

Supervised Learning

Neural Networks

Jonathan Mwaura

Khoury Institute of Computer Science

December 1, 2022

Introduction

Acknowledgements

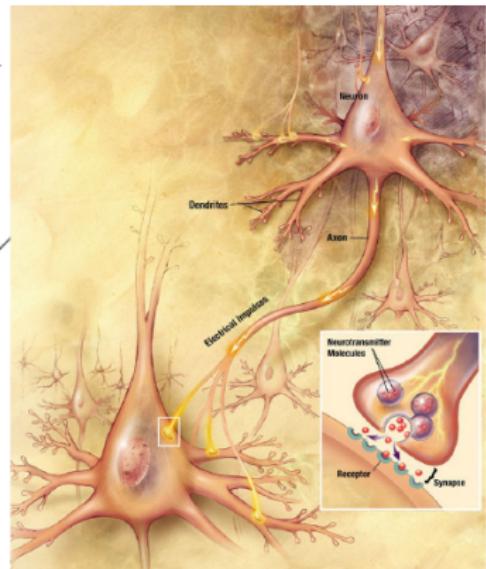
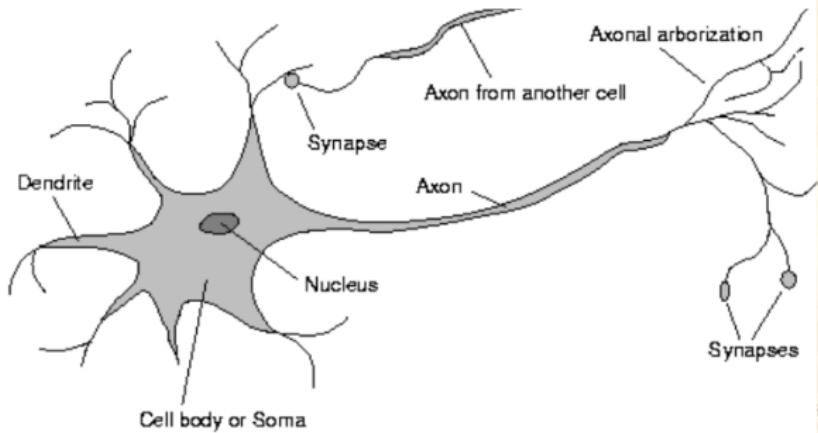
These slides have been adapted from the following Professors:

- 1) Andrew Ng - Stanford
- 2) Eric Eaton - UPenn
- 3) David Sontag - MIT
- 4) Alina Oprea - Northeastern

Neural Function

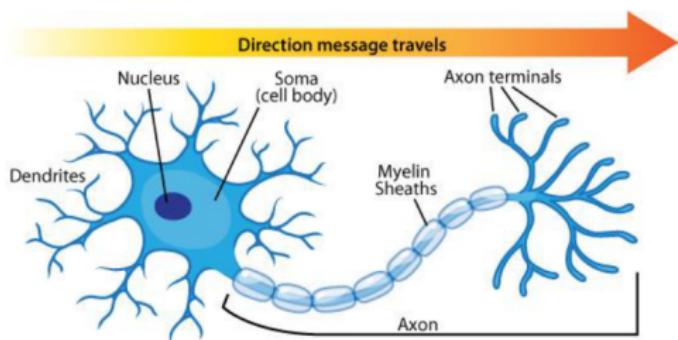
- Brain function (thought) occurs as the result of the firing of **neurons**
- Neurons connect to each other through **synapses**, which propagate **action potential** (electrical impulses) by releasing **neurotransmitters**
 - Synapses can be **excitatory** (potential-increasing) or **inhibitory** (potential-decreasing), and have varying **activation thresholds**
 - Learning occurs as a result of the synapses' **plasticity**: They exhibit long-term changes in connection strength
- There are about 10^{11} neurons and about 10^{14} synapses in the human brain!

Biology of a Neuron

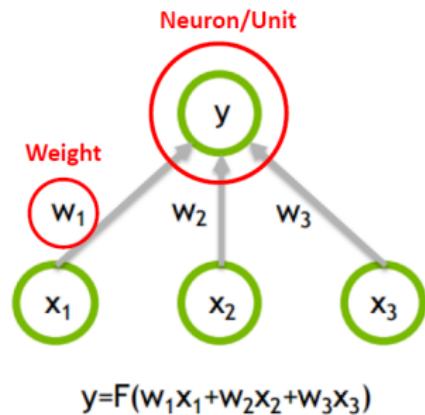


Analogy to Human Brain

Human Brain



Biological Neuron



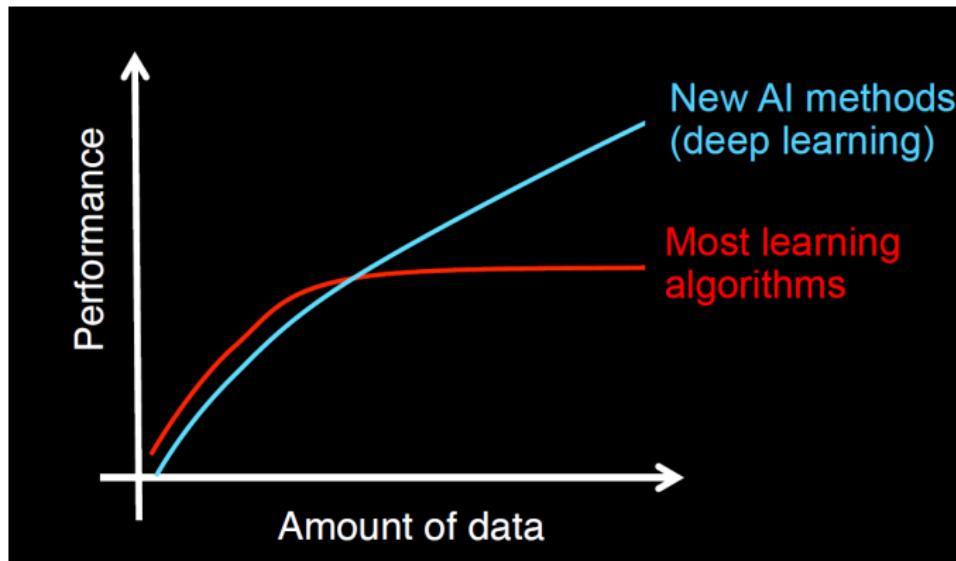
$$F(x) = \max(0, x)$$

Artificial Neuron

Neural Networks

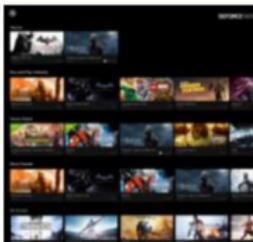
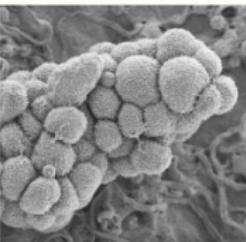
- Origins: Algorithms that try to mimic the brain.
- Very widely used in 80s and early 90s; popularity diminished in late 90s.
- Recent resurgence: State-of-the-art technique for many applications
- Artificial neural networks are not nearly as complex or intricate as the actual brain structure

Performance of Deep Learning



Deep Learning Applications

DEEP LEARNING EVERYWHERE



INTERNET & CLOUD

Image Classification
Speech Recognition
Language Translation
Language Processing
Sentiment Analysis
Recommendation

MEDICINE & BIOLOGY

Cancer Cell Detection
Diabetic Grading
Drug Discovery

MEDIA & ENTERTAINMENT

Video Captioning
Video Search
Real Time Translation

SECURITY & DEFENSE

Face Detection
Video Surveillance
Satellite Imagery

AUTONOMOUS MACHINES

Pedestrian Detection
Lane Tracking
Recognize Traffic Sign

Success stories: Speech recognition

www.technewsworld.com/story/84013.html

40 maps that explain Amazon Web Services Primers | Math n' Pro deeplearning.net/tutorials Deep Learning Tutorials deep learning PHILIPS - Golden Ears Language Technology MyIDCare - Dashboard Other bookmarks

TECHNEWSWORLD

EMERGING TECH

SEARCH

Computing Internet IT Mobile Tech Reviews Security Technology Tech Blog Reader Services

Microsoft AI Beats Humans at Speech Recognition

By Richard Adhikari Oct 20, 2016 11:40 AM PT

Print Email

5 Tweet 25 Share 45 Share 11 Share 0 Share 104



Image: Adobe Stock

How do you feel about Black Friday and Cyber Monday?

- They're great -- I get a lot of bargains!
- The deals are too spread out -- I'd prefer just one day.
- They're a fun way to kick off the holiday season.
- I don't like the commercialization of Thanksgiving Day.
- They're crucial for the retail industry and the economy.
- The deals typically aren't that good.

[Vote to See Results](#)

E-Commerce Times

[Black Friday Shoppers Hungry for New Experiences, New Tech](#)

[Pay TV's Newest Innovation: Giving Users Control](#)

[Apple Celebrates Itsself in \\$300 Coffee Table Tome](#)

[AWS Enjoys Top Perch in IaaS, PaaS Markets](#)

[US Comptroller Gears Up for Blockchain and](#)

Success stories: Machine Translation

The screenshot shows a web browser window with the URL <https://blog.google/products/translate/found-translation-more-accurate-fluent-sentences-google-translate/>. The page title is "Found in translation: More accurate, fluent sentences in Google Translate". The author is Barak Turovsky, Product Lead, Google Translate. The date is NOV 15, 2016. The content discusses the evolution of Google Translate from supporting just a few languages to 103, connecting strangers, reaching across language barriers, and even helping.

← → 🔍 https://blog.google/products/translate/found-translation-more-accurate-fluent-sentences-google-translate/

top 40 maps that explain Amazon Web Services Primers | Math n Pro... deeplearning.net/tutorials Deep Learning Tutor... deep learning PHILIPS - Golden Ears Language Technology MyIDCare - Dashboard Other bookmarks

G The Keyword Latest Stories Product News Topics

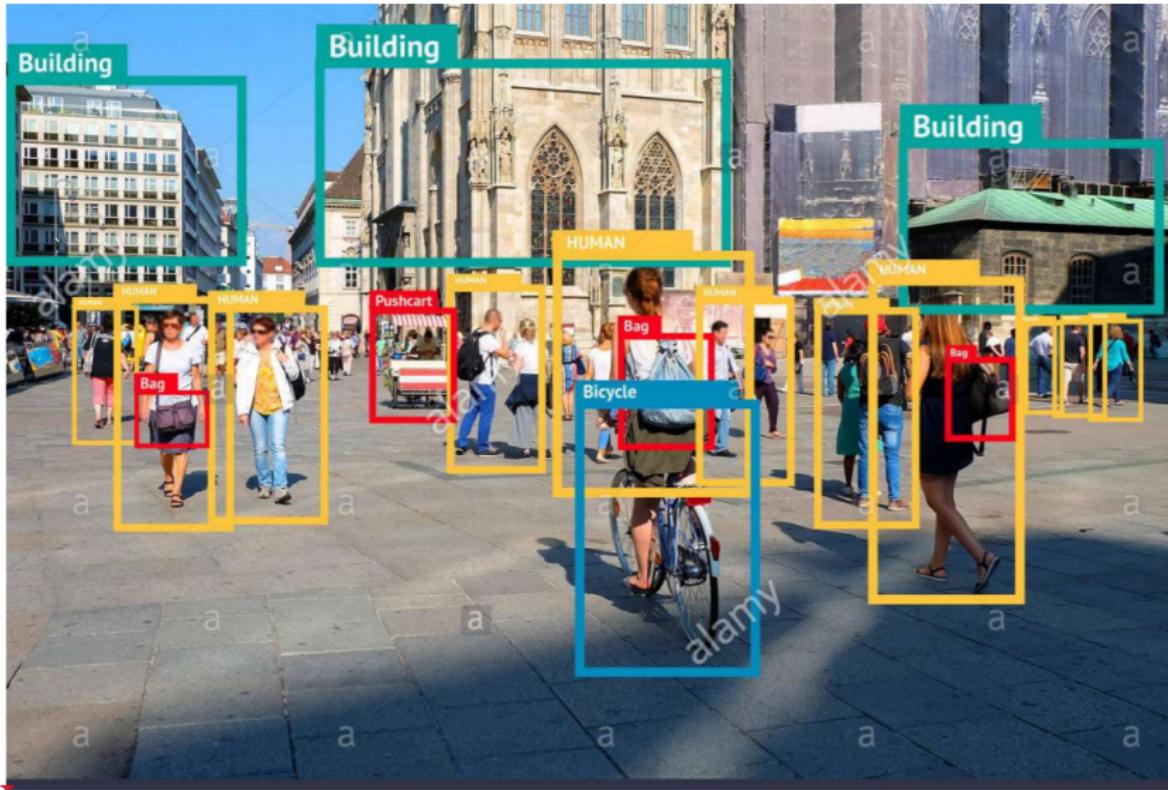
TRANSLATE NOV 15, 2016

Found in translation: More accurate, fluent sentences in Google Translate

Barak Turovsky
PRODUCT LEAD, GOOGLE TRANSLATE

In 10 years, Google Translate has gone from supporting just a few languages to 103, connecting strangers, reaching across language barriers and even helping

Success stories: Image segmentation



Success stories: Image captioning

The screenshot shows a web browser window with the URL <http://cs.stanford.edu/people/karpathy/deepimagesent/>. The page displays eight image caption pairs, each consisting of a small image thumbnail followed by a caption generated by the system.

- "man in black shirt is playing guitar."
- "construction worker in orange safety vest is working on road."
- "two young girls are playing with lego toy."
- "boy is doing backflip on wakeboard."
- "girl in pink dress is jumping in air."
- "black and white dog jumps over bar."
- "young girl in pink shirt is swinging on swing."
- "man in blue wetsuit is surfing on wave."

References

- Deep Learning books
 - <https://d2l.ai/> (D2L)
 - <https://www.deeplearningbook.org/> (advanced)
- Stanford notes on deep learning
 - http://cs229.stanford.edu/notes/cs229-notes-deep_learning.pdf
- History of Deep Learning
 - https://beamandrew.github.io/deeplearning/2017/02/23/deep_learning_101_part1.html

Example

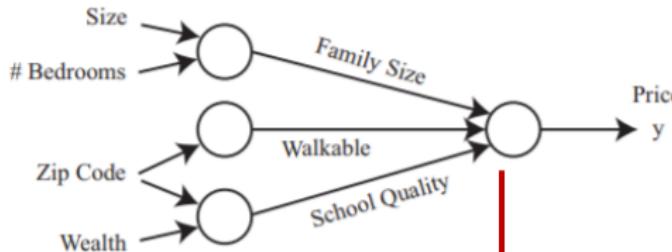
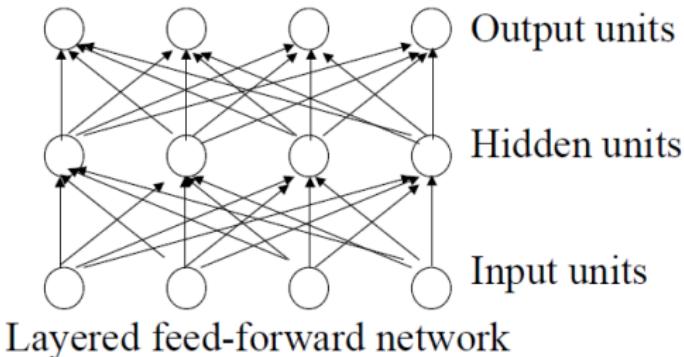


Figure 2: Diagram of a small neural network for predicting housing prices.

- Provide as input only training data: input and label
- Neural Networks automatically learn intermediate features!

Neural Networks



Training labels

Learned
during training

Training data

- Neural networks are made up of **nodes** or **units**, connected by **links**
- Each link has an associated **weight** and **activation level**
- Each node has an **input function** (typically summing over weighted inputs), an **activation function**, and an **output**

Recall: The Perceptron

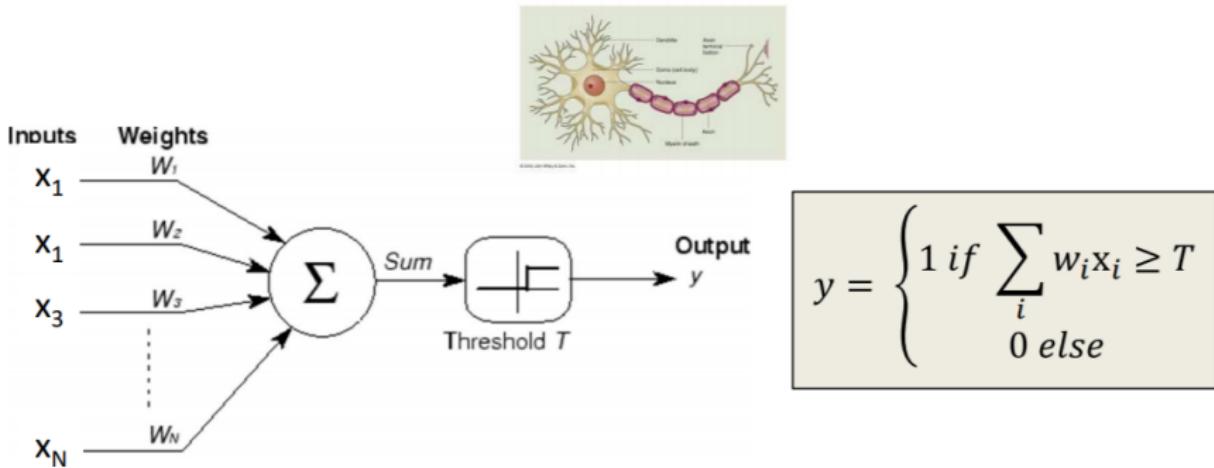
$$h(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^T \mathbf{x}) \quad \text{where} \quad \text{sign}(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

- The perceptron uses the following update rule each time it receives a new training instance $(\mathbf{x}^{(i)}, y^{(i)})$

$$\theta_j \leftarrow \theta_j - \frac{1}{2} \underbrace{\left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)}}_{\text{either 2 or -2}}$$

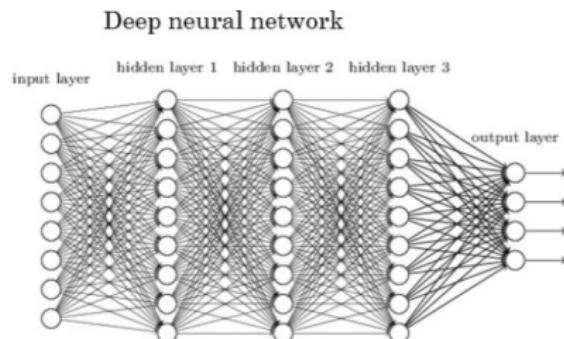
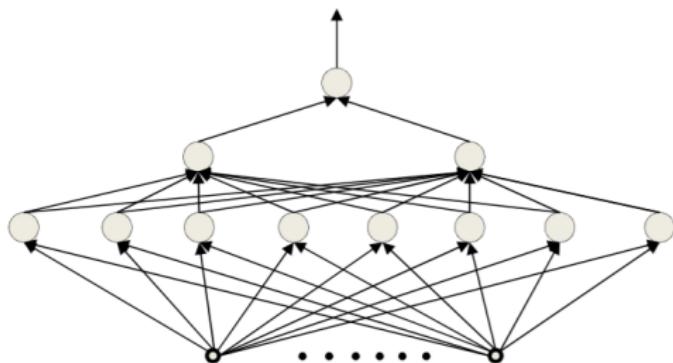
- If the prediction matches the label, make no change
- Otherwise, adjust $\boldsymbol{\theta}$

Perceptron



- A threshold unit
 - “Fires” if the weighted sum of inputs exceeds a threshold

Multi-Layer Perceptron



- A network of perceptrons
 - Generally “layered”

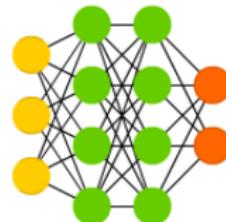


Neural Network Architectures

Feed-Forward Networks

- Neurons from each layer connect to neurons from next layer

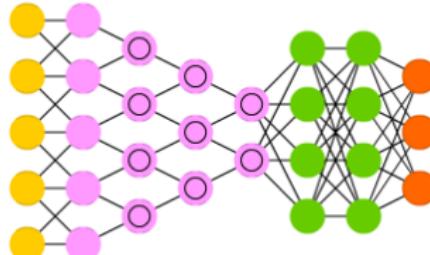
Deep Feed Forward (DFF)



Convolutional Networks

- Includes convolution layer for feature reduction
- Learns hierarchical representations

Deep Convolutional Network (DCN)

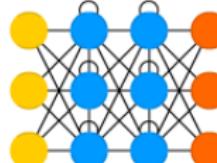


Recurrent Networks

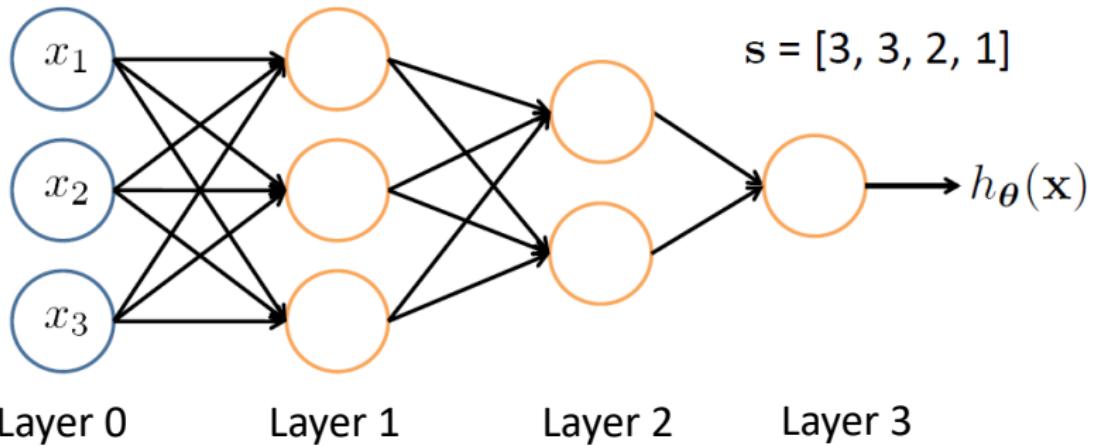
- Keep hidden state
- Have cycles in

N The Roux Institute
at Northeastern University

Recurrent Neural Network (RNN)



Feed-Forward Networks



L denotes the number of layers

$s \in \mathbb{N}^{+L}$ contains the numbers of nodes at each layer

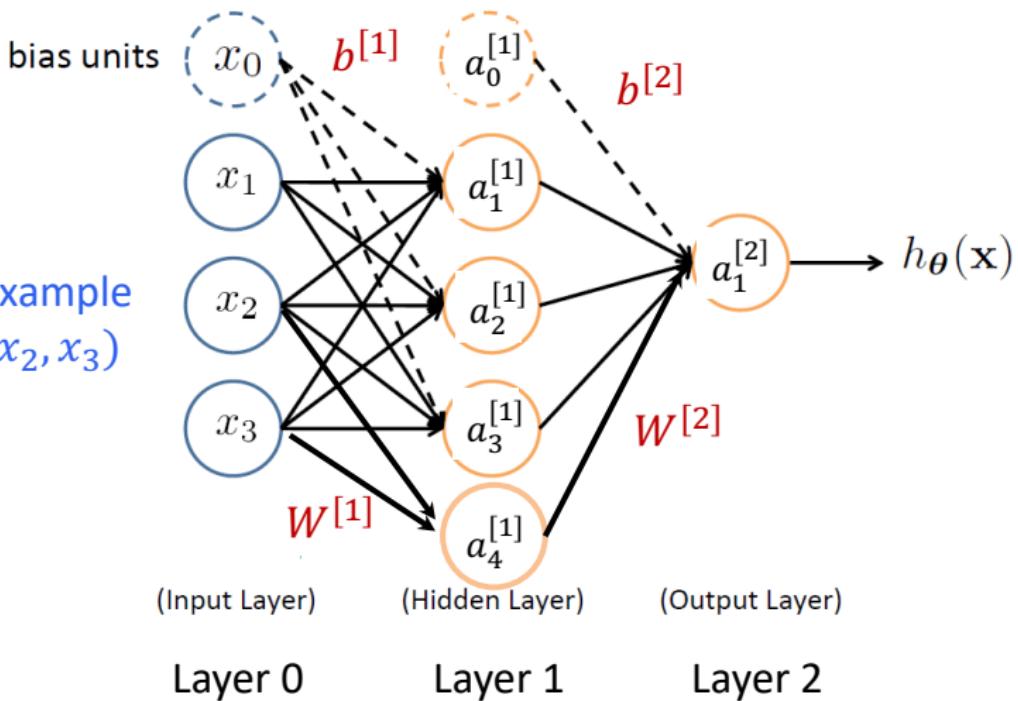
- Not counting bias units
- Typically, $s_0 = d$ (# input features) and $s_{L-1} = K$ (# classes)

Feed-Forward NN

- Hyper-parameters
 - Number of layers
 - Architecture (how layers are connected)
 - Number of hidden units per layer
 - Number of units in output layer
 - Activation functions
- Other
 - Initialization
 - Regularization

Feed-Forward Neural Network

Training example
 $x = (x_1, x_2, x_3)$



Vectorization

$$z_1^{[1]} = W_1^{[1]T} x + b_1^{[1]} \quad \text{and} \quad a_1^{[1]} = g(z_1^{[1]})$$

 \vdots \vdots \vdots

$$z_4^{[1]} = W_4^{[1]T} x + b_4^{[1]} \quad \text{and} \quad a_4^{[1]} = g(z_4^{[1]})$$

$$\underbrace{\begin{bmatrix} z_1^{[1]} \\ \vdots \\ z_4^{[1]} \end{bmatrix}}_{z^{[1]} \in \mathbb{R}^{4 \times 1}} = \underbrace{\begin{bmatrix} -W_1^{[1]T} \\ -W_2^{[1]T} \\ \vdots \\ -W_4^{[1]T} \end{bmatrix}}_{W^{[1]} \in \mathbb{R}^{4 \times 3}} \underbrace{x}_{x \in \mathbb{R}^{3 \times 1}} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ \vdots \\ b_4^{[1]} \end{bmatrix}}_{b^{[1]} \in \mathbb{R}^{4 \times 1}}$$

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = g(z^{[1]})$$

Vectorization

Output layer

$$z_1^{[2]} = W_1^{[2]T} a^{[1]} + b_1^{[2]} \quad \text{and} \quad a_1^{[2]} = g(z_1^{[2]})$$

- - - - -

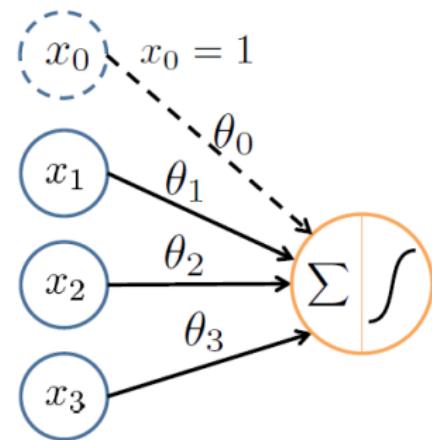
$$\underbrace{z^{[2]}}_{1 \times 1} = \underbrace{W^{[2]}}_{1 \times 4} \underbrace{a^{[1]}}_{4 \times 1} + \underbrace{b^{[2]}}_{1 \times 1} \quad \text{and} \quad \underbrace{a^{[2]}}_{1 \times 1} = g(\underbrace{z^{[2]}}_{1 \times 1})$$

Hidden Units

- Layer 1
 - First hidden unit:
 - Linear: $z_1^{[1]} = W_1^{[1] T} x + b_1^{[1]}$
 - Non-linear: $a_1^{[1]} = g(z_1^{[1]})$
 - ...
 - Fourth hidden unit:
 - Linear: $z_4^{[1]} = W_4^{[1] T} x + b_4^{[1]}$
 - Non-linear: $a_4^{[1]} = g(z_4^{[1]})$
- Terminology
 - $a_i^{[j]}$ - Activation of unit i in layer j
 - g - Activation function
 - $W^{[j]}$ - Weight vector controlling mapping from layer $j-1$ to j
 - $b^{[j]}$ - Bias vector from layer $j-1$ to j

Logistic Unit: A simple NN

“bias unit”

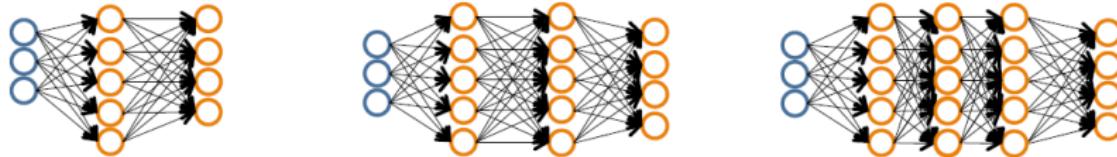


$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$
$$h_{\boldsymbol{\theta}}(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x}) = g(\mathbf{w}^T \mathbf{x} + b)$$
$$= \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}}$$

Sigmoid (logistic) activation function: $g(z) = \frac{1}{1 + e^{-z}}$

How to pick architecture?

Pick a network architecture (connectivity pattern between nodes)



- # input units = # of features in dataset
- # output units = # classes

Reasonable default: 1 hidden layer

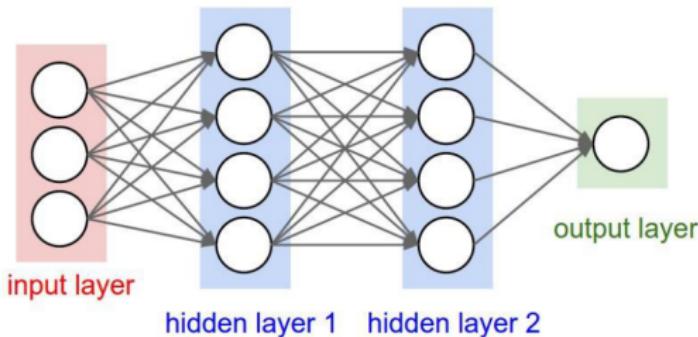
- or if >1 hidden layer, have same # hidden units in every layer (usually the more the better)

Training Neural Networks

- Input training dataset D
 - Number of features: d
 - Labels from K classes
- First layer has $d+1$ units (one per feature and bias)
- Output layer has K units
- Training procedure determines parameters that optimize loss function
 - Backpropagation
 - Learn optimal $W^{[i]}, b^{[i]}$ at layer i
- Testing done by forward propagation

Forward Propagation

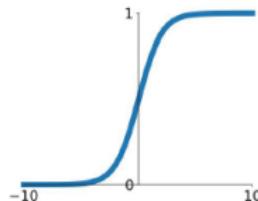
- The input neurons first receive the data features of the object. After processing the data, they send their output to the first hidden layer.
- The hidden layer processes this output and sends the results to the next hidden layer.
- This continues until the data reaches the final output layer, where the output value determines the object's classification.
- This entire process is known as **Forward Propagation**, or **Forward prop**.



Activation Functions

Sigmoid

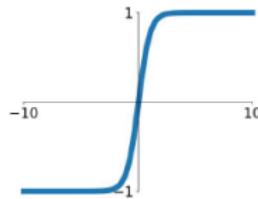
$$\sigma(x) = \frac{1}{1+e^{-x}}$$



Binary
Classification

tanh

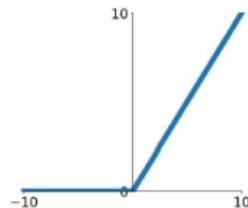
$$\tanh(x)$$



Regression

ReLU

$$\max(0, x)$$



Intermediary
layers

Why Non-Linear Activations?

- Assume g is linear: $g(z) = Uz$
 - At layer 1: $z^{[1]} = W^{[1]}x + b^{[1]}$
 - $a^{[1]} = g(z^{[1]}) = Uz^{[1]} = UW^{[1]}x + Ub^{[1]}$
- Layer 2:
 - $a^{[2]} = g(z^{[2]}) = Uz^{[2]} = UW^{[2]}a^{[1]} + Ub^{[2]} = UW^{[2]}UW^{[1]}x + UW^{[2]}Ub^{[1]} + Ub^{[2]}$
- Last layer
 - Output is linear in input!
 - Then NN will only learn linear functions

Multiple Output Units: One-vs-Rest



Pedestrian



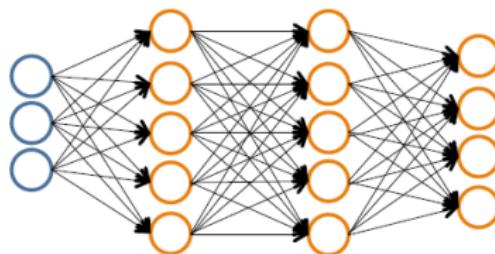
Car



Motorcycle



Truck



$$h_{\Theta}(\mathbf{x}) \in \mathbb{R}^K$$

We want:

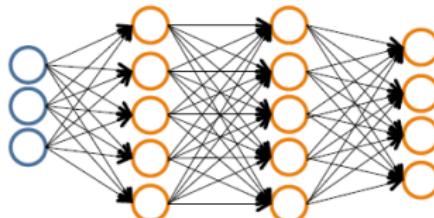
$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \text{when pedestrian}$$

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \text{when car}$$

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \text{when motorcycle}$$

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \quad \text{when truck}$$

Multiple Output Units: One-vs-Rest



$$h_{\Theta}(\mathbf{x}) \in \mathbb{R}^K$$

We want:

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

when pedestrian

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

when car

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

when motorcycle

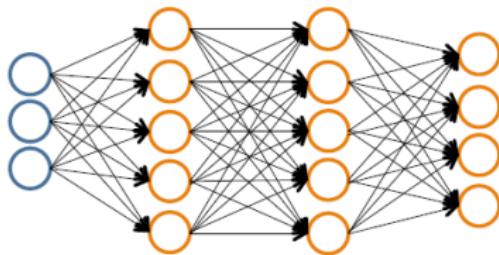
$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

when truck

- Given $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- Must convert labels to 1-of- K representation

– e.g., $y_i = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ when motorcycle, $y_i = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ when car, etc.

Neural Network Classification



Given:

$$\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$$

$\mathbf{s} \in \mathbb{N}^{+L}$ contains # nodes at each layer

- $s_0 = d$ (# features)

Binary classification

$y = 0 \text{ or } 1$

1 output unit ($s_{L-1} = 1$)

Sigmoid

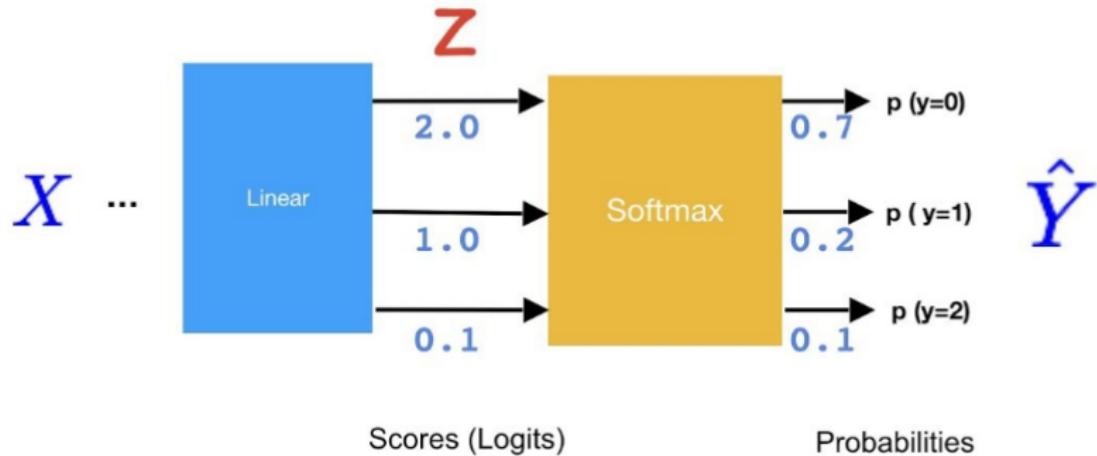
Multi-class classification (K classes)

$\mathbf{y} \in \mathbb{R}^K$ e.g. $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

K output units ($s_{L-1} = K$)

Softmax

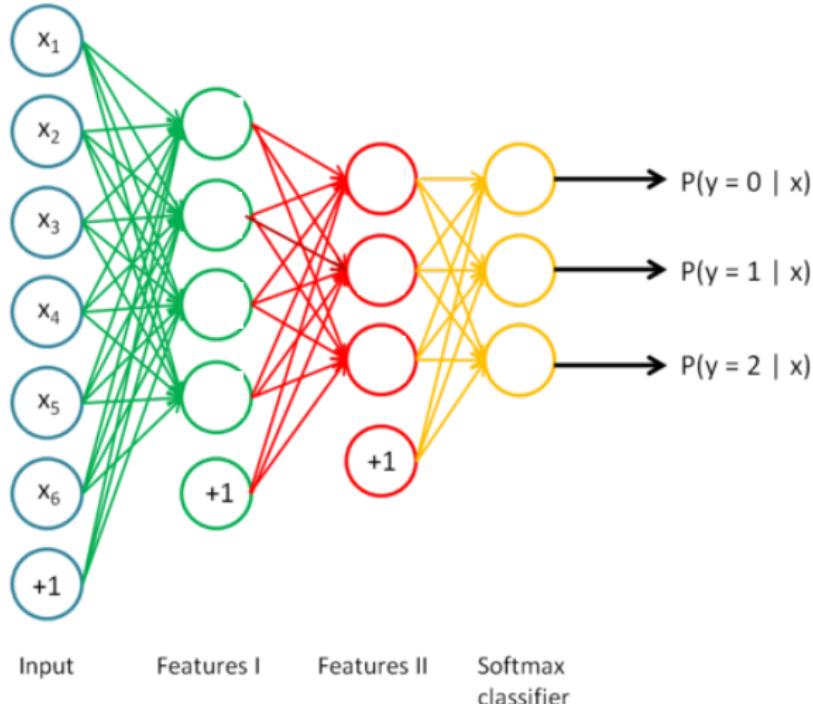
Softmax classifier



$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$

- Predict the class with highest probability
- Generalization of sigmoid/logistic regression to multi-class

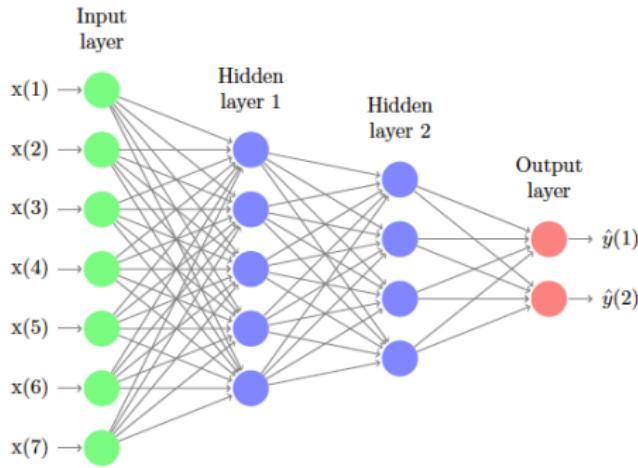
Multi-class classification



Review Feed-Forward Neural Networks

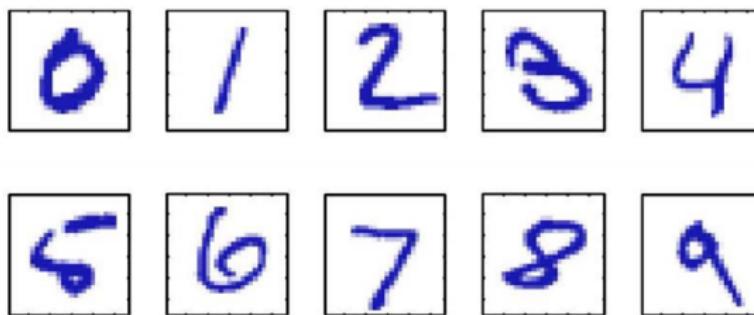
- Simplest architecture of NN
- Neurons from one layer are connected to neurons from next layer
 - Input layer has feature space dimension
 - Output layer has number of classes
 - Hidden layers use linear operations, followed by non-linear activation function
 - Multi-Layer Perceptron (MLP): fully connected layers
- Activation functions
 - Sigmoid for binary classification in last layer
 - Softmax for multi-class classification in last layer
 - ReLU for hidden layers
- Forward propagation is the computation of the network output given an input
- Back propagation is the training of a network
 - Determine weights and biases at every layer

FFNN Architectures



- Input and Output Layers are completely specified by the problem domain
- In the Hidden Layers, number of neurons in Layer $i+1$ is always smaller than number of neurons in Layer i

MNIST: Handwritten digit recognition



Images are 28 x 28 pixels

Represent input image as a vector $x \in \mathbb{R}^{784}$

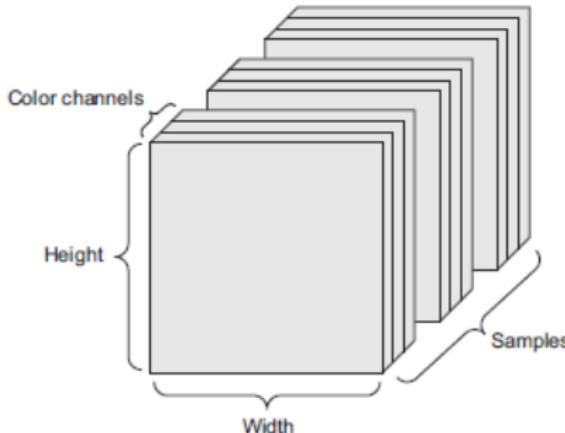
Learn a classifier $f(x)$ such that,

$$f : x \rightarrow \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

Predict the digit
Multi-class classifier

Image Representation

- Image is 3D “tensor”: height, width, color channel (RGB)
- Black-and-white images are 2D matrices: height, width
 - Each value is pixel intensity



Lab – Feed Forward NN

```
import time
import numpy as np
from keras.utils import np_utils
import keras.callbacks as cb
from keras.models import Sequential
from keras.layers.core import Dense, Dropout, Activation
from keras.optimizers import RMSprop
from keras.datasets import mnist

import matplotlib
matplotlib.use('agg')
import matplotlib.pyplot as plt
```

Import modules

```
def load_data():
    print("Loading data")
    (X_train, y_train), (X_test, y_test) = mnist.load_data()

    X_train = X_train.astype('float32')
    X_test = X_test.astype('float32')

    # Normalize
    X_train /= 255
    X_test /= 255

    y_train = np_utils.to_categorical(y_train, 10)
    y_test = np_utils.to_categorical(y_test, 10)

    X_train = np.reshape(X_train, (60000, 784))
    X_test = np.reshape(X_test, (10000, 784))

    print("Data Loaded")
    return [X_train, X_test, y_train, y_test]
```

Load MNIST data
Processing

Vector
representation

Neural Network Architecture

```
def init_model():
    start_time = time.time()

    print("Compiling Model")
    model = Sequential()
    model.add(Dense(10, input_dim=784))
    model.add(Activation('relu'))           ← 10 hidden units
                                            ← ReLU activation

    model.add(Dense(10))
    model.add(Activation('softmax'))       ← Output Layer
                                            ← Softmax activation

    rms = RMSprop()
    model.compile(loss='categorical_crossentropy', optimizer=rms, metrics=[ 'accuracy'])

    print("Model finished"+format(time.time() - start_time))
    return model

↓                               ↓
Loss function      Optimizer
```

Feed-Forward Neural Network Architecture

- 1 Hidden Layer (“Dense” or Fully Connected)
- 10 neurons
- Output layer uses softmax activation

Train and evaluate

```
>def run_network(data=None, model=None, epochs=10, batch=256):
try:
    start_time = time.time()
    if data is None:
        X_train, X_test, y_train, y_test = load_data()
    else:
        X_train, X_test, y_train, y_test = data

    if model is None:
        model = init_model()

    print("Training model")
    history = model.fit(X_train, y_train, nb_epoch=epochs, batch_size=batch,
                         validation_data=(X_test, y_test), verbose=2)

    print("Training duration:"+format(time.time() - start_time))
    score = model.evaluate(X_test, y_test, batch_size=16)

    print("\nNetwork's test Loss and accuracy:"+format(score))
    return model, history
except KeyboardInterrupt:
    print("KeyboardInterrupt")
    return model, history
```

Training/testing results

```
2s - loss: 0.9114 - acc: 0.7649 - val_loss: 0.4499 - val_acc: 0.8819
Epoch 2/10
0s - loss: 0.3935 - acc: 0.8907 - val_loss: 0.3378 - val_acc: 0.9049
Epoch 3/10
0s - loss: 0.3296 - acc: 0.9063 - val_loss: 0.3042 - val_acc: 0.9128
Epoch 4/10
0s - loss: 0.3036 - acc: 0.9132 - val_loss: 0.2889 - val_acc: 0.9181
Epoch 5/10
0s - loss: 0.2888 - acc: 0.9189 - val_loss: 0.2874 - val_acc: 0.9185
Epoch 6/10
0s - loss: 0.2785 - acc: 0.9210 - val_loss: 0.2703 - val_acc: 0.9257
Epoch 7/10
0s - loss: 0.2705 - acc: 0.9241 - val_loss: 0.2718 - val_acc: 0.9239
Epoch 8/10
0s - loss: 0.2649 - acc: 0.9257 - val_loss: 0.2694 - val_acc: 0.9240
Epoch 9/10
0s - loss: 0.2601 - acc: 0.9264 - val_loss: 0.2616 - val_acc: 0.9261
Epoch 10/10
0s - loss: 0.2561 - acc: 0.9277 - val_loss: 0.2607 - val_acc: 0.9274
Training duration:10.31288456916809
 9840/10000 [=====>.] - ETA: 0s
Network's test loss and accuracy:[0.26067940444946291, 0.9274]
```

Epoch Output

Metrics

- Loss
- Accuracy

Reported on both training and validation

Changing Number of Neurons

```
def init_model():
    start_time = time.time()

    print("Compiling Model")
    model = Sequential()
    model.add(Dense(500, input_dim=784)) → 500 hidden units
    model.add(Activation('relu'))

    model.add(Dense(10))
    model.add(Activation('softmax'))

    rms = RMSprop()
    model.compile(loss='categorical_crossentropy', optimizer=rms, metrics=['accuracy'])

    print("Model finished"+format(time.time() - start_time))
    return model
```

```
2s - loss: 0.3169 - acc: 0.9088 - val_loss: 0.1652 - val_acc: 0.9502
Epoch 2/10
0s - loss: 0.1277 - acc: 0.9626 - val_loss: 0.1071 - val_acc: 0.9679
Epoch 3/10
0s - loss: 0.0847 - acc: 0.9749 - val_loss: 0.0861 - val_acc: 0.9731
Epoch 4/10
0s - loss: 0.0607 - acc: 0.9822 - val_loss: 0.0746 - val_acc: 0.9767
Epoch 5/10
0s - loss: 0.0471 - acc: 0.9863 - val_loss: 0.0655 - val_acc: 0.9796
Epoch 6/10
0s - loss: 0.0359 - acc: 0.9895 - val_loss: 0.0636 - val_acc: 0.9813
Epoch 7/10
0s - loss: 0.0280 - acc: 0.9920 - val_loss: 0.0599 - val_acc: 0.9810
Epoch 8/10
0s - loss: 0.0223 - acc: 0.9937 - val_loss: 0.0678 - val_acc: 0.9795
Epoch 9/10
0s - loss: 0.0174 - acc: 0.9952 - val_loss: 0.0607 - val_acc: 0.9815
Epoch 10/10
0s - loss: 0.0134 - acc: 0.9964 - val_loss: 0.0672 - val_acc: 0.9806
Training duration:10.458189249038696
9456/10000 [=====>...] - ETA: 0s
Network's test loss and accuracy:[0.067179036217656543, 0.9806000000000003]
```

Two Layers

```
#def init_model():
    start_time = time.time()

    print("Compiling Model")
    model = Sequential()

    # Hidden Layer 1
    model.add(Dense(500, input_dim=784)) → Layer 1
    model.add(Activation('relu'))

    # Hidden Layer 2
    model.add(Dense(300)) → Layer 2
    model.add(Activation('relu'))

    model.add(Dense(10))
    model.add(Activation('softmax')) → Output Softmax Layer

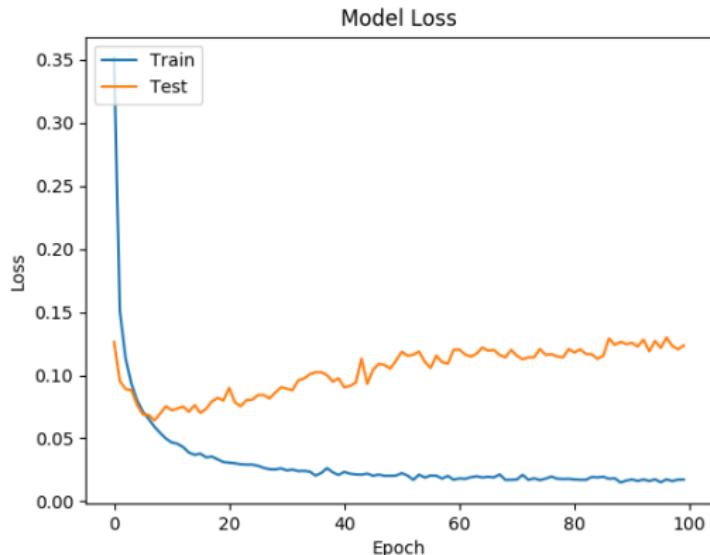
    rms = RMSprop()
    model.compile(loss='categorical_crossentropy', optimizer=rms, metrics=['accuracy'])

    print("Model finished"+format(time.time() - start_time))
    return model
```

```
2s - loss: 0.2800 - acc: 0.9132 - val_loss: 0.1821 - val_acc: 0.9409
Epoch 2/10
1s - loss: 0.0974 - acc: 0.9699 - val_loss: 0.0951 - val_acc: 0.9703
Epoch 3/10
0s - loss: 0.0616 - acc: 0.9803 - val_loss: 0.0843 - val_acc: 0.9754
Epoch 4/10
0s - loss: 0.0429 - acc: 0.9862 - val_loss: 0.0670 - val_acc: 0.9809
Epoch 5/10
0s - loss: 0.0303 - acc: 0.9904 - val_loss: 0.0820 - val_acc: 0.9777
Epoch 6/10
0s - loss: 0.0233 - acc: 0.9922 - val_loss: 0.0794 - val_acc: 0.9783
Epoch 7/10
0s - loss: 0.0180 - acc: 0.9941 - val_loss: 0.0866 - val_acc: 0.9792
Epoch 8/10
0s - loss: 0.0137 - acc: 0.9956 - val_loss: 0.0788 - val_acc: 0.9814
Epoch 9/10
0s - loss: 0.0116 - acc: 0.9963 - val_loss: 0.0901 - val_acc: 0.9795
Epoch 10/10
1s - loss: 0.0100 - acc: 0.9966 - val_loss: 0.0812 - val_acc: 0.9827
Training duration:11.816290140151978
9744/10000 [=====>.] - ETA: 0s
```

Monitor Loss

```
def plot_losses(history):
    plt.plot(history.history['Loss'])
    plt.plot(history.history['val_Loss'])
    plt.title('Model Loss')
    plt.ylabel('Loss')
    plt.xlabel('Epoch')
    plt.legend(['Train', 'Test'], loc='upper left')
    plt.show()
    plt.savefig('output.png')
```



Model Parameters

```
model.summary()
```

```
Using TensorFlow backend.  
Loading data  
Data loaded  
Compiling Model  


| Layer (type)              | Output Shape | Param # |
|---------------------------|--------------|---------|
| dense_1 (Dense)           | (None, 500)  | 392500  |
| activation_1 (Activation) | (None, 500)  | 0       |
| dense_2 (Dense)           | (None, 300)  | 150300  |
| activation_2 (Activation) | (None, 300)  | 0       |
| dense_3 (Dense)           | (None, 10)   | 3010    |
| activation_3 (Activation) | (None, 10)   | 0       |

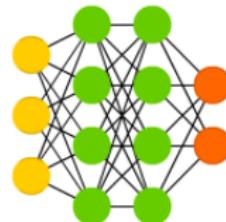
  
Total params: 545,810  
Trainable params: 545,810  
Non-trainable params: 0
```

Neural Network Architectures

Feed-Forward Networks

- Neurons from each layer connect to neurons from next layer

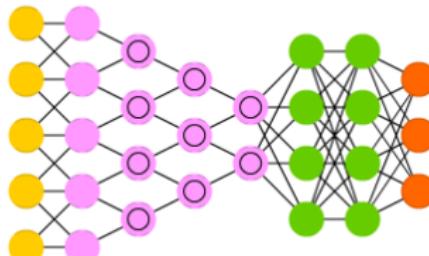
Deep Feed Forward (DFF)



Convolutional Networks

- Includes convolution layer for feature reduction
- Learns hierarchical representations

Deep Convolutional Network (DCN)



Recurrent Networks

- Keep hidden state
- Have cycles in

N The Roux Institute
at Northeastern University

Recurrent Neural Network (RNN)

