# Deep Learning: A first look

Exebit 2018, IITM

Sachin Sridhar, Department of CSE

April 15, 2018

Indian Institute of Technology Madras

# Table of contents

1

# Machine Learning: Some basics

- *"The science of getting computers to act without being explicitly programmed."* - Andrew Ng

- *"The science of getting computers to act without being explicitly programmed."* - Andrew Ng
- Uses statistical procedures to get computer systems to be able to "learn" and "predict" an output given an input

- *"The science of getting computers to act without being explicitly programmed."* - Andrew Ng
- Uses statistical procedures to get computer systems to be able to "learn" and "predict" an output given an input
- Machine learning algorithms are different from conventional algorithms; they are statistical in nature and often have no correctness guarantees

- *"The science of getting computers to act without being explicitly programmed."* - Andrew Ng
- Uses statistical procedures to get computer systems to be able to "learn" and "predict" an output given an input
- Machine learning algorithms are different from conventional algorithms; they are statistical in nature and often have no correctness guarantees
- Seen in daily life!

- Predicting an output given an input is mathematically formulated as a probability query $P(Y = y | X = x)$.

- Predicting an output given an input is mathematically formulated as a probability query $P(Y = y | X = x)$.
- We need a measure of "how good" our estimate is. So we try to minimize a *loss function* (mean squared error, negative log-likelihood, etc)

- Predicting an output given an input is mathematically formulated as a probability query $P(Y = y | X = x)$.
- We need a measure of "how good" our estimate is. So we try to minimize a *loss function* (mean squared error, negative log-likelihood, etc)
- This can be optimized using techniques from calculus, convex analysis, etc. A popular example of such a technique is gradient descent.

Consider the following sequence of coin tosses:

<div align="center">

H T T H H H T T T T

</div>

Consider the following sequence of coin tosses:

H T T H H H T T T T

What could the bias of the coin be?

Consider the following sequence of coin tosses:

H   T   T   H   H   H   T   T   T   T

What could the bias of the coin be? (*We can never say for sure*)

Consider the following sequence of coin tosses:

<p align="center">H   T   T   H   H   H   T   T   T   T</p>

What could the bias of the coin be? (*We can never say for sure*)

What is a good guess for the bias of the coin?

Consider the following sequence of coin tosses:

$$H \quad T \quad T \quad H \quad H \quad H \quad T \quad T \quad T \quad T$$

What could the bias of the coin be? (*We can never say for sure*)

What is a good guess for the bias of the coin?

Intuitively we would say $\frac{\text{number of heads}}{\text{total number of coin tosses}} = \frac{4}{10}$

Consider the following sequence of coin tosses:

<div align="center">

H  T  T  H  H  H  T  T  T  T

</div>

What could the bias of the coin be? (*We can never say for sure*)

What is a good guess for the bias of the coin?

Intuitively we would say $\frac{\text{number of heads}}{\text{total number of coin tosses}} = \frac{4}{10}$

We can use this for predicting further coin tosses!

## Coin bias: The math

Let $D$ be the given data, and $\theta$ be a random variable representing the bias of the coin.

### Definition

The *likelihood function* $L(\theta)$ is given by $L(\theta) = P(D|\theta)$.

Here we have $L(\theta) = \binom{10}{4}\theta^4(1-\theta)^6$ (assuming i.i.d coin tosses).

Setting the derivative wrt $\theta$ to zero to maximize the likelihood, we get:

$$\frac{dL(\theta)}{d\theta} = \binom{10}{4}(\theta^4 \cdot 6(1-\theta)^5 + 4\theta^3 \cdot (1-\theta)^6) = 0$$

$$\implies \boxed{\theta = 0.4}$$

In some sense, this is the best estimate we can get!

## Example 2: Firm profits

Suppose we want to predict the profit of a firm given some simple data:

| Units produced (x) | Cost of production (y) |
|---|---|
| 10 | 592 |
| 20 | 1090 |
| 30 | 1604 |
| 40 | 2122 |
| 50 | 2620 |
| 70 | ? |

We can assume that there is some fixed cost for setting up production (*a*), plus a marginal cost (the cost per unit) required to produce the items (*b*):

$y = f(x) = a + bx$                                    (This is a linear model)

## Example 2: Firm profits

We require a good estimate for $a$ and $b$.

## Example 2: Firm profits

We require a good estimate for $a$ and $b$.

### Definition

Consider a dataset with $n$ input and output pairs, denoted $(x_i, y_i)$. The *mean squared error* (MSE) is defined as $\frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i))^2$.

# Example 2: Firm profits

We require a <span style="color:orange">good</span> estimate for *a* and *b*.

### Definition
Consider a dataset with *n* input and output pairs, denoted $(x_i, y_i)$. The *mean squared error* (MSE) is defined as $\frac{1}{n}\sum_{i=1}^{n}(y_i - f(x_i))^2$.

In this case, the $(x_i, y_i)$ pairs are given by $\{(10, 592), (20, 1090), (30, 1604), (40, 2122), (50, 2620)\}$. We look for values of *a* and *b* that minimize the mean squared error $L = \frac{1}{5}\sum_{i=1}^{5}((a + bx_i) - y_i)^2$.

We require a good estimate for $a$ and $b$.

### Definition

Consider a dataset with $n$ input and output pairs, denoted $(x_i, y_i)$. The *mean squared error* (MSE) is defined as $\frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i))^2$.

In this case, the $(x_i, y_i)$ pairs are given by $\{(10, 592), (20, 1090), (30, 1604), (40, 2122), (50, 2620)\}$. We look for values of $a$ and $b$ that minimize the mean squared error $L = \frac{1}{5} \sum_{i=1}^{5} ((a + bx_i) - y_i)^2$.

We can do this by setting $\frac{\partial L}{\partial a} = 0$, $\frac{\partial L}{\partial b} = 0$ to get $\boxed{a = 79.2, b = 50.8}$.

## Firm profits: The math

$$L = \frac{1}{5} \sum_{i=1}^{5} (a + bx_i - y_i)^2$$

$$\frac{\partial L}{\partial a} = \frac{1}{5} \sum_{i=1}^{5} 2x_i(a + bx_i - y_i) = 0 \tag{1}$$

$$\frac{\partial L}{\partial b} = \frac{1}{5} \sum_{i=1}^{5} 2(a + bx_i - y_i) = 0 \tag{2}$$
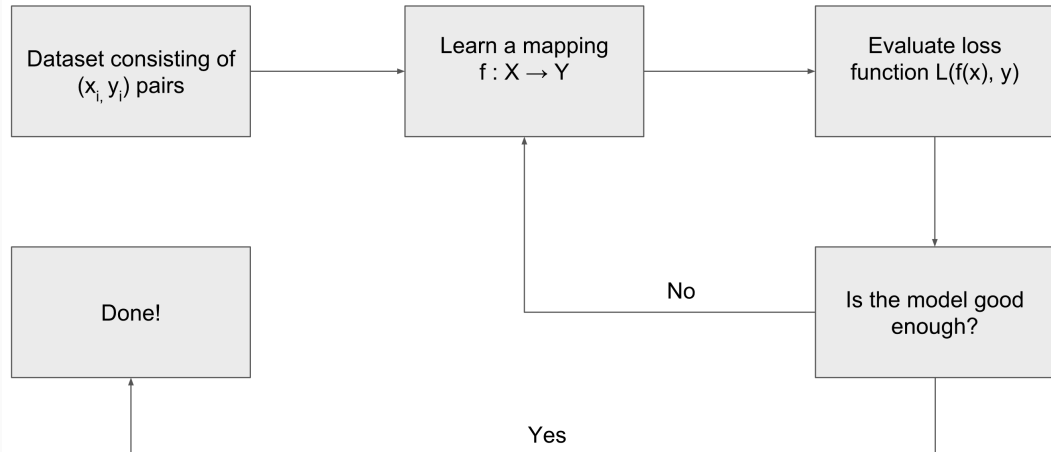
Solving, we get

$$b = \frac{\sum_{i=1}^{5}(x_i - \overline{x})(y_i - \overline{y})}{\sum_{i=1}^{5}(x_i - \overline{x})^2}, \qquad a = \overline{y} - b\overline{x}$$

Where $\overline{x} = \frac{1}{5} \sum_{i=1}^{5} x_i$ and $\overline{y} = \frac{1}{5} \sum_{i=1}^{5} y_i$

- Tries to predict an output *y* given an input *x* by learning a mapping $f : X \rightarrow Y$
- Quality of the estimate defined by a loss function $L(f(x), y)$
- ML algorithms try to minimize the loss function (closed form, gradient descent, etc)

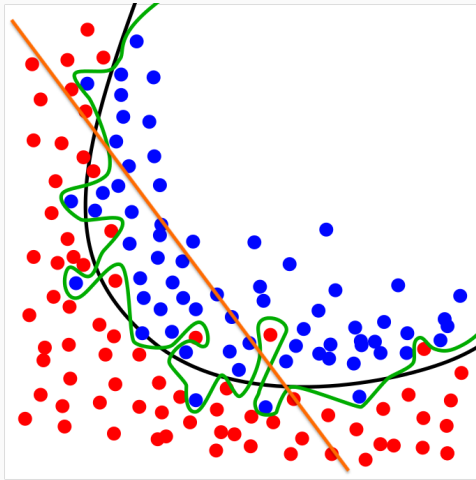Is it really that simple?

## Machine Learning: Some issues

*Overfitting*: The model (function) fits too strongly to the data.
*Underfitting*: The model does not fit well enough to the data.

# Machine Learning: Some issues

*Overfitting*: The model (function) fits too strongly to the data.
*Underfitting*: The model does not fit well enough to the data.

# Deep Learning
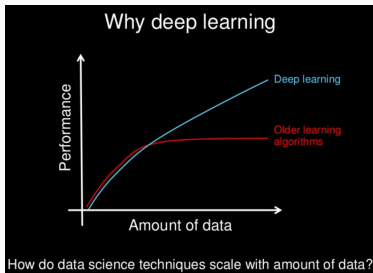
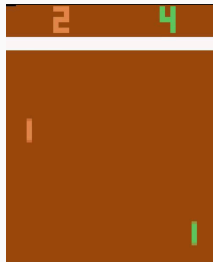- Subfield of Machine Learning inspired by a crude model of the brain, called *artificial neural networks*

- Subfield of Machine Learning inspired by a crude model of the brain, called *artificial neural networks*
- Typically works with large amounts of data and outperforms other machine learning algorithms

- Subfield of Machine Learning inspired by a crude model of the brain, called *artificial neural networks*
- Typically works with large amounts of data and outperforms other machine learning algorithms
- Can be used to learn powerful representations and highly non linear functions



How do data science techniques scale with amount of data?

- In the 1940s, Frank Rosenblatt developed an artificial neuron model that was capable of learning (the "perceptron")

- In the 1940s, Frank Rosenblatt developed an artificial neuron model that was capable of learning (the "perceptron")
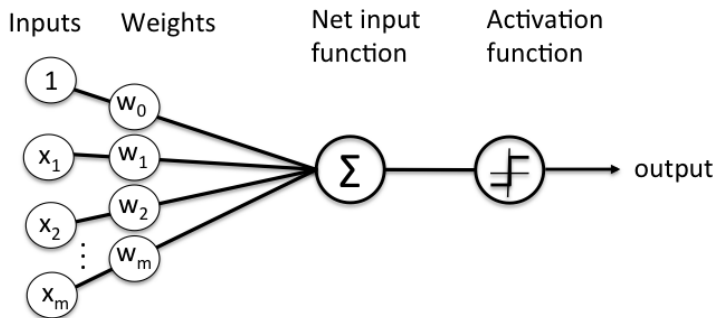- Worked well on small problems that were *linearly separable* (what's that?)

- In the 1940s, Frank Rosenblatt developed an artificial neuron model that was capable of learning (the "perceptron")
- Worked well on small problems that were *linearly separable* (what's that?)
- Was eventually superseded by support vector machines (SVMs) which were more stable and performed better on real-world data
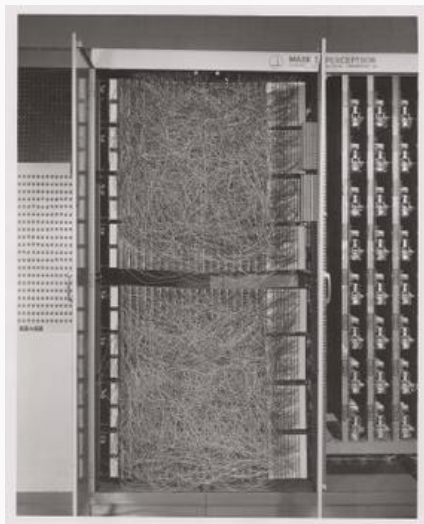
**Schematic of Rosenblatt's perceptron.**

Mathematically, if an input $x = (x_1, ..., x_m)$ is given to the perceptron, the output $y$ is calculated by

$$\hat{y} = f(w_0 + w_1 x_1 + ... + w_m x_m)$$

Where $f(\cdot)$ is the *activation function*, which is a non-linear function of its input. The weights $w_i$ are parameters of the perceptron, and they are *learned* by the perceptron.

Some commonly used activation functions are step, sigmoid, ReLU, tanh, etc.

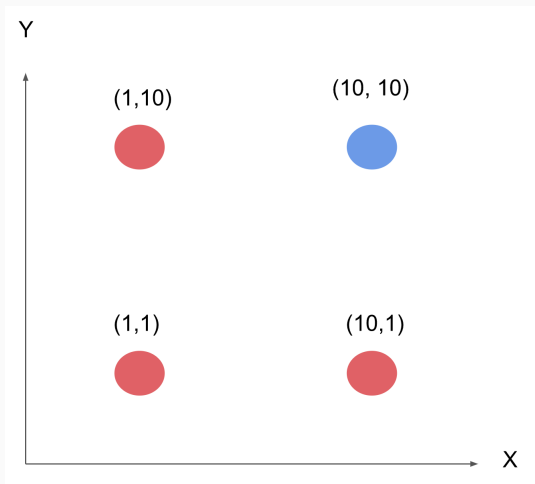Consider a perceptron that uses a step function for its activation:

$$\hat{y} = f(w_0 + w_1 x_1 + ... + w_m x_m)$$

$$= \begin{cases} 1, & w_0 + w_1 x_1 + ... + w_m x_m > 0 \\ 0, & \text{otherwise} \end{cases}$$

What about the boundary region defined by $w_0 + w_1 x_1 + ... + w_m x_m = 0$?

It defines a *plane* in (m-1) dimensions.

Perceptrons find a *separating plane* to try and classify the data.

Perceptrons find a *separating plane* to try and classify the data.

Consider a dataset $\{(x^1, y^1), (x^2, y^2), ..., (x^n, y^n)\}$. The perceptron learning algorithm is as follows:

## Perceptrons: How do they learn?

Consider a dataset $\{(x^1, y^1), (x^2, y^2), ..., (x^n, y^n)\}$. The perceptron learning algorithm is as follows:

1. Initialize all weights $w_i$ to 0

## Perceptrons: How do they learn?

Consider a dataset $\{(x^1, y^1), (x^2, y^2), ..., (x^n, y^n)\}$. The perceptron learning algorithm is as follows:

1. Initialize all weights $w_i$ to 0
2. For each training example $(x^j, y^j)$, calculate the output

$$\hat{y}^j = f(w_0 + w_1 x_1^j + + w_2 x_2^j + ..., + w_m x_m^j)$$

## Perceptrons: How do they learn?

Consider a dataset $\{(x^1, y^1), (x^2, y^2), ..., (x^n, y^n)\}$. The perceptron learning algorithm is as follows:

1. Initialize all weights $w_i$ to 0

2. For each training example $(x^j, y^j)$, calculate the output

$$\hat{y}^j = f(w_0 + w_1 x_1^j + + w_2 x_2^j + ..., + w_m x_m^j)$$

3. Each of the weights $w_i$ are updated as

$$w_i \rightarrow w_i + (y^j - \hat{y}^j) x_i^j$$

Can a linear classifier correctly classify all 4 points?

# Feed forward networks

The perceptron is unable to learn non linear functions like XOR. Is there a way to modify it so that it can learn these functions?

The perceptron is unable to learn non linear functions like XOR. Is there a way to modify it so that it can learn these functions?

The idea is to chain together multiple perceptrons to learn more complex functions. This is done by feeding the output of a perceptron as the input of another perceptron.

$a = f(w_{10} + w_{11}x + w_{12}y)$

$c = f(w_{30} + w_{31}a + w_{32}b)$

output

$b = f(w_{20} + w_{21}x + w_{22}y)$

The parameters to be learned are $w_{10}, w_{11}, w_{12}, w_{20}, w_{21}, w_{22}, w_{30}, w_{31}, w_{32}$. This successfully manages to learn the XOR function!

# General feed forward networks

This idea can be extended to many layers and many perceptron units (neurons) in each layer.



input layer

hidden layer 1   hidden layer 2

output layer

This is called a feed forward neural network.

Perceptrons can represent linear functions, and a few perceptrons chained together can represent non linear functions like XOR. What could a general feed forward neural network represent?

Perceptrons can represent linear functions, and a few perceptrons chained together can represent non linear functions like XOR. What could a general feed forward neural network represent?

### Theorem (Universal Approximation)

Any continuous function can be approximated to arbitrary accuracy by a feed forward neural network with the appropriate architecture.

This is great, we can learn anything we need!

Perceptrons can represent linear functions, and a few perceptrons chained together can represent non linear functions like XOR. What could a general feed forward neural network represent?

### Theorem (Universal Approximation)
Any continuous function can be approximated to arbitrary accuracy by a feed forward neural network with the appropriate architecture.

This is great, we can learn anything we need! (right?)

- (Stochastic) gradient descent is used to optimize the loss function for learning.

- (Stochastic) gradient descent is used to optimize the loss function for learning.
- The gradient is calculated using the *backpropagation algorithm*, a dynamic programming algorithm that keeps track of various derivatives that are required to construct the final gradient.

- (Stochastic) gradient descent is used to optimize the loss function for learning.
- The gradient is calculated using the *backpropagation algorithm*, a dynamic programming algorithm that keeps track of various derivatives that are required to construct the final gradient.
- Need not write code for this, already done for you by software

## The backpropagation algorithm (*)

- $L$: The total number of layers in the network (0 indexed)
- $n_l$: The number of neurons in the layer $l$
- $x_i^l$: The input to the $i^{th}$ neuron in the $l^{th}$ layer
- $o_i^l$: The output of the $i^{th}$ neuron in the $l^{th}$ layer
- $W_{ij}^l$: The weight connecting neuron $i$ of layer $l-1$ to neuron $j$ of layer $l$
- $b_i^l$: The bias of the neuron $i$ in layer $l$
- $f$: The activation function used at each hidden neuron (including the output layer)
- $E$: The error (loss) function. Here we assume that the squared error is used.

## The backpropagation algorithm (*)

Since the input $x_j^l$ only affects neurons after it and $x_j^l$ can be written as a function of $W_{ij}^l$ and other variables, we can write:

$$\frac{\partial E}{\partial W_{ij}^l} = \frac{\partial E}{\partial x_j^l} \frac{\partial x_j^l}{\partial W_{ij}^l}$$

$$\frac{\partial E}{\partial b_j^l} = \frac{\partial E}{\partial x_j^l} \frac{\partial b_j^l}{\partial b_j^l}$$

Defining $\delta_j^l \equiv \frac{\partial E}{\partial x_j^l}$, we have

$$\frac{\partial E}{\partial W_{ij}^l} = \delta_j^l \frac{\partial x_j^l}{\partial W_{ij}^l} \tag{1}$$

$$\frac{\partial E}{b_j^l} = \delta_j^l \frac{\partial x_j^l}{b_j^l} \tag{2}$$

## The backpropagation algorithm (*)

By the construction of the network,

$$x_j^l = \left[ \sum_{i=1}^{n_{l-1}} W_{ij}^l o_i^{l-1} \right] + b_j^l$$

And so

$$\frac{\partial x_j^l}{\partial W_{ij}^l} = o_i^{l-1} \tag{3}$$

$$\frac{\partial x_j^l}{\partial b_j^l} = 1 \tag{4}$$

## The backpropagation algorithm (*)

Now, if $l$ is an output layer we can write

$$\delta_j^l = \frac{\partial E}{\partial x_j^l} = \frac{\partial E}{\partial o_j^l} \frac{\partial o_j^l}{\partial x_j^l}$$

Since we are using the squared loss, $E = \frac{1}{2}(\sum_{j=1}^{n_L}(t - o_j^L)^2)$ where $t$ is the target output. This gives an easy way to calculate the first term. For the second term, note that the derivative of the output of a neuron with respect to the input is just the derivative of the activation function. Putting these two together, we get

$$\delta_j^l = (t - o_j^L)f'(x_j^l), \quad \text{If } l \text{ is an output layer} \tag{5}$$

## The backpropagation algorithm (*)

If $l$ is not an output layer, we can determine $\delta_j^l$ in terms of the $\delta_k^{l+1}$ values. Since changing the input to a neuron affects all inputs of the next layer of neurons, we can write

$$\delta_j^l = \frac{\partial E}{\partial x_j^l} = \sum_{k=0}^{n_{l+1}} \frac{\partial E}{\partial x_k^{l+1}} \frac{\partial x_k^{l+1}}{\partial o_j^l} \frac{\partial o_j^l}{\partial x_j^l}$$

Finally we have $x_k^{l+1} = \sum_{j=1}^{n_l} W_{jk}^{l+1} o_j^l$, so $\frac{\partial x_k^{l+1}}{\partial o_j^l} = W_{jk}^{l+1}$, and

$$\delta_j^l = \sum_{k=0}^{n_{l+1}} \delta_k^{l+1} W_{jk}^{l+1} f'(x_j^l), \quad \text{If } l \text{ is not an output layer} \tag{6}$$

## The backpropagation algorithm (*)

Putting together (1), (2), (3), (4), (5), (6), we get

$$\frac{\partial E}{\partial W_{ij}^l} = \delta_j^l o_i^{l-1}$$

$$\frac{\partial E}{\partial b_j^l} = \delta_j^l$$

$$\delta_j^l = (t - o_j^l) f'(x_j^l), \quad \text{If } l \text{ is an output layer}$$

$$\delta_j^l = \sum_{k=0}^{n_{l+1}} \delta_k^{l+1} W_{jk}^{l+1} f'(x_j^l), \quad \text{If } l \text{ is not an output layer}$$

## The backpropagation algorithm (*)

This can be written in matrix/vector form as

### Backpropagation equations

$$\frac{\partial E}{\partial W^l} = \delta^l (o^{l-1})^T$$

$$\frac{\partial E}{\partial b^l} = \delta^l$$

$$\delta^l = (t - o^L) \circ f'(x^l), \quad \text{If } l \text{ is an output layer}$$

$$\delta^l = (W^{l+1} \delta^{l+1}) \circ f'(x^l), \quad \text{If } l \text{ is not an output layer}$$

### Update equation

$$W \leftarrow W - \eta \frac{\partial E}{\partial W}$$

## The backpropagation algorithm (*)

This can be written in matrix/vector form as

### Backpropagation equations

$$\frac{\partial E}{\partial W^l} = \delta^l (o^{l-1})^T$$

$$\frac{\partial E}{\partial b^l} = \delta^l$$

$$\delta^l = (t - o^L) \circ f'(x^l), \quad \text{If } l \text{ is an output layer}$$

$$\delta^l = (W^{l+1} \delta^{l+1}) \circ f'(x^l), \quad \text{If } l \text{ is not an output layer}$$

### Update equation

$$W \leftarrow W - \eta \frac{\partial E}{\partial W}$$

# The backpropagation algorithm (*)

## Update equation

$$W \leftarrow W - \eta \frac{\partial E}{\partial W}$$

Where $\eta$ is a parameter called the learning rate.

## Update equation

$$W \leftarrow W - \eta \frac{\partial E}{\partial W}$$

Where $\eta$ is a parameter called the learning rate.

This is way too complicated! There has to be an easier way!

## Update equation

$$W \leftarrow W - \eta \frac{\partial E}{\partial W}$$

Where $\eta$ is a parameter called the learning rate.

This is way too complicated! There has to be an easier way!

Enter Tensorflow.

# Tensorflow

- An open source framework for ML/DL developed by Google
- Does all the heavy math for you, no need to code gradient descent. Specifying the model architecture is good enough
- Uses a computational graph internally to accomplish all of this. The graph can be manipulated too
- Has a Python library available (that we are going to use)

A placeholder in Tensorflow is an object that allows use to define computations on it without giving it a value. It is initialized as and when required, for example when giving it as input to a function.

A placeholder in Tensorflow is an object that allows use to define computations on it without giving it a value. It is initialized as and when required, for example when giving it as input to a function.

```python
import tensorflow as tf
x = tf.placeholder(tf.float32, (None, 1))
```

This declares a placeholder variable *x* that has size (unknown ×1).

## Variables

A Variable in Tensorflow is an object that holds a value, and is typically something that we are optimizing over (like the weights of a neural network). They must be initialized in code.

A Variable in Tensorflow is an object that holds a value, and is typically something that we are optimizing over (like the weights of a neural network). They must be initialized in code.

```python
import tensorflow as tf
m = tf.Variable(tf.zeros, ([1]))
```

This declares a Variable *m* that has size 1 (think of it like an array of size 1).

A Variable in Tensorflow is an object that holds a value, and is typically something that we are optimizing over (like the weights of a neural network). They must be initialized in code.

```python
import tensorflow as tf
m = tf.Variable(tf.zeros, ([1]))
```

This declares a Variable *m* that has size 1 (think of it like an array of size 1).

In short, placeholders are input data that does not change as the model trains, and Variables are model parameters that change with time.

# Sessions

A Session in Tensorflow is where code actually gets executed. Outside of a session, only variables and the computational graph is defined. Code execution requires that a Tensorflow Session be declared.

# Sessions

A Session in Tensorflow is where code actually gets executed. Outside of a session, only variables and the computational graph is defined. Code execution requires that a Tensorflow Session be declared.

```python
import tensorflow as tf
sess = tf.Session() # Open a session
sess.run(something) # Do stuff
sess.close() # Close the session
```

Or,

```python
import tensorflow as tf
with tf.Session() as sess: # Open a session
    sess.run(something) # Do stuff
```

| Units produced (x) | Cost of production (y) |
|:---:|:---:|
| 10 | 592 |
| 20 | 1090 |
| 30 | 1604 |
| 40 | 2122 |
| 50 | 2620 |

How would we find a good linear model for this problem using Tensorflow?

# Firm costs revisited

(Download this at goo.gl/5yayn5)

```python
import numpy as np
import tensorflow as tf

# Model linear regression y_hat = a + bx

# Placeholder variables for the input data
x = tf.placeholder(tf.float32, [None, 1])
y = tf.placeholder(tf.float32, [None, 1])

# Model parameters are Variables in Tensorflow
a = tf.Variable(tf.zeros([1]))
b = tf.Variable(tf.zeros([1, 1]))

# Calculating the prediction of the linear model
y_hat = a + tf.matmul(x,b)

# Loss function =  sum((y_hat-y)^2)
cost = tf.reduce_mean(tf.square(y_hat - y))

# Training using Gradient Descent to minimize cost
train_step = tf.train.GradientDescentOptimizer(0.000001).minimize(cost)
```

```
22
23    sess = tf.Session()
24    init = tf.global_variables_initializer()
25    sess.run(init)
26    steps = 20000
27    for i in range(steps):
28        # Input data for the model
29        xs = np.array([[10], [20], [30], [40], [50]])
30        ys = np.array([[592], [1090], [1604], [2122], [2620]])
31
32        # Defining a feed dictionary; this is the set of
33        # (placeholder : value) pairs for calculations
34        feed = { x : xs, y : ys }
35
36        # Running a single training step
37        sess.run(train_step, feed_dict = feed)
38
39    # Printing the results
40    print("a: %f" % sess.run(b))
41    print("b: %f" % sess.run(a))
```

# Computational graph

```python
1   import numpy as np
2   import tensorflow as tf
3
4   # Model linear regression y_hat = a + bx
5
6   # Placeholder variables for the input data
7   x = tf.placeholder(tf.float32, [None, 1])
8   y = tf.placeholder(tf.float32, [None, 1])
9
10  # Model parameters are Variables in Tensorflow
11  a = tf.Variable(tf.zeros([1]))
12  b = tf.Variable(tf.zeros([1, 1]))
13
14  # Calculating the prediction of the linear model
15  y_hat = a + tf.matmul(x,b)
16
17  # Loss function =  sum((y_hat-y)^2)
18  cost = tf.reduce_mean(tf.square(y_hat - y))
19
20  # Training using Gradient Descent to minimize cost
21  train_step = tf.train.GradientDescentOptimizer(0.000001).minimize(cost)
```

43

- *All* computations happen within a session. This includes printing the value of a variable!

## A few important things about Tensorflow

- *All* computations happen within a session. This includes printing the value of a variable!

- To build the computational graph, all arithmetic operations for the model must be done by Tensorflow. For example, tf.reduce_sum must be used instead of Python's builtin sum() function.

## A few important things about Tensorflow

- *All* computations happen within a session. This includes printing the value of a variable!

- To build the computational graph, all arithmetic operations for the model must be done by Tensorflow. For example, tf.reduce_sum must be used instead of Python's builtin sum() function.

- When declaring a placeholder, if a dimension is "None", it simply means that it is unknown.

Thank you!