

# Machine Learning: Overview

1 / 111

## Speaker Qualifications

- Specialize in next-generation technologies
- Author of O'Reilly Videos on Hypermedia, Linking Data, Security and Encryption
- Author of 'Resource-Oriented Architecture Patterns for Webs of Data'
- Teaches and speaks internationally about REST, Semantic Web, Data Science, Machine Learning, GPU Computing, Security, Visualization, Architecture
- Worked in Defense, Finance, Retail, Hospitality, Video Game, Health Care, Telecommunications and Publishing Industries
- International Pop Recording Artist

2 / 111

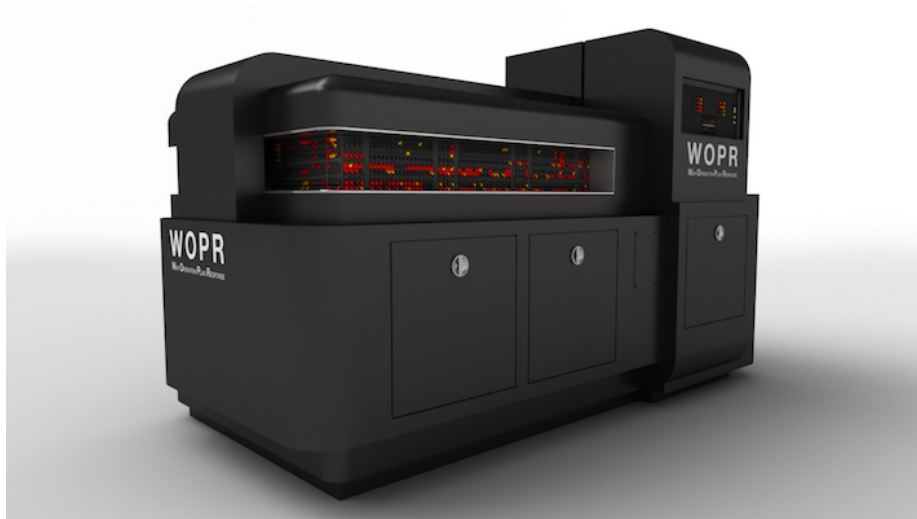
# Agenda

- Introduction
- Algorithms
- Neural Networks

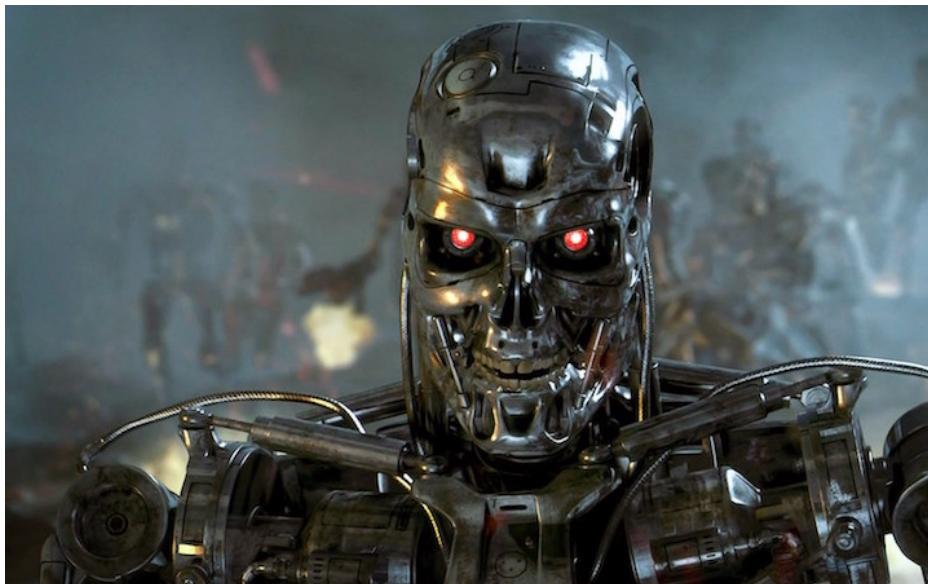
3 / 111

## Introduction

4 / 111



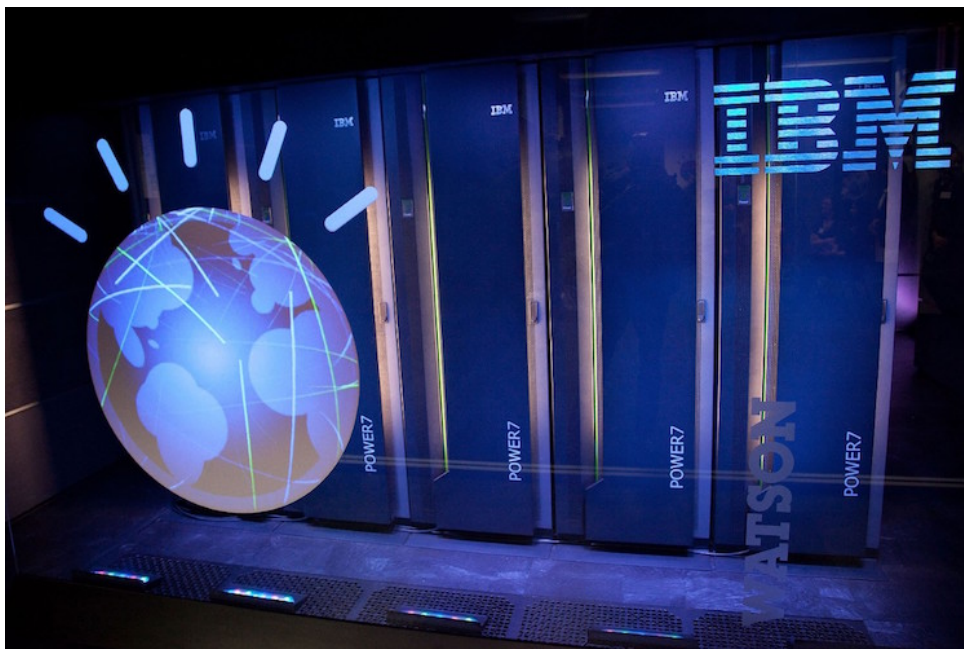
5 / 111



6 / 111

<https://youtu.be/rVlhMGQgDkY?t=85>

7 / 111



8 / 111

"The term machine learning refers to the automated detection of meaningful patterns in data."

*Source: Shavel-Shwartz and Ben-David, "Understanding Machine Learning: From Theory to Algorithms"*

9 / 111

## Advancement of Machine Learning

10 / 111

# Advancement of Machine Learning

- Stock market prediction in the 1980s

11 / 111

# Advancement of Machine Learning

- Stock market prediction in the 1980s
- Mining corporate databases in the 1990s (direct marketing, CRM, credit scoring, fraud detection)

12 / 111

# Advancement of Machine Learning

- Stock market prediction in the 1980s
- Mining corporate databases in the 1990s (direct marketing, CRM, credit scoring, fraud detection)
- E-commerce (personalization, click analysis)

13 / 111

# Advancement of Machine Learning

- Stock market prediction in the 1980s
- Mining corporate databases in the 1990s (direct marketing, CRM, credit scoring, fraud detection)
- E-commerce (personalization, click analysis)
- 9/11 brought interest to applying ML to fighting terror

14 / 111

# Advancement of Machine Learning

- Stock market prediction in the 1980s
- Mining corporate databases in the 1990s (direct marketing, CRM, credit scoring, fraud detection)
- E-commerce (personalization, click analysis)
- 9/11 brought interest to applying ML to fighting terror
- Web 2.0 (social networks, sentiment analysis, etc.)

15 / 111

# Advancement of Machine Learning

- Stock market prediction in the 1980s
- Mining corporate databases in the 1990s (direct marketing, CRM, credit scoring, fraud detection)
- E-commerce (personalization, click analysis)
- 9/11 brought interest to applying ML to fighting terror
- Web 2.0 (social networks, sentiment analysis, etc.)
- Science (molecular biologists and astronomers were early adopters)

16 / 111



# Advancement of Machine Learning

- Stock market prediction in the 1980s
- Mining corporate databases in the 1990s (direct marketing, CRM, credit scoring, fraud detection)
- E-commerce (personalization, click analysis)
- 9/11 brought interest to applying ML to fighting terror
- Web 2.0 (social networks, sentiment analysis, etc.)
- Science (molecular biologists and astronomers were early adopters)
- Housing bust freed up a lot of talent

17 / 111

# Advancement of Machine Learning

- Stock market prediction in the 1980s
- Mining corporate databases in the 1990s (direct marketing, CRM, credit scoring, fraud detection)
- E-commerce (personalization, click analysis)
- 9/11 brought interest to applying ML to fighting terror
- Web 2.0 (social networks, sentiment analysis, etc.)
- Science (molecular biologists and astronomers were early adopters)
- Housing bust freed up a lot of talent
- Big Data

18 / 111



David Hume

19 / 111

## Inductivist Turkey



CC BY 2.0, <https://www.flickr.com/photos/29311691@N05/3577002803>

20 / 111



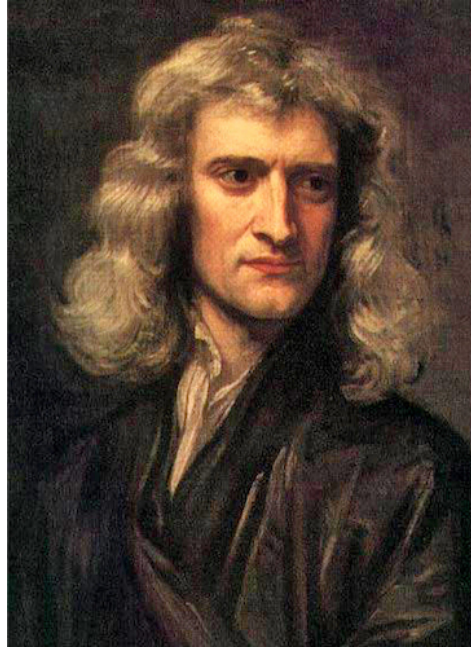
Tycho Brahe

21 / 111



Johannes Kepler

22 / 111

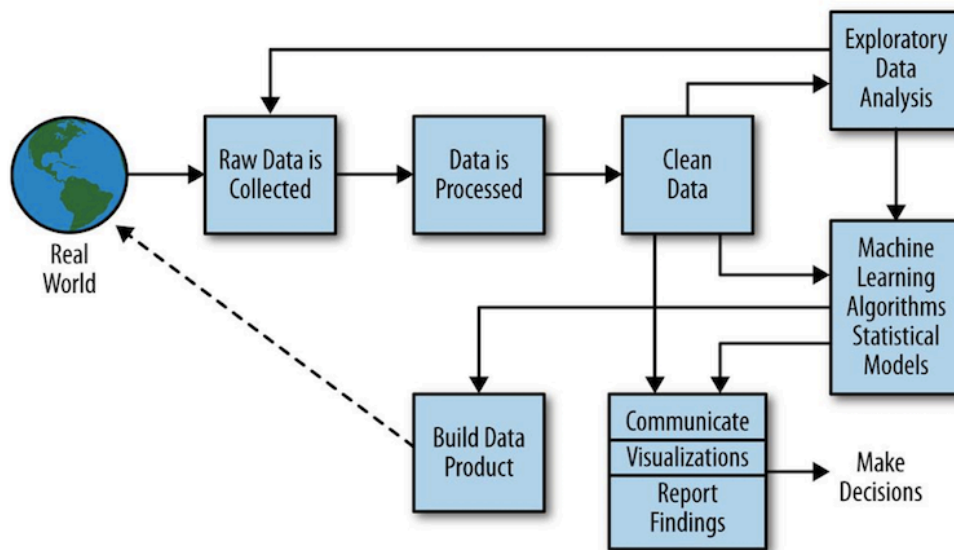


Isaac Newton

23 / 111

# Algorithms

24 / 111

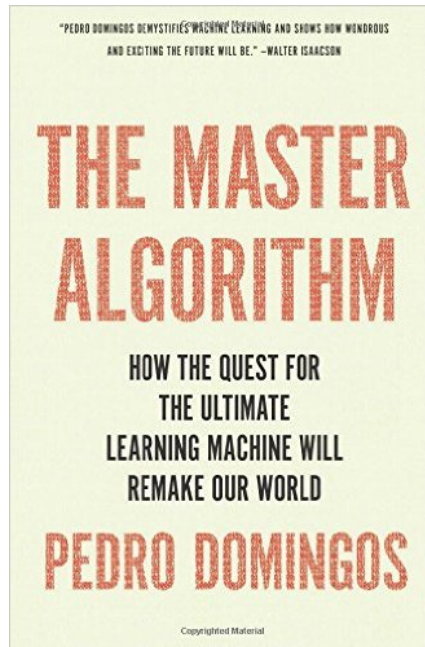


O'Neil and Schutt, "Doing Data Science"

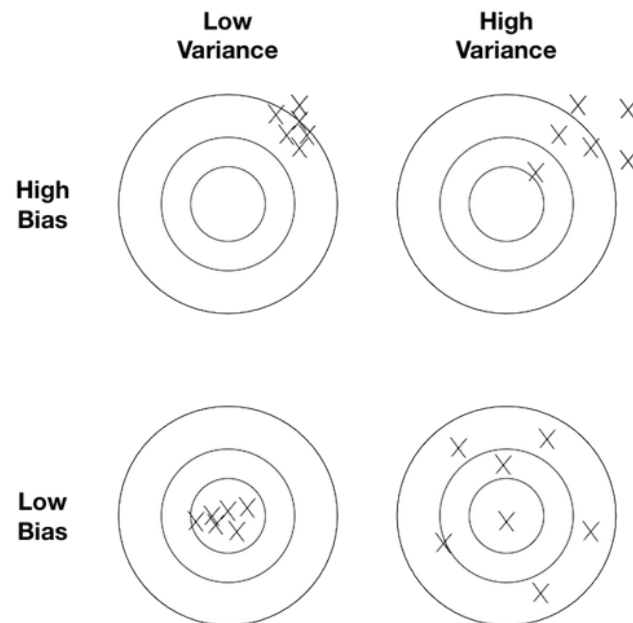
25 / 111



26 / 111



27 / 111



Domingos, "Master Algorithm : : How the Quest for the Ultimate Learning Machine Will Remake Our World"

28 / 111

<b>Tribe</b>	<b>Approach</b>	<b>Master Algorithm</b>
Evolutionaries	Natural Selection/Deriving Learning Structure	Genetic Programming
Connectionists	Reverse engineer the brain	Backpropagation
Symbolists	Symbol manipulation from initial knowledge	Inverse Deduction
Bayesians	Managing uncertainty, noisy, incomplete and contradictory information	Probabilistic Inference (Bayes' Theorem)
Analogizers	Recognizing similarities	Support Vector Machines

29 / 111

## Sample Algorithms

30 / 111

# Sample Algorithms

- Linear Regression

31 / 111

# Sample Algorithms

- Linear Regression
- k-Nearest Neighbors

32 / 111



# Sample Algorithms

- Linear Regression
- k-Nearest Neighbors
- Naive Bayes

33 / 111

# Sample Algorithms

- Linear Regression
- k-Nearest Neighbors
- Naive Bayes
- Decision Trees

34 / 111

# Sample Algorithms

- Linear Regression
- k-Nearest Neighbors
- Naive Bayes
- Decision Trees
- Support Vector Machines

35 / 111

## Linear Regression

36 / 111

"Linear Regression: In statistics, linear regression is an approach for modeling the relationship between a scalar dependent variable  $y$  and one or more explanatory variables denoted  $X$ ."

*Source: [https://en.wikipedia.org/wiki/Linear\\_regression](https://en.wikipedia.org/wiki/Linear_regression)*

37 / 111

## Ordinary Least Squares

38 / 111

# Ordinary Least Squares

- Simple Linear Regression

39 / 111

# Ordinary Least Squares

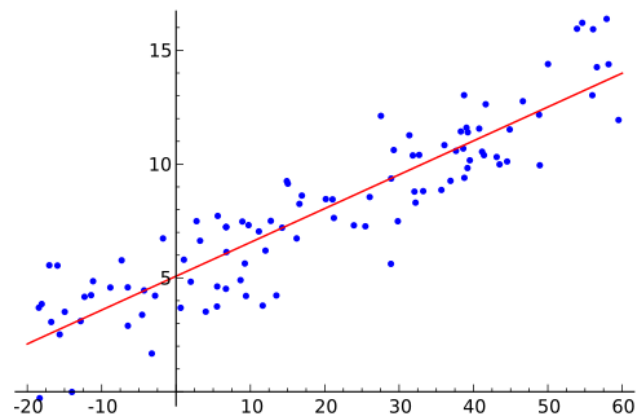
- Simple Linear Regression
- Fit a straight line through the observed points

40 / 111

# Ordinary Least Squares

- Simple Linear Regression
- Fit a straight line through the observed points
- Minimizes the sum of square residuals of the model

41 / 111



[https://commons.wikimedia.org/wiki/File:Linear\\_regression.svg](https://commons.wikimedia.org/wiki/File:Linear_regression.svg)

42 / 111

# Linear Regression

Pros	Cons
Common approach for numeric data	Strong assumptions about the data
Can handle most modeling tasks	Model form must be specified in advance
Estimates the strength and size of the relationships among features and outcomes	Does not handle missing data
	Only numeric features
	Requires statistical knowledge to understand the model

43 / 111

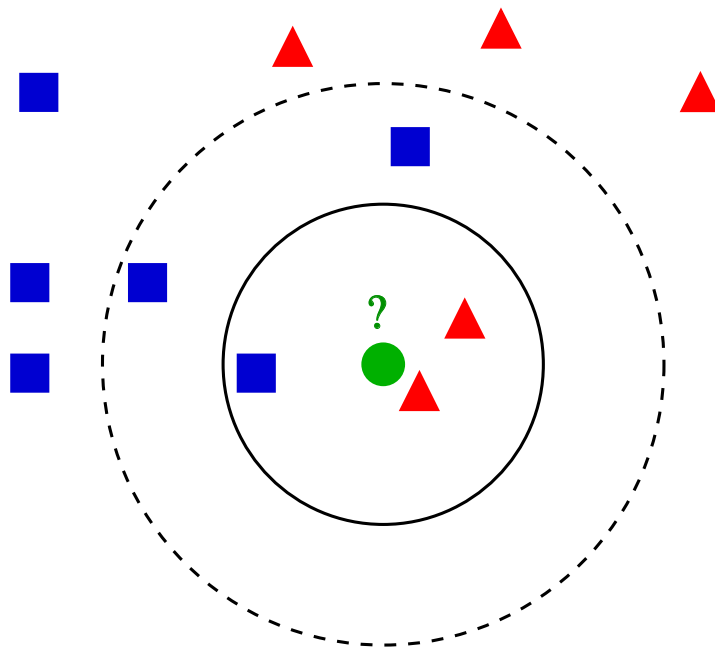
## k-Nearest Neighbors

44 / 111

"k-Nearest Neighbor: a non-parametric method used for classification and regression. In both cases, the input consists of the k closest training examples in the feature space. The output depends on whether k-NN is used for classification or regression."

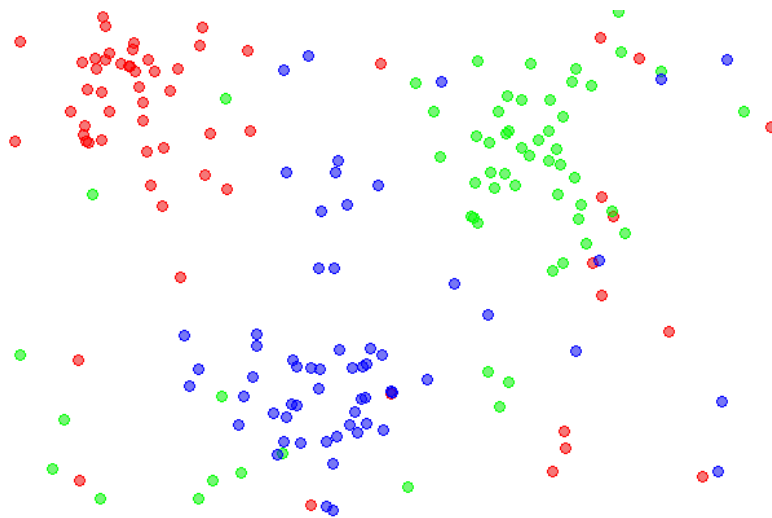
Source: [http://en.wikipedia.org/wiki/K-nearest\\_neighbors\\_algorithm](http://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm)

45 / 111



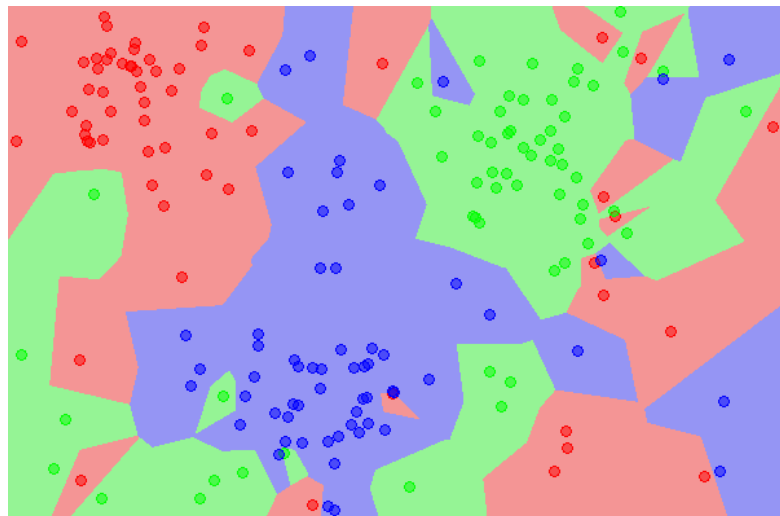
[http://en.wikipedia.org/wiki/K-nearest\\_neighbors\\_algorithm#/media/File:KnnClassification.svg](http://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm#/media/File:KnnClassification.svg)

46 / 111



<http://commons.wikimedia.org/wiki/File:Data3classes.png>

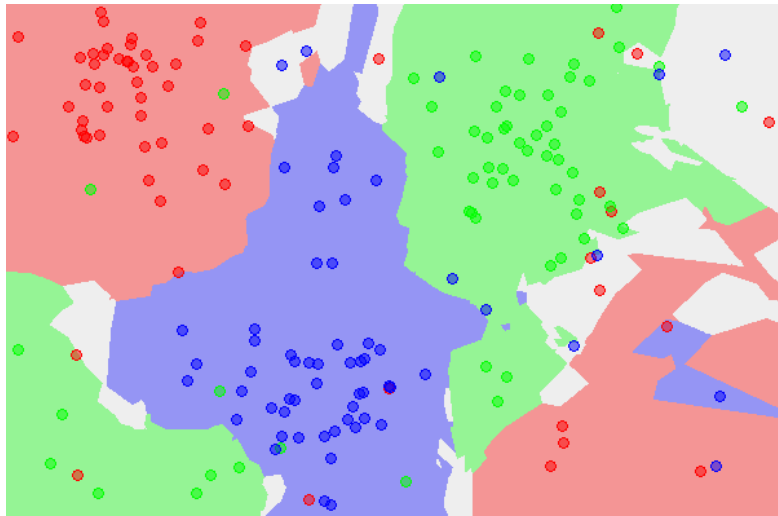
47 / 111



<http://commons.wikimedia.org/wiki/File:Map1NN.png>

48 / 111





<http://commons.wikimedia.org/wiki/File:Map5NN.png>

49 / 111

Pros	Cons
Simple and effective	Does not produce a model
Makes no assumptions about the data distribution	Efficacy affected by choice of $k$
Fast training phase	Slow classification phase

50 / 111

# Naive Bayes

51 / 111

## Naive Bayes

52 / 111

# Naive Bayes

- Family of algorithms to produce probabilistic classifiers based on Bayes Theorem

53 / 111

# Naive Bayes

- Family of algorithms to produce probabilistic classifiers based on Bayes Theorem
- Requires relatively little training data

54 / 111

# Naive Bayes

- Family of algorithms to produce probabilistic classifiers based on Bayes Theorem
- Requires relatively little training data
- Often used for text/document classification

55 / 111

# Naive Bayes

- Family of algorithms to produce probabilistic classifiers based on Bayes Theorem
- Requires relatively little training data
- Often used for text/document classification
- Assumes independence of the features

56 / 111

$$P(A), P(B)$$

$$P(A), P(B)$$

$$P(A \cap B) = P(A) \cdot P(B)$$

59 / 111

$$P(A), P(B)$$

$$P(A \cap B) = P(A) \cdot P(B)$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

60 / 111

$$P(A), P(B)$$

$$P(A \cap B) = P(A) \cdot P(B)$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$

61 / 111

$$P(A), P(B)$$

$$P(A \cap B) = P(A) \cdot P(B)$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$

$$P(spam|Viagra) = \frac{P(Viagra|spam) \cdot P(spam)}{P(Viagra)}$$

62 / 111

Pros	Cons
Simple and effective	Assumption of the independence of features is usually wrong
Does well with noisy and missing data	Doesn't work well with lots of numeric features
Works well with arbitrary sizes of training data	Estimated probabilities aren't as reliable as the classifications
Easy to produce the estimated probability for predictions	

63 / 111

# Decision Trees

64 / 111



# Decision Trees

65 / 111

# Decision Trees

- Tree structure-based classifier

66 / 111

# Decision Trees

- Tree structure-based classifier
- Models relationships between features and outputs

67 / 111

# Decision Trees

- Tree structure-based classifier
- Models relationships between features and outputs
- Easy to explain to users

68 / 111

# Decision Trees

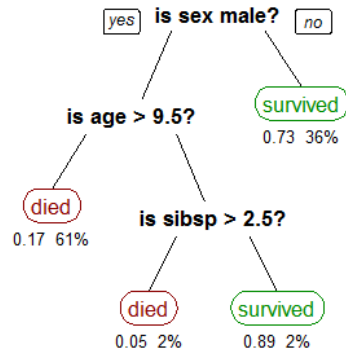
- Tree structure-based classifier
- Models relationships between features and outputs
- Easy to explain to users
- Can be turned into external representation beyond the classifier

69 / 111

# Decision Trees

- Tree structure-based classifier
- Models relationships between features and outputs
- Easy to explain to users
- Can be turned into external representation beyond the classifier
- Supports both classification and regression

70 / 111



[https://en.wikipedia.org/wiki/Decision\\_tree\\_learning](https://en.wikipedia.org/wiki/Decision_tree_learning)

71 / 111

Pros	Cons
Useful classifier for most problems	Can be biased toward feature splits with several levels
Automated learning process	Easy to misfit the model
Supports numeric, nominal and missing data	Small changes in the training data can have been impact on decision logic
Works with large and small data sets	Large trees may be hard to interpret
Easily interpreted	
Efficient	

72 / 111

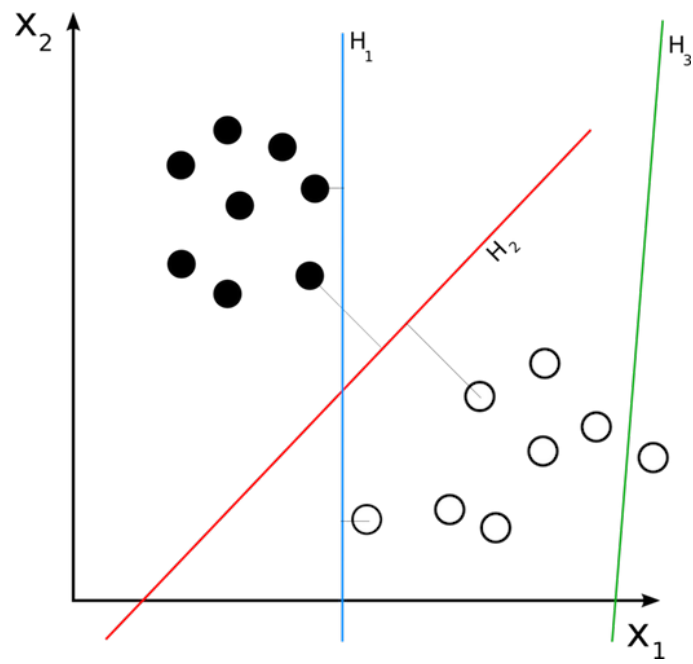
# Support Vector Machines

73 / 111

"[a] support vector machine constructs a hyperplane or set of hyperplanes in a high- or infinite-dimensional space, which can be used for classification, regression, or other tasks"

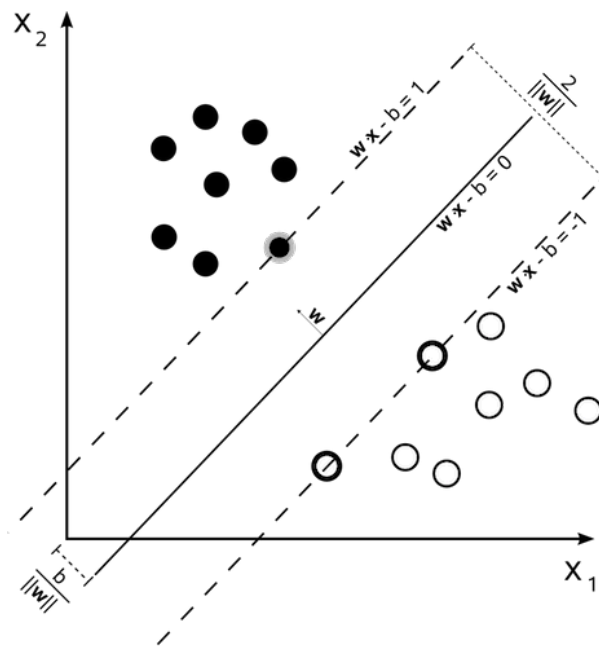
*Source:* [https://en.wikipedia.org/wiki/Support\\_vector\\_machine](https://en.wikipedia.org/wiki/Support_vector_machine)

74 / 111



[https://en.wikipedia.org/wiki/Support\\_vector\\_machine](https://en.wikipedia.org/wiki/Support_vector_machine)

75 / 111



[https://en.wikipedia.org/wiki/Support\\_vector\\_machine](https://en.wikipedia.org/wiki/Support_vector_machine)

76 / 111

# Neural Networks

77 / 111

## Modeling Neurons

78 / 111

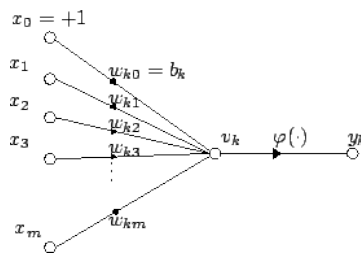
# Modeling Neurons

- Walter Pitts and Warren McCulloch in 1943

79 / 111

# Modeling Neurons

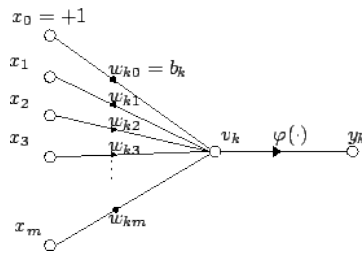
- Walter Pitts and Warren McCulloch in 1943





# Modeling Neurons

- Walter Pitts and Warren McCulloch in 1943



$$y_k = \phi \left( \sum_{j=0}^m w_{kj} x_j \right)$$

CC BY-SA 2.0, <https://commons.wikimedia.org/w/index.php?curid=125222>

81 / 111

# Perceptrons

82 / 111

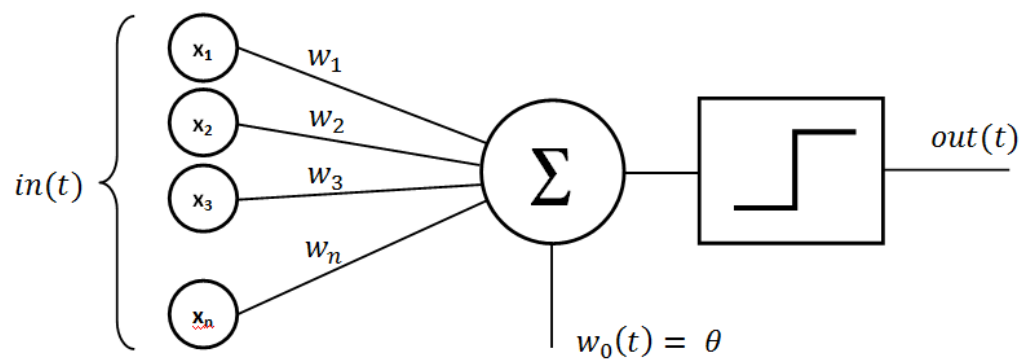
# Perceptrons

- Frank Rosenblatt in late 1950s

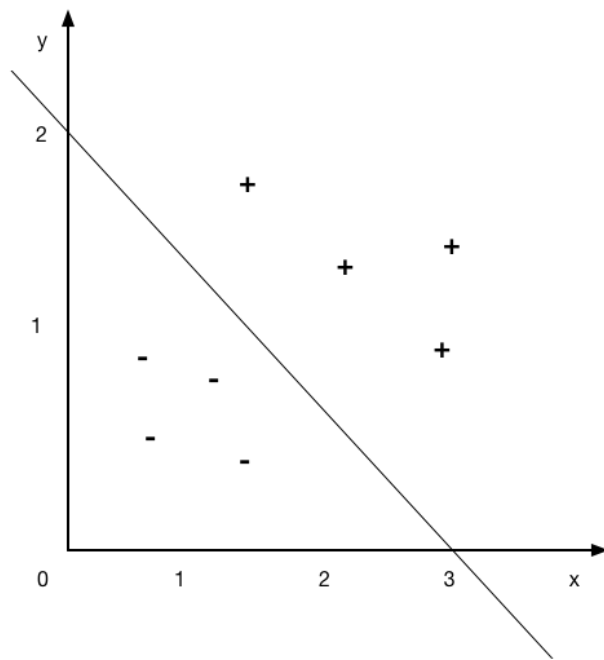
83 / 111

# Perceptrons

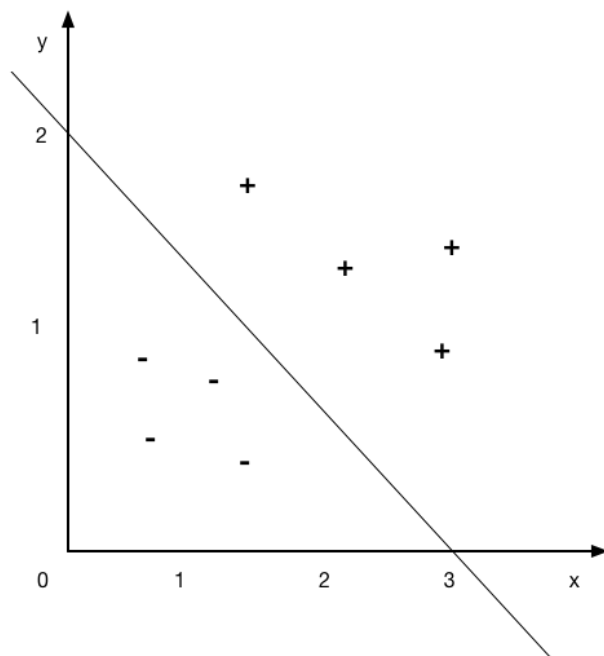
- Frank Rosenblatt in late 1950s



84 / 111

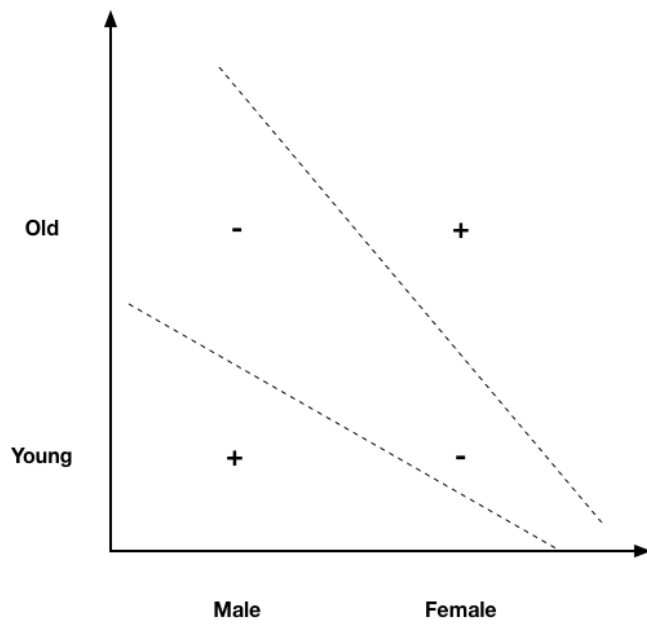


85 / 111

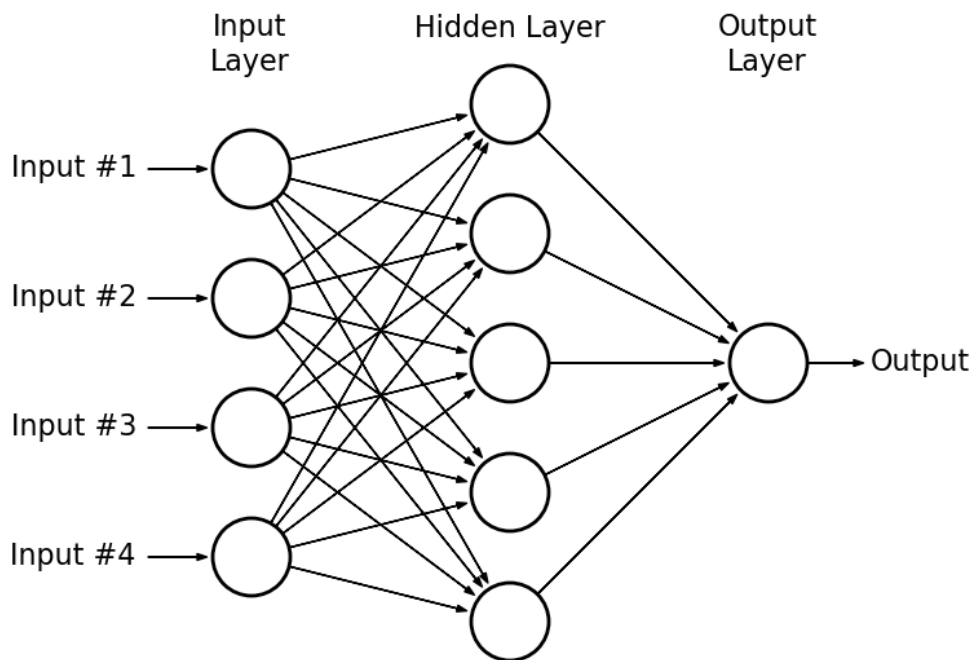


$$2x + 3y = 6$$

86 / 111



87 / 111



88 / 111

<http://neuralnetworksanddeeplearning.com/chap1.html>

89 / 111

<http://colah.github.io/posts/2015-08-Backprop/>

<http://cs231n.github.io/optimization-2/>

90 / 111

<https://github.com/dennybritz/nn-from-scratch>

91 / 111

```
# Create and activate new virtual environment (optional)  
virtualenv venv  
source venv/bin/activate  
# Install requirements  
pip install -r requirements.txt  
# Start the notebook server  
jupyter notebook
```

92 / 111

```
# Create and activate new virtual environment (optional)
virtualenv venv
source venv/bin/activate
# Install requirements
pip install -r requirements.txt
# Start the notebook server
jupyter notebook
```

```
# Package imports
import matplotlib.pyplot as plt
import numpy as np
import sklearn
import sklearn.datasets
import sklearn.linear_model
import matplotlib
```

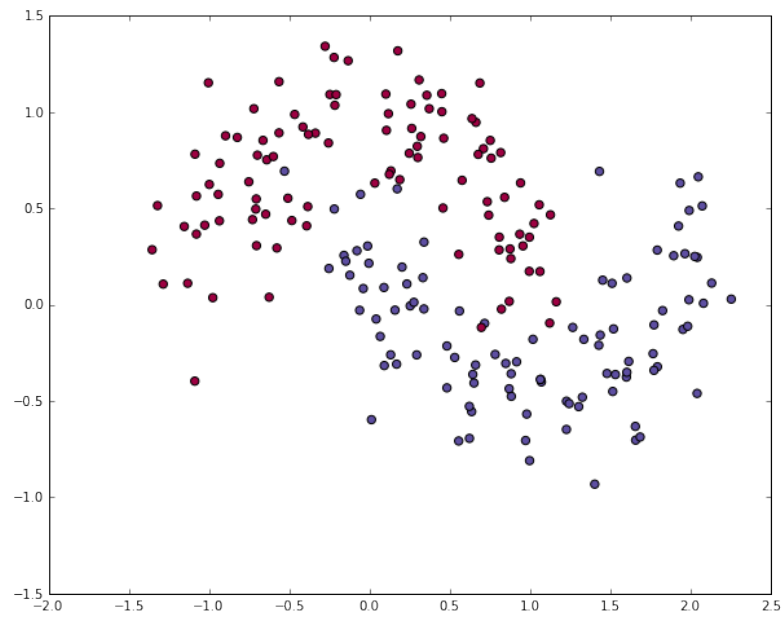
93 / 111

```
# Create and activate new virtual environment (optional)
virtualenv venv
source venv/bin/activate
# Install requirements
pip install -r requirements.txt
# Start the notebook server
jupyter notebook
```

```
# Package imports
import matplotlib.pyplot as plt
import numpy as np
import sklearn
import sklearn.datasets
import sklearn.linear_model
import matplotlib
```

```
# Generate a dataset and plot it
np.random.seed(0)
X, y = sklearn.datasets.make_moons(200, noise=0.20)
plt.scatter(X[:,0], X[:,1], s=40, c=y, cmap=plt.cm.Spectral)
```

94 / 111



95 / 111

```
# Train the logistic regression classifier  
clf = sklearn.linear_model.LogisticRegressionCV()  
clf.fit(X, y)
```

96 / 111



```
# Train the logistic regression classifier
clf = sklearn.linear_model.LogisticRegressionCV()
clf.fit(X, y)
```

```
# Helper function to plot a decision boundary.
# If you don't fully understand this function don't worry, it just generates the c
def plot_decision_boundary(pred_func):
    # Set min and max values and give it some padding
    x_min, x_max = X[:, 0].min() - .5, X[:, 0].max() + .5
    y_min, y_max = X[:, 1].min() - .5, X[:, 1].max() + .5
    h = 0.01
    # Generate a grid of points with distance h between them
    xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
    # Predict the function value for the whole grid
    Z = pred_func(np.c_[xx.ravel(), yy.ravel()])
    Z = Z.reshape(xx.shape)
    # Plot the contour and training examples
    plt.contourf(xx, yy, Z, cmap=plt.cm.Spectral)
    plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Spectral)
```

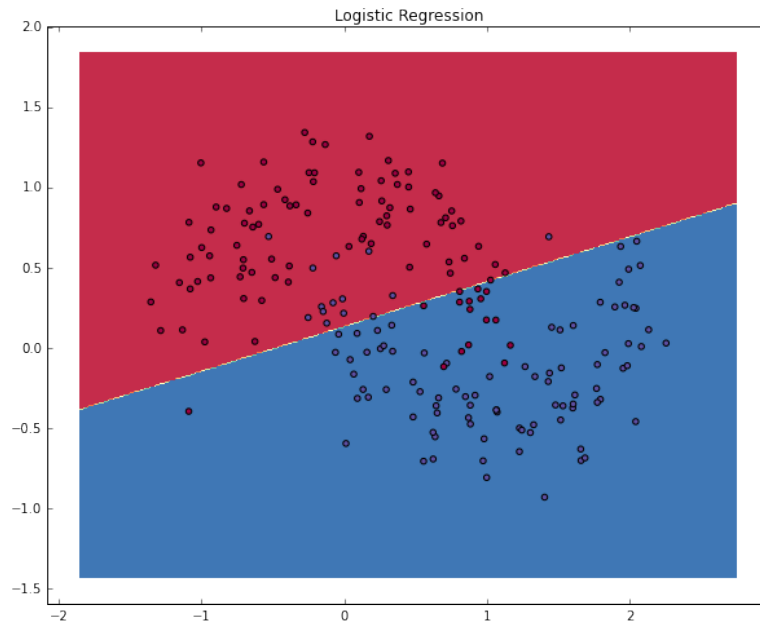
97 / 111

```
# Train the logistic regression classifier
clf = sklearn.linear_model.LogisticRegressionCV()
clf.fit(X, y)
```

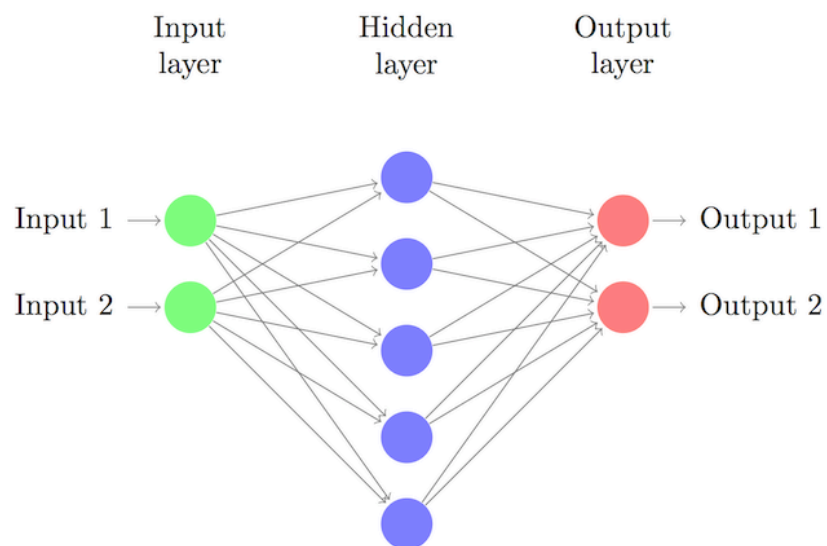
```
# Helper function to plot a decision boundary.
# If you don't fully understand this function don't worry, it just generates the c
def plot_decision_boundary(pred_func):
    # Set min and max values and give it some padding
    x_min, x_max = X[:, 0].min() - .5, X[:, 0].max() + .5
    y_min, y_max = X[:, 1].min() - .5, X[:, 1].max() + .5
    h = 0.01
    # Generate a grid of points with distance h between them
    xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
    # Predict the function value for the whole grid
    Z = pred_func(np.c_[xx.ravel(), yy.ravel()])
    Z = Z.reshape(xx.shape)
    # Plot the contour and training examples
    plt.contourf(xx, yy, Z, cmap=plt.cm.Spectral)
    plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Spectral)
```

```
# Plot the decision boundary
plot_decision_boundary(lambda x: clf.predict(x))
plt.title("Logistic Regression")
```

98 / 111



99 / 111



100 / 111

```

num_examples = len(X) # training set size
nn_input_dim = 2 # input layer dimensionality
nn_output_dim = 2 # output layer dimensionality

# Gradient descent parameters (I picked these by hand)
epsilon = 0.01 # learning rate for gradient descent
reg_lambda = 0.01 # regularization strength

```

101 / 111

```

num_examples = len(X) # training set size
nn_input_dim = 2 # input layer dimensionality
nn_output_dim = 2 # output layer dimensionality

# Gradient descent parameters (I picked these by hand)
epsilon = 0.01 # learning rate for gradient descent
reg_lambda = 0.01 # regularization strength

```

```

# Helper function to evaluate the total loss on the dataset
def calculate_loss(model):
    W1, b1, W2, b2 = model['W1'], model['b1'], model['W2'], model['b2']
    # Forward propagation to calculate our predictions
    z1 = X.dot(W1) + b1
    a1 = np.tanh(z1)
    z2 = a1.dot(W2) + b2
    exp_scores = np.exp(z2)
    probs = exp_scores / np.sum(exp_scores, axis=1, keepdims=True)
    # Calculating the loss
    correct_logprobs = -np.log(probs[range(num_examples), y])
    data_loss = np.sum(correct_logprobs)
    # Add regularization term to loss (optional)
    data_loss += reg_lambda/2 * (np.sum(np.square(W1)) + np.sum(np.square(W2)))
    return 1./num_examples * data_loss

```

102 / 111

```

# Helper function to predict an output (0 or 1)
def predict(model, x):
    W1, b1, W2, b2 = model['W1'], model['b1'], model['W2'], model['b2']
    # Forward propagation
    z1 = x.dot(W1) + b1
    a1 = np.tanh(z1)
    z2 = a1.dot(W2) + b2
    exp_scores = np.exp(z2)
    probs = exp_scores / np.sum(exp_scores, axis=1, keepdims=True)
    return np.argmax(probs, axis=1)

```

103 / 111

```

# This function learns parameters for the neural network and returns the model.
# - nn_hdim: Number of nodes in the hidden layer
# - num_passes: Number of passes through the training data for gradient descent
# - print_loss: If True, print the loss every 1000 iterations
def build_model(nn_hdim, num_passes=20000, print_loss=False):

    # Initialize the parameters to random values. We need to learn these.
    np.random.seed(0)
    W1 = np.random.randn(nn_input_dim, nn_hdim) / np.sqrt(nn_input_dim)
    b1 = np.zeros((1, nn_hdim))
    W2 = np.random.randn(nn_hdim, nn_output_dim) / np.sqrt(nn_hdim)
    b2 = np.zeros((1, nn_output_dim))

    # This is what we return at the end
    model = {}

    # Gradient descent. For each batch...
    for i in xrange(0, num_passes):

        # Forward propagation
        z1 = X.dot(W1) + b1
        a1 = np.tanh(z1)
        z2 = a1.dot(W2) + b2
        exp_scores = np.exp(z2)
        probs = exp_scores / np.sum(exp_scores, axis=1, keepdims=True)

        ...

```

104 / 111

```

...

# Backpropagation
delta3 = probs
delta3[range(num_examples), y] -= 1
dW2 = (a1.T).dot(delta3)
db2 = np.sum(delta3, axis=0, keepdims=True)
delta2 = delta3.dot(W2.T) * (1 - np.power(a1, 2))
dW1 = np.dot(X.T, delta2)
db1 = np.sum(delta2, axis=0)

# Add regularization terms (b1 and b2 don't have regularization terms)
dW2 += reg_lambda * W2
dW1 += reg_lambda * W1

# Gradient descent parameter update
W1 += -epsilon * dW1
b1 += -epsilon * db1
W2 += -epsilon * dW2
b2 += -epsilon * db2

# Assign new parameters to the model
model = { 'W1': W1, 'b1': b1, 'W2': W2, 'b2': b2}

# Optionally print the loss.
# This is expensive because it uses the whole dataset, so we don't want to
if print_loss and i % 1000 == 0:
    print "Loss after iteration %i: %f" %(i, calculate_loss(model))

return model

```

105 / 111

```

# Build a model with a 3-dimensional hidden layer
model = build_model(3, print_loss=True)

```

```

Loss after iteration 0: 0.432387
Loss after iteration 1000: 0.068947
Loss after iteration 2000: 0.068936
Loss after iteration 3000: 0.071218
Loss after iteration 4000: 0.071253
Loss after iteration 5000: 0.071278
Loss after iteration 6000: 0.071293
Loss after iteration 7000: 0.071303
Loss after iteration 8000: 0.071308
Loss after iteration 9000: 0.071312
Loss after iteration 10000: 0.071314
Loss after iteration 11000: 0.071315
Loss after iteration 12000: 0.071315
Loss after iteration 13000: 0.071316
Loss after iteration 14000: 0.071316
Loss after iteration 15000: 0.071316
Loss after iteration 16000: 0.071316
Loss after iteration 17000: 0.071316
Loss after iteration 18000: 0.071316
Loss after iteration 19000: 0.071316

```

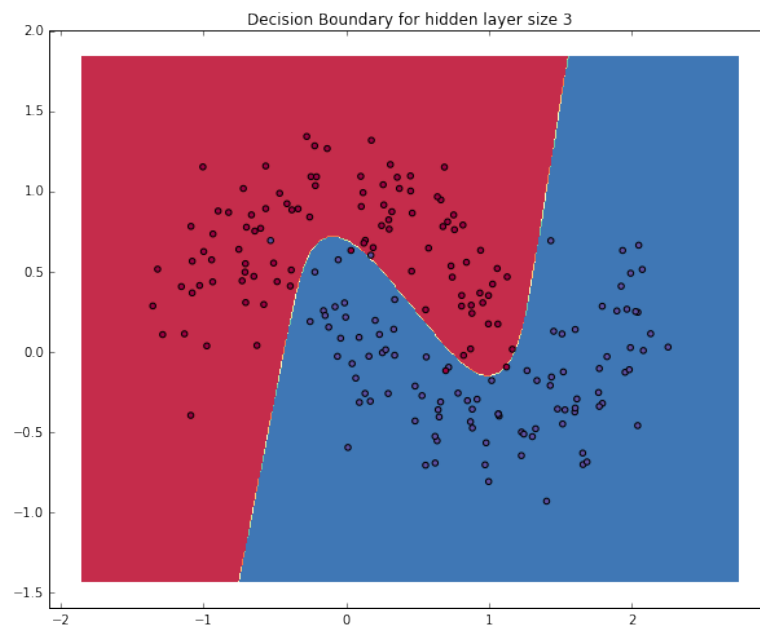
106 / 111

```
# Build a model with a 3-dimensional hidden layer
model = build_model(3, print_loss=True)
```

```
Loss after iteration 0: 0.432387
Loss after iteration 1000: 0.068947
Loss after iteration 2000: 0.068936
Loss after iteration 3000: 0.071218
Loss after iteration 4000: 0.071253
Loss after iteration 5000: 0.071278
Loss after iteration 6000: 0.071293
Loss after iteration 7000: 0.071303
Loss after iteration 8000: 0.071308
Loss after iteration 9000: 0.071312
Loss after iteration 10000: 0.071314
Loss after iteration 11000: 0.071315
Loss after iteration 12000: 0.071315
Loss after iteration 13000: 0.071316
Loss after iteration 14000: 0.071316
Loss after iteration 15000: 0.071316
Loss after iteration 16000: 0.071316
Loss after iteration 17000: 0.071316
Loss after iteration 18000: 0.071316
Loss after iteration 19000: 0.071316
```

```
# Plot the decision boundary
plot_decision_boundary(lambda x: predict(model, x))
plt.title("Decision Boundary for hidden layer size 3")
```

107 / 111



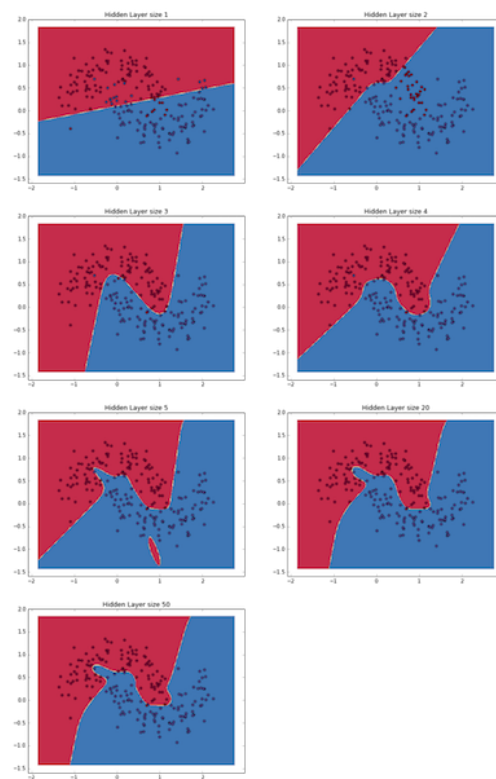
108 / 111

```

plt.figure(figsize=(16, 32))
hidden_layer_dimensions = [1, 2, 3, 4, 5, 20, 50]
for i, nn_hdim in enumerate(hidden_layer_dimensions):
    plt.subplot(5, 2, i+1)
    plt.title('Hidden Layer size %d' % nn_hdim)
    model = build_model(nn_hdim)
    plot_decision_boundary(lambda x: predict(model, x))
plt.show()

```

109 / 111



110 / 111

A large, golden Buddha statue is shown from the chest up, positioned on the left side of the frame. The statue has a serene expression and is set against a bright blue sky filled with fluffy white clouds. The lighting suggests a sunny day.

# Questions

✉ [brian@bosatsu.net](mailto:brian@bosatsu.net)

🐦 [@bsletten](https://twitter.com/bsletten)

🔗 [bsletten](https://bsletten.com)