

Concepts

Unit refers to the activation function in a layer by which the inputs are transformed via a nonlinear activation function (e.g. sigmoid, ReLU, tanh, etc.). Usually, a unit has several incoming connections and several outgoing connections.

Input Layer: Comprised of multiple Real-Valued inputs. Each input must be linearly independent from each other.

Hidden Layers: Layers other than the input and output layers. A layer is the highest-level building block in deep learning. A layer is a container that usually receives weighted input, transforms it with a set of mostly non-linear functions and then passes these values as output to the next layer.

Batch Normalization: Using mini-batches of examples, as opposed to one example at a time, is helpful in several ways. First, the gradient of the loss over a mini-batch is an estimate of the gradient over the entire batch, whose quality improves as the batch size increases. Second, computation over a batch can be much more efficient than in computers for individual examples, due to the parallelism afforded by the modern computing platforms.

Maximum Likelihood Estimation (MLE): In general, for a fixed set of data and underlying statistical model, the method of maximum likelihood selects the set of values of the model parameters that maximizes the likelihood function. Intuitively, this maximizes the "alignment" of the statistical model with the observed data and for discrete random variables it indeed maximizes the probability of the observed data under the resulting distribution. Maximum-likelihood estimation gives a unified approach to estimation, which is well defined in the case of the normal distribution and many other problems.

Cross-Entropy: Cross entropy can be used to define the loss function in machine learning and optimization. The true probability p is the true label, and the given distribution q is the predicted value of the current model.

Logistic: The logistic loss function is defined as: $V(f(\theta; x), y) = \frac{1}{2} \ln(1 + e^{2f(\theta; x) - y})$

Quadratic: The use of a quadratic loss function is common, for example when using least squares techniques. It is often more mathematically tractable than other loss functions because of the properties of variances, as well as being symmetric, an error above the target causes the same loss as the same magnitude of error below the target. If the target is 1, then a quadratic loss function is: $L(x) = C(1 - x)^2$

0-1 Loss: In statistics and decision theory, a frequently used loss function is the 0-1 loss function: $L(\hat{y}, y) = I(\hat{y} \neq y)$

Hinge Loss: The hinge loss is a loss function used for training classifiers. For an intended output $t = \pm 1$ and a classifier score x , the hinge loss of the prediction is defined as: $L(y) = \max(0, 1 - ty)$

Exponential: $L(x, y) = \exp(\frac{1}{2} \sum_{i=1}^n (x_i - y_i)^2)$

Hellinger Distance: It is used to quantify the similarity between two probability distributions. It is a type of divergence.

Kullback-Leibler Divergence: It is a measure of how one probability distribution diverges from a second expected probability distribution. Applications include characterizing the relative information entropy in information systems, randomness in continuous time-series, and information gain when comparing statistical models of interest.

Itakura-Saito distance: It is a measure of the difference between an original spectrum $P(\omega)$ and an approximation $Q(\omega)$ of that spectrum. Although it is not a perceptual measure, it is intended to reflect perceptual dissimilarity.

L1 norm: Manhattan Distance. L1-norm is also known as least absolute deviations (LAD), least absolute errors (LAE). It is basically minimizing the sum of the absolute differences (b) between the target value and the estimated values. $S = \sum_{i=1}^n |y_i - f(x_i)|$

L2 norm: Euclidean Distance. L2-norm is also known as least squares. It is basically minimizing the sum of the square of the differences (b) between the target value and the estimated values. $S = \sum_{i=1}^n (y_i - f(x_i))^2$

Early Stopping: Early stopping rules provide guidance as to how many iterations can be run before the learner begins to overfit and stop the algorithm then.

Dropout: Is a regularization technique for reducing overfitting in neural networks by preventing complex co-adaptations on training data. It is a very efficient way of performing model averaging with neural networks. The term "dropout" refers to dropping out units (both hidden and visible) in a neural network.

Sparsity regularizer on columns: This regularizer defines an L2 norm on each column and an L1 norm over all columns. It can be solved by proximal methods. $R(w) = \sum_{i=1}^n |w_{ij}|$

Nuclear norm regularization: $R(W) = \|W\|_*$, where $\|W\|_*$ is the approximation to the singular value decomposition of W .

Mean-constrained regularization: This regularizer constrains the functions learned for each task to be similar to the overall average of the functions across all tasks. It is useful for expressing prior information that each task is expected to share similarities with each other task. An example is predicting blood iron levels measured at different times of the day, where each task represents a different portion.

Clustered mean-constrained regularization: This regularizer is similar to the mean-constrained regularizer, but instead enforces similarity between tasks within the same cluster. This can capture more complex prior information. This technique has been used to predict Netflix recommendations.

Graph-based similarity: More general than above, similarity between tasks can be defined by a function. The regularizer encourages the model to learn similar functions for similar tasks. $R(f_1, f_2) = \sum_{i,j} \sum_{k,l} \sum_{m,n} \sum_{p,q} \sum_{r,s} \sum_{t,u} \sum_{v,w} \sum_{x,y} \sum_{z,\dots} \dots$

Activation Functions

Define the output of that node given an input or set of inputs.

- ReLU**: $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$
- Sigmoid / Logistic**: $f(x) = \frac{1}{1 + e^{-x}}$
- Binary**: $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
- Tanh**: $f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
- Softplus**: $f(x) = \ln(1 + e^x)$
- Softmax**: $f_i(x) = \frac{e^{x_i}}{\sum_{j=1}^J e^{x_j}}$ for $i = 1, \dots, J$
- Maxout**: $f(x) = \max_k z_k$

Types: Leaky ReLU, PReLU, RReLU, ELU, SELU, and others.

Backpropagation

In this method, we reuse partial derivatives computed for higher layers in lower layers, for efficiency.

Neural Network taking 4 dimension vector representation of words.

$x = U^T(Wx + b)$, $x \in \mathbb{R}^{20 \times 1}$, $W \in \mathbb{R}^{8 \times 20}$, $b \in \mathbb{R}^{8 \times 1}$

$z = U^T(Wx + b)$, $z \in \mathbb{R}^{8 \times 1}$

$u = \sigma(z)$

$y = W'u + b'$, $y \in \mathbb{R}^{10 \times 1}$

$W' = U^T W$, $b' = U^T b + b'$

Intuition for backpropagation

Simple Example (Carroll)

Simple Example (Flowgraph)

However, if you actually try that, the weights will change far too much each iteration, which will make them "overcorrect" and the loss will actually increase/diverge. So in practice, people usually multiply each derivative by a small value called the "learning rate" before they subtract it from its corresponding weight.

Neural networks are often trained by gradient descent on the weights. This means at each iteration we use backpropagation to calculate the derivative of the loss function with respect to each weight and subtract it from that weight.

Learning Rate

Simplest recipe: keep it fixed and use the same for all parameters.

Better results by allowing learning rates to decrease: Options:

- Reduce by 0.5 when validation error stops improving.
- Convergence guarantees, with hyperparameters η and τ is iteration numbers.
- Better: No hard-stopping learning of rates by using Adagrad.

Tricks:

- $\alpha = \frac{\alpha_0 \tau}{1 + \tau}$

Optimization

Gradient Descent: A first-order iterative optimization algorithm for finding the minimum of a function. To find a local minimum of a function uses gradient descent one takes steps proportional to the negative of the gradient (or of the approximate gradient) of the function at the current point. If instead one takes steps proportional to the positive of the gradient, one approaches a local maximum of that function; this then known as gradient ascent.

Stochastic Gradient Descent (SGD): Gradient descent uses total gradient over all examples per update, SGD updates after every 1 or few examples.

Mini-batch Stochastic Gradient Descent (SGD): Gradient descent uses total gradient over all examples per update, SGD updates after every 1 example.

Momentum: Idea: Add a fraction v of previous update to current one. When the gradient keeps pointing in the same direction, this will increase the size of the steps taken towards the minimum.

Adagrad: Adaptive learning rates for each parameter.

Weight Initialization

Alz0 Initialization: In the ideal situation, with proper data normalization, it is reasonable to assume that approximately half of the weights will be positive and half of them will be negative. A reasonable-sounding idea then might be to set all the initial weights to zero, which you expect to be the "best guess" in expectation.

Initialization with Small Random Numbers: Thus, you still want the weights to be very close to zero, but not identically zero. In this way, you can randomize those neurons to small numbers which are very close to zero, and it is treated as symmetry breaking. The idea is that the neurons are all random and unique in the beginning so that each can contribute differently to the updates and integrate themselves as diverse parts of the full network.

Calibrating the Variances: One problem with the above suggestion is that the distribution of the outputs from a randomly initialized neuron is not a normal distribution and empirically improves the rate of convergence. The detailed derivations can be found from Page 18 to 25 of the slides. Please note that in the derivations, it does not consider the influence of ReLU neurons.

Weight Initialization: The implementation for weights might simply draw values from a normal distribution with zero mean, and unit standard deviation. It is also possible to use small numbers drawn from a uniform distribution, but this seems to have relatively little impact on the final performance (I practice).

ReLU: $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$

Sigmoid / Logistic: $f(x) = \frac{1}{1 + e^{-x}}$

Binary: $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$

Tanh: $f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

Softplus: $f(x) = \ln(1 + e^x)$

Softmax: $f_i(x) = \frac{e^{x_i}}{\sum_{j=1}^J e^{x_j}}$ for $i = 1, \dots, J$

Maxout: $f(x) = \max_k z_k$

Leaky ReLU, PReLU, RReLU, ELU, SELU, and others.

Neural Network taking 4 dimension vector representation of words.

$x = U^T(Wx + b)$, $x \in \mathbb{R}^{20 \times 1}$, $W \in \mathbb{R}^{8 \times 20}$, $b \in \mathbb{R}^{8 \times 1}$

$z = U^T(Wx + b)$, $z \in \mathbb{R}^{8 \times 1}$

$u = \sigma(z)$

$y = W'u + b'$, $y \in \mathbb{R}^{10 \times 1}$

$W' = U^T W$, $b' = U^T b + b'$

Simple Example (Carroll)

Simple Example (Flowgraph)

However, if you actually try that, the weights will change far too much each iteration, which will make them "overcorrect" and the loss will actually increase/diverge. So in practice, people usually multiply each derivative by a small value called the "learning rate" before they subtract it from its corresponding weight.

Neural networks are often trained by gradient descent on the weights. This means at each iteration we use backpropagation to calculate the derivative of the loss function with respect to each weight and subtract it from that weight.

Learning Rate

Simplest recipe: keep it fixed and use the same for all parameters.

Better results by allowing learning rates to decrease: Options:

- Reduce by 0.5 when validation error stops improving.
- Convergence guarantees, with hyperparameters η and τ is iteration numbers.
- Better: No hard-stopping learning of rates by using Adagrad.

Tricks:

- $\alpha = \frac{\alpha_0 \tau}{1 + \tau}$

Optimization

Gradient Descent: A first-order iterative optimization algorithm for finding the minimum of a function. To find a local minimum of a function uses gradient descent one takes steps proportional to the negative of the gradient (or of the approximate gradient) of the function at the current point. If instead one takes steps proportional to the positive of the gradient, one approaches a local maximum of that function; this then known as gradient ascent.

Stochastic Gradient Descent (SGD): Gradient descent uses total gradient over all examples per update, SGD updates after every 1 or few examples.

Mini-batch Stochastic Gradient Descent (SGD): Gradient descent uses total gradient over all examples per update, SGD updates after every 1 example.

Momentum: Idea: Add a fraction v of previous update to current one. When the gradient keeps pointing in the same direction, this will increase the size of the steps taken towards the minimum.

Adagrad: Adaptive learning rates for each parameter.

Weight Initialization

Alz0 Initialization: In the ideal situation, with proper data normalization, it is reasonable to assume that approximately half of the weights will be positive and half of them will be negative. A reasonable-sounding idea then might be to set all the initial weights to zero, which you expect to be the "best guess" in expectation.

Initialization with Small Random Numbers: Thus, you still want the weights to be very close to zero, but not identically zero. In this way, you can randomize those neurons to small numbers which are very close to zero, and it is treated as symmetry breaking. The idea is that the neurons are all random and unique in the beginning so that each can contribute differently to the updates and integrate themselves as diverse parts of the full network.

Calibrating the Variances: One problem with the above suggestion is that the distribution of the outputs from a randomly initialized neuron is not a normal distribution and empirically improves the rate of convergence. The detailed derivations can be found from Page 18 to 25 of the slides. Please note that in the derivations, it does not consider the influence of ReLU neurons.

Weight Initialization: The implementation for weights might simply draw values from a normal distribution with zero mean, and unit standard deviation. It is also possible to use small numbers drawn from a uniform distribution, but this seems to have relatively little impact on the final performance (I practice).

Architectures

Strategy

1. Select Network Structure appropriate for problem

Structure: Single words, fixed windows, sentence based, document level; bag of words, recursive vs. recurrent, CNN

Nonlinearity (Activation Functions)

1. Implement your gradient
2. Implement a finite difference computation by looping through the parameters of your network, adding and subtracting a small epsilon (~10^-4) and estimate derivatives
3. Compare the two and make sure they are almost the same

$$f'(\theta) \approx \frac{J(\theta^{(t+\epsilon)}) - J(\theta^{(t-\epsilon)})}{2\epsilon} \quad \theta^{(t+\epsilon)} = \theta + \epsilon \times x_i$$

2. Check for implementation bugs with gradient checks

If you gradient fails and you don't know why?

Using Gradient Checks

- Simplify your model until you have no bug!
- What now? Create a very tiny synthetic model and dataset
- Only softmax on fixed input
- Backprop into word vectors and softmax
- Add single unit single hidden layer
- Add multi unit single layer
- Add second layer single unit, add multiple units, bias - Add one softmax on top, then two softmax layers
- Add bias

3. Parameter initialization

Initialize hidden layer biases to 0 and output (or reconstruction) biases to optimal value if weights were 0 (e.g., mean target or inverse sigmoid of mean target).

Initialize weights ~ Uniform(-r, r), r inversely proportional to fan-in (previous layer size) and fan-out (next layer size):

$$\sqrt{6 / (\text{fan-in} + \text{fan-out})}$$

Is a first-order iterative optimization algorithm for finding the minimum of a function. To find a local minimum of a function using gradient descent, one takes steps proportional to the negative of the gradient (or of the approximate gradient) of the function at the current point. If instead one takes steps proportional to the positive of the gradient, one approaches a local maximum of that function; the procedure is then known as gradient ascent.

Gradient Descent

Gradient descent uses total gradient over all examples per update, SGD updates after only 1 or few examples:

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J_t(\theta)$$

Stochastic Gradient Descent (SGD)

Ordinary gradient descent as a batch method is very slow, should never be used. Use 2nd order batch method such as L-BFGS.

On large datasets, SGD usually wins over all batch methods. On smaller datasets L-BFGS or Conjugate Gradients win. Large batch L-BFGS extends the reach of L-BFGS [Le et al. ICLR 2011].

Mini-batch Stochastic Gradient Descent (SGD)

Gradient descent uses total gradient over all examples per update, SGD updates after only 1 example

Most commonly used now. Size of each mini batch B: 20 to 1000

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J_{t:t+B}(\theta)$$

Helps parallelizing any model by computing gradients for multiple elements of the batch in parallel

Momentum

Idea: Add a fraction v of previous update to current one. When the gradient keeps pointing in the same direction, this will increase the size of the steps taken towards the minimum.

$$v = \mu v - \alpha \nabla_{\theta} J_t(\theta)$$

$$\theta^{new} = \theta^{old} + v$$

v is initialized at 0
Momentum often increased after some epochs (0.5 to 0.99)

Adagrad

Adaptive learning rates for each parameter
Learning rate is adapting differently for each parameter and rare parameters get larger updates than frequently occurring parameters. Word vectors!

$$\text{Let } g_{t,i} = \frac{\partial}{\partial \theta_i} J_t(\theta), \text{ then: } \theta_{t,i} = \theta_{t-1,i} - \frac{\alpha}{\sqrt{\sum_{\tau=1}^t g_{\tau,i}^2}} g_{t,i}$$

5. Check if the model is powerful enough to overfit

If not, change model structure or make model "larger"

- Simple first step: Reduce model size by lowering number of units and layers and other parameters
- Standard L1 or L2 regularization on weights
- Early Stopping: Use parameters that gave best validation error.
- Sparsity constraints on hidden activations, e.g., add to cost.

If you can overfit: Regularize to prevent overfitting:

- Training time: at each instance of evaluation (in online SGD-training), randomly set 50% of the inputs to each neuron to 0
- Test time: halve the model weights (now twice as many). This prevents feature co-adaptation: A feature cannot only be useful in the presence of particular other features
- In a single layer: A kind of middle-ground between Naive Bayes (where all feature weights are set independently) and logistic regression models (where weights are set in the context of all others)
- Can be thought of as a form of model bagging
- It also acts as a strong regularizer

Feed Forward

Single-Layer Perceptron

Is an artificial neural network wherein connections between the units do not form a cycle. In this network, the information moves in only one direction, forward, from the input nodes, through the hidden nodes (if any) and to the output nodes. There are no cycles or loops in the network.

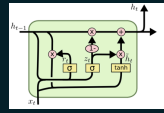
The inputs are fed directly to the outputs via a series of weights. By adding an Logistic activation function to the outputs, the model is identical to a classical Logistic Regression model.

Multi-Layer Perceptron

This class of networks consists of multiple layers of computational units, usually interconnected in a feed-forward way. Each neuron in one layer has directed connections to the neurons of the subsequent layer. In many applications the units of these networks apply a sigmoid function as an activation function.

Kinds

LSTMs



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

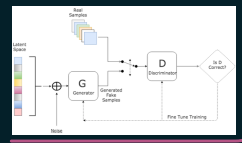
$$\tilde{h}_t = \tanh(W_{\tilde{h}} \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Long short-term memory - It is a type of recurrent (RNN), allowing data to flow both forwards and backwards within the network.

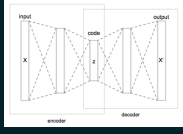
An LSTM is well-suited to learn from experience to classify, process and predict time series given time lags of unknown size and bound between important events. Relative insensitivity to gap length gives an advantage to LSTM over alternative RNNs, hidden Markov models and other sequence learning methods in numerous applications.

GANs



GANs or Generative Adversarial Networks are a class of artificial intelligence algorithms used in unsupervised machine learning, implemented by a system of two neural networks contesting with each other in a zero-sum game framework.

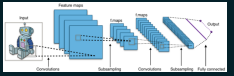
Auto-Encoders



Is an artificial neural network used for unsupervised learning of efficient codings.

The aim of an autoencoder is to learn a representation (encoding) for a set of data, typically for the purpose of dimensionality reduction. Recently, the autoencoder concept has become more widely used for learning generative models of data.

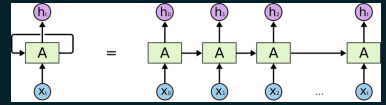
Convolutional Neural Networks (CNN)



They have applications in image and video recognition, recommender systems and natural language processing.

- Pooling
- Convolution
- Subsampling

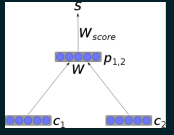
RNNs (Recurrent)



Is a class of artificial neural network where connections between units form a directed cycle. This allows it to exhibit dynamic temporal behavior. Unlike feedforward neural networks, RNNs can use their internal memory to process arbitrary sequences of inputs.

This makes them applicable to tasks such as unsegmented, connected handwriting recognition or speech recognition.

RNNs (Recursive)



Is a kind of deep neural network created by applying the same set of weights recursively over a structure, to produce a structured prediction over variable-size input structures, or a scalar prediction on it, by traversing a given structure in topological order.

RNNs have been successful for instance in learning sequence and tree structures in natural language processing, mainly phrase and sentence continuous representations based on word embedding.



TensorFlow is a deep learning library recently open-sourced by Google. It provides primitives for defining functions on tensors and automatically computing their derivatives, expressed as a graph.

The Tensorflow Graph is built to contain all placeholders for X and y, all variables for W's and b's, all mathematical operations, the cost function, and the optimisation procedure. Then, at runtime, the values for the data are fed into that Graph, by placing the data batches in the placeholders and running the Graph.

Each node in the Graph can then be connected to each other node over the network, and thus running Tensorflow models can be parallelised.

TensorFlow has some neat built-in visualization tools (TensorBoard)

Assembles a computational graph

The computation graph has no numerical value until evaluated.

All computations add nodes to global default graph

A Session object encapsulates the environment in which Tensor objects are evaluated

Uses a session to execute ops in the graph

Declared variables must be initialised before they have values.

When you train a model you use variables to hold and update parameters. Variables are in-memory buffers containing tensors.

Stateful nodes that output their current value, their state is retained across multiple executions of the graph.

Mostly Parameters we're interested in tuning, such as Weights (W) and Biases (b).

Variables can be shared by Explicitly passing tf.Variable objects around, or...

Provides simple name spacing to avoid cases when querying

Creates/Access variables from a variable scope

Nodes whose value is fed at execution time.

Inputs, Features (X) and Labels (y).

MathMul, Add, ReLU, etc.

They are Operations, containing any number of inputs and outputs.

The tensors that flow between the nodes.

It's a binding to a particular execution context: CPU, GPU.

List of graph nodes. Returns the output of these nodes.

Dictionary mapping from graph nodes to concrete values.

Specified the value of each graph node given in the dictionary.

Does lazy evaluation. Need to build the graph, and then run it in a session.

```

import tensorflow as tf

# Create a variable named 'W'
W = tf.Variable([.3], dtype=tf.float32)

# Create a variable named 'b'
b = tf.Variable([.1], dtype=tf.float32)

# Create a variable named 'x'
x = tf.placeholder(tf.float32)

# Operation to add W and b
addition = tf.add(W, b)

# Operation to multiply addition with x
multiply = tf.multiply(addition, x)

# Print out addition and multiply
print(addition)
print(multiply)

# Launch the graph in a session
sess = tf.Session()

# Run the operations
print(sess.run(addition))
print(sess.run(multiply))

```

```

import tensorflow as tf

# Create a variable named 'W'
W = tf.Variable([.3], dtype=tf.float32)

# Create a variable named 'b'
b = tf.Variable([.1], dtype=tf.float32)

# Create a variable named 'x'
x = tf.placeholder(tf.float32)

# Operation to add W and b
addition = tf.add(W, b)

# Operation to multiply addition with x
multiply = tf.multiply(addition, x)

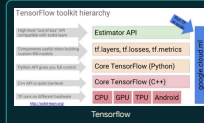
# Print out addition and multiply
print(addition)
print(multiply)

# Launch the graph in a session
sess = tf.Session()

# Run the operations
print(sess.run(addition))
print(sess.run(multiply))

```

Comparison to Numpy



Packages

tf estimator

Main Steps

- 1. Create the Model
- 2. Define Target
- 3. Define Loss function and Optimizer
- 4. Define the Session and Initialise Variables
- 5. Train the Model
- 6. Test Trained Model

```

# 1. Create the Model
w = tf.Variable(tf.random_normal([1, 10]), name='w')
b = tf.Variable(tf.zeros([1, 1]), name='b')
x = tf.placeholder(tf.float32, name='x')
y = tf.placeholder(tf.float32, name='y')

# 2. Define Target
x_ = tf.placeholder(tf.float32, name='x_')

# 3. Define Loss function and Optimizer
cross_entropy = tf.nn.softmax_cross_entropy_with_logits(labels=y_, logits=y)
train_step = tf.train.GradientDescentOptimizer(0.5).minimize(cross_entropy)

# 4. Define the Session and Initialise Variables
sess = tf.InteractiveSession()
tf.global_variables_initializer().run()

# 5. Train the Model
for i in range(1000):
    w_init, b_init = w.get_shape().as_list()
    w_init = tf.random_normal([w_init, 10])
    b_init = tf.zeros([1, 1])
    w.assign(w_init)
    b.assign(b_init)

# 6. Test Trained Model
correct_prediction = tf.equal(tf.argmax(y, 1), tf.argmax(x_))
accuracy = tf.reduce_mean(tf.cast(correct_prediction, tf.float32))
print('accuracy: ', sess.run(accuracy, {x_: sess.test_data()}))

```

- tf.estimator.LinearClassifier: Constructs a linear classification model.
- tf.estimator.LinearRegressor: Constructs a linear regression model.
- tf.estimator.DNNClassifier: Construct a neural network classification model.
- tf.estimator.DNNRegressor: Construct a neural network regression model.
- tf.estimator.DNNLinearCombinedClassifier: Construct a neural network and linear combined classification model.
- tf.estimator.DNNLinearCombinedRegressor: Construct a neural network and linear combined regression model.

TensorFlow's high-level machine learning API (tf.estimator) makes it easy to configure, train, and evaluate a variety of machine learning models.

FeatureColumns are the primary way of encoding features for pre-trained tf.learn Estimators.

- Continuous Features: Can be represented by real_valued_column
- Categorical Features: Can be represented by any sparse_column_with_... column (sparse_column_with_keys, sparse_column_with_vocabulary_list, sparse_column_with_hash_buckets, sparse_column_with_interspersed_feature)

When using FeatureColumns with tf.learn models, the type of feature column you should choose depends on the feature type and the model type.

Using a pre-built Logistic Regression Classifier

```

import tensorflow as tf

# Create a variable named 'W'
W = tf.Variable([.3], dtype=tf.float32)

# Create a variable named 'b'
b = tf.Variable([.1], dtype=tf.float32)

# Create a variable named 'x'
x = tf.placeholder(tf.float32)

# Operation to add W and b
addition = tf.add(W, b)

# Operation to multiply addition with x
multiply = tf.multiply(addition, x)

# Print out addition and multiply
print(addition)
print(multiply)

# Launch the graph in a session
sess = tf.Session()

# Run the operations
print(sess.run(addition))
print(sess.run(multiply))

```

This function holds the actual data (features and labels). Features is a python dictionary.

Using the fit function, on the input fit. Notice that the feature columns are fed to the model as arguments.

Using the eval_input_fn defined previously.

```

import tensorflow as tf

# Create a variable named 'W'
W = tf.Variable([.3], dtype=tf.float32)

# Create a variable named 'b'
b = tf.Variable([.1], dtype=tf.float32)

# Create a variable named 'x'
x = tf.placeholder(tf.float32)

# Operation to add W and b
addition = tf.add(W, b)

# Operation to multiply addition with x
multiply = tf.multiply(addition, x)

# Print out addition and multiply
print(addition)
print(multiply)

# Launch the graph in a session
sess = tf.Session()

# Run the operations
print(sess.run(addition))
print(sess.run(multiply))

```

```

import tensorflow as tf

# Create a variable named 'W'
W = tf.Variable([.3], dtype=tf.float32)

# Create a variable named 'b'
b = tf.Variable([.1], dtype=tf.float32)

# Create a variable named 'x'
x = tf.placeholder(tf.float32)

# Operation to add W and b
addition = tf.add(W, b)

# Operation to multiply addition with x
multiply = tf.multiply(addition, x)

# Print out addition and multiply
print(addition)
print(multiply)

# Launch the graph in a session
sess = tf.Session()

# Run the operations
print(sess.run(addition))
print(sess.run(multiply))

```

Comparison to Numpy

```

import tensorflow as tf

# Create a variable named 'W'
W = tf.Variable([.3], dtype=tf.float32)

# Create a variable named 'b'
b = tf.Variable([.1], dtype=tf.float32)

# Create a variable named 'x'
x = tf.placeholder(tf.float32)

# Operation to add W and b
addition = tf.add(W, b)

# Operation to multiply addition with x
multiply = tf.multiply(addition, x)

# Print out addition and multiply
print(addition)
print(multiply)

# Launch the graph in a session
sess = tf.Session()

# Run the operations
print(sess.run(addition))
print(sess.run(multiply))

```

```

import tensorflow as tf

# Create a variable named 'W'
W = tf.Variable([.3], dtype=tf.float32)

# Create a variable named 'b'
b = tf.Variable([.1], dtype=tf.float32)

# Create a variable named 'x'
x = tf.placeholder(tf.float32)

# Operation to add W and b
addition = tf.add(W, b)

# Operation to multiply addition with x
multiply = tf.multiply(addition, x)

# Print out addition and multiply
print(addition)
print(multiply)

# Launch the graph in a session
sess = tf.Session()

# Run the operations
print(sess.run(addition))
print(sess.run(multiply))

```

Comparison to Numpy

```

import tensorflow as tf

# Create a variable named 'W'
W = tf.Variable([.3], dtype=tf.float32)

# Create a variable named 'b'
b = tf.Variable([.1], dtype=tf.float32)

# Create a variable named 'x'
x = tf.placeholder(tf.float32)

# Operation to add W and b
addition = tf.add(W, b)

# Operation to multiply addition with x
multiply = tf.multiply(addition, x)

# Print out addition and multiply
print(addition)
print(multiply)

# Launch the graph in a session
sess = tf.Session()

# Run the operations
print(sess.run(addition))
print(sess.run(multiply))

```

```

import tensorflow as tf

# Create a variable named 'W'
W = tf.Variable([.3], dtype=tf.float32)

# Create a variable named 'b'
b = tf.Variable([.1], dtype=tf.float32)

# Create a variable named 'x'
x = tf.placeholder(tf.float32)

# Operation to add W and b
addition = tf.add(W, b)

# Operation to multiply addition with x
multiply = tf.multiply(addition, x)

# Print out addition and multiply
print(addition)
print(multiply)

# Launch the graph in a session
sess = tf.Session()

# Run the operations
print(sess.run(addition))
print(sess.run(multiply))

```

Comparison to Numpy