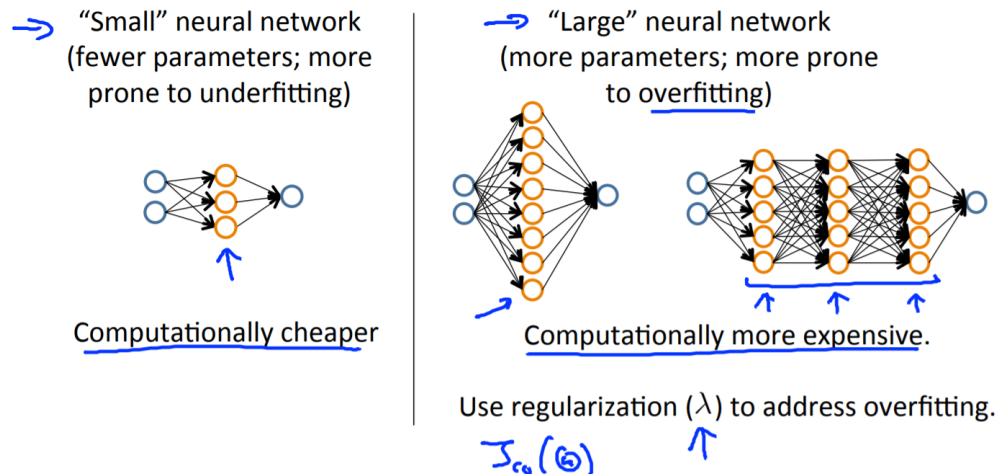2. **Trying smaller sets of features:** Fixes high variance.

3. **Adding features:** Fixes high bias.

4. **Adding polynomial features:** Fixes high bias.

5. **Decreasing $\lambda$:** Fixes high bias.

6. **Increasing $\lambda$:** Fixes high variance.

**Diagnosing Neural Networks**

1. A neural network with fewer parameters is **prone to underfitting**. It is also **computationally cheaper.**

2. A large neural network with more parameters is **prone to overfitting.** It is also **computationally expensive.** In this case you can use regularization (increase $\lambda$) to address the overfitting.

Using a single hidden layer is a good starting default. You can train your neural network on a number of hidden layers using your cross validation set. You can then select the one that performs best.

## Neural networks and overfitting



**Machine learning diagnostic**

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

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

**Model Complexity Effects:**

1. Lower-order polynomials (low model complexity) have high bias and low variance. In this case, the model fits poorly consistently.

2. 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.

3. In reality, we would want to choose a model somewhere in between, that can generalize well but also fits the data reasonably well.

### 2.1.8   Building a spam classifier