

# Machine Learning & MLLib



# Agenda



- Introduction to Machine Learning
  - Machine Learning
  - Algorithms
- Machine Learning Using Spark MLLib

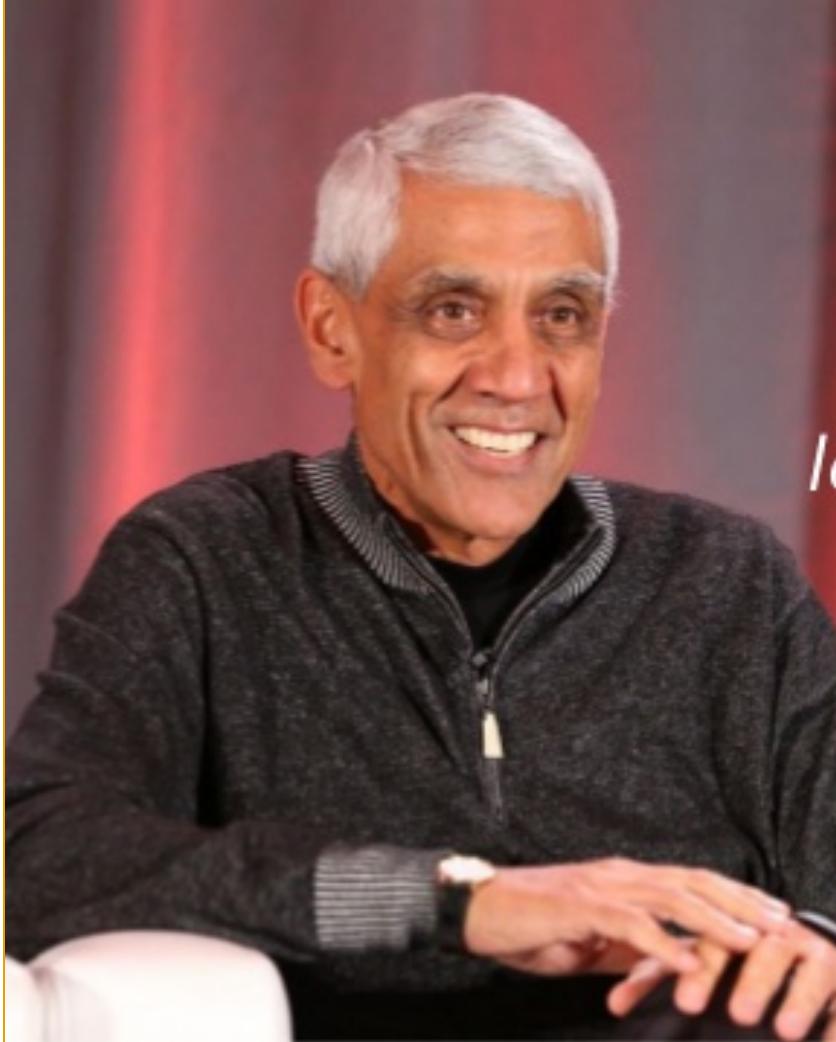




*“A breakthrough in machine learning would be worth ten Microsofts”*

–Bill Gates

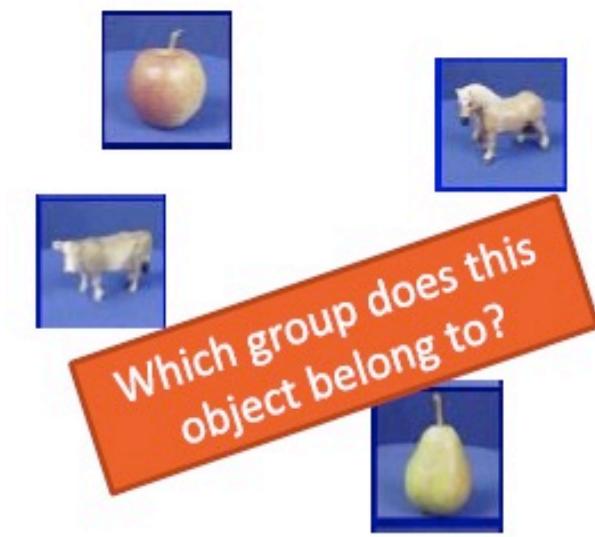


A portrait of Vinod Khosla, an elderly man with white hair, smiling and wearing a dark grey zip-up sweater. He is seated in front of a background with vertical red and grey stripes.

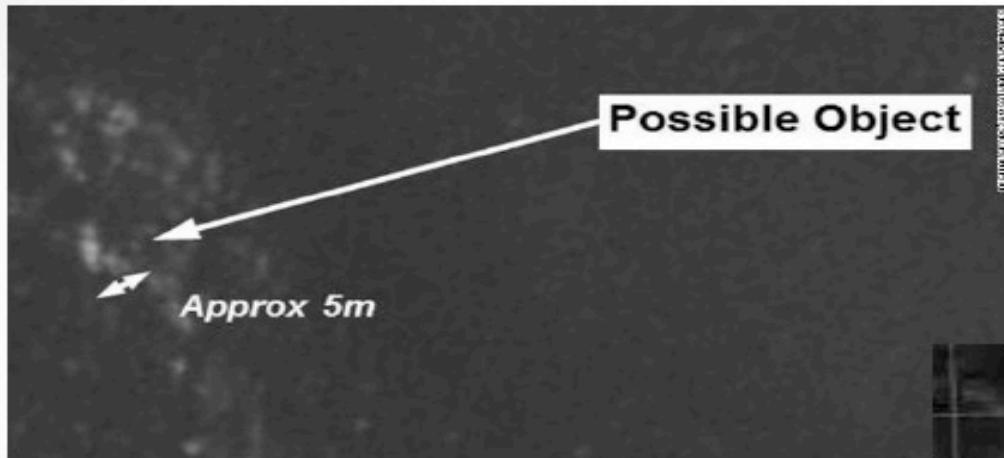
*"In the next 20 years, machine learning will have more impact than mobile has."*

—Vinod Khosla

# Simple Example



# How about this?





16 Mar 2014  
04162/1216.local  
© DigitalGlobe  
Panchromatic

UNCLASSIFIED

## Object 2 Possibly Associated With MH370 Search

Indian Ocean COORDS: 44°03'025" S 091°13'27"E



16 Mar 2014  
04162/1216.local  
© DigitalGlobe  
Multispectral

Possible Object

Approx 5m

Possible Object

Approx 5m

N

DIGO-00718-02-14

# Questions to Ponder



- How did your brain process the images, and group them?
- We call that part as learning
  - Learning could be visual, hearing, sense, touch etc based
  - After learning, we look at new images, and compare with groups classified during earlier phase
    - This is prediction/forecasting
- This is how humans have evolved...we are natural at this
- Machines..not so much.
- ML is that entire branch that attempts to make it possible

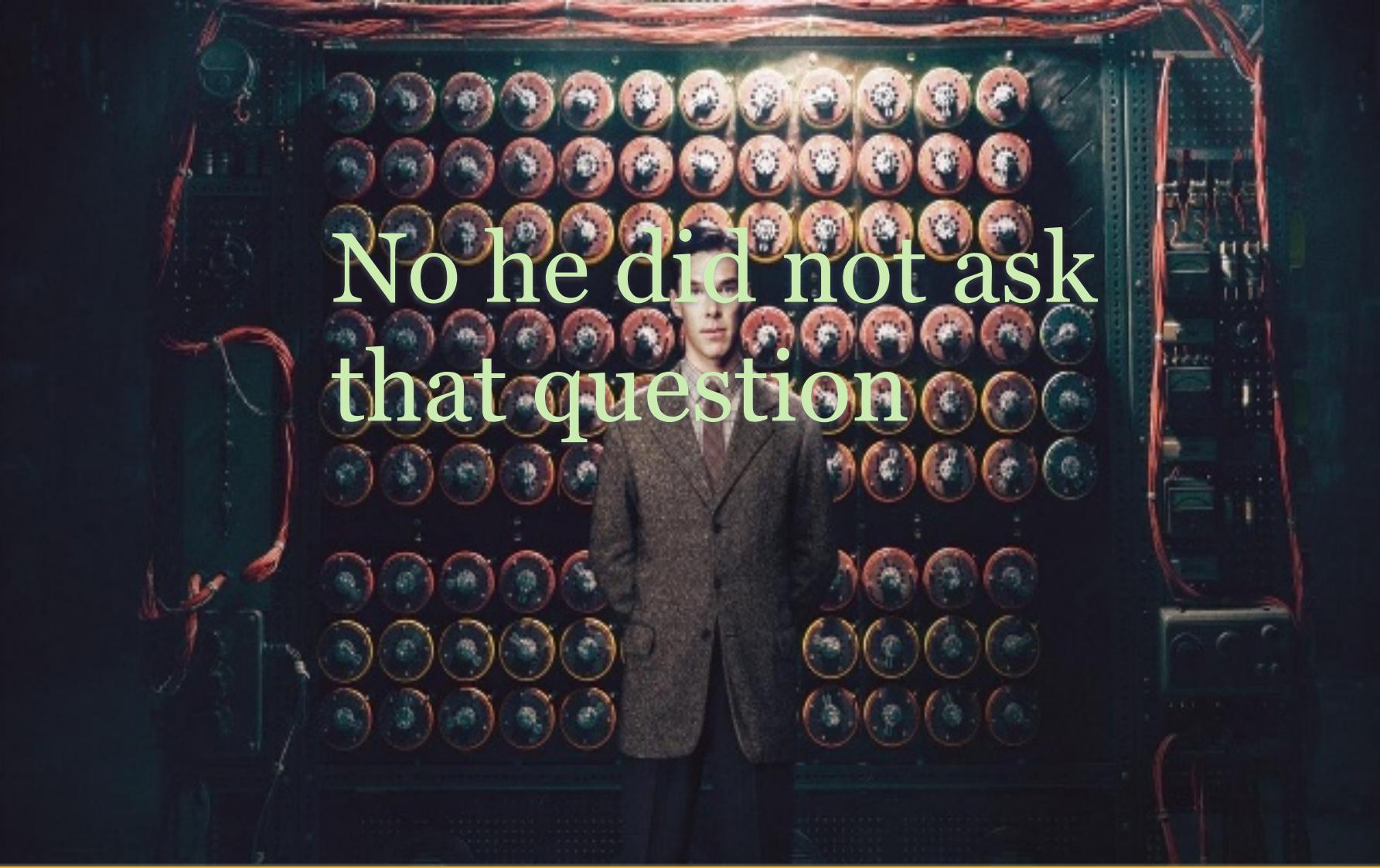
# Idea behind ML



- So how old is ML?
  - Literally as old as computing itself
- Idea originated with a simple question:
- CAN MACHINES THINK?

Remember this guy?

No he did not ask  
that question

A photograph of a man in a brown tweed suit standing in front of a large, dark wooden control panel. The panel is covered in numerous circular control knobs, each with a small glowing blue light in the center. The man has short brown hair and is looking directly at the camera with a neutral expression. He is wearing a white shirt and a dark tie. To his left, a portion of a red leather chair is visible. The background is dark, making the glowing knobs stand out.

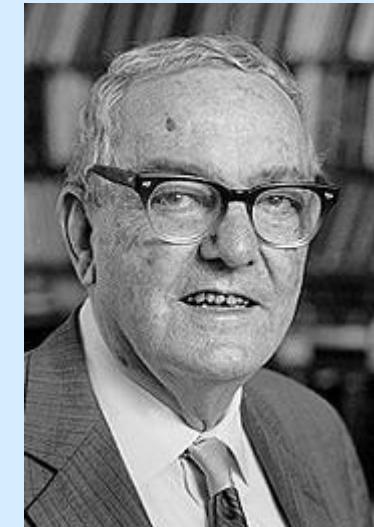
He did....in 1950!!! (Alan Turing)



# Machine Learning



- **Herbert Alexander Simon:**  
“Learning is any process by which a system improves performance from experience.”
- “Machine Learning is concerned with computer programs that automatically improve their performance through experience.  
“



**Herbert Simon**  
Turing Award 1975  
Nobel Prize in Economics  
1978

# Why Machine Learning?



- Develop systems that can automatically adapt and customize themselves to individual users.
  - Personalized news or mail filter
- Discover new knowledge from large databases (**data mining**).
  - Market basket analysis (e.g. diapers and beer)
- Ability to mimic human behavior and replace certain monotonous tasks - which require some intelligence.
  - ★ like recognizing handwritten characters

# Machine Learning



- Systems learn with data
  - Different from traditional paradigm
  - Logic learnt from data
- Computation intensive
- Ubiquitous
  - Apple (Siri), Google (Search), Amazon (Recommendations), FB
- Eventual goal is to predict or classify

# Why now?



- Flood of available data (especially with the advent of the Internet)
- Increasing computational power
- Growing progress in available algorithms and theory developed by researchers
- Increasing support from industries

# ML Applications

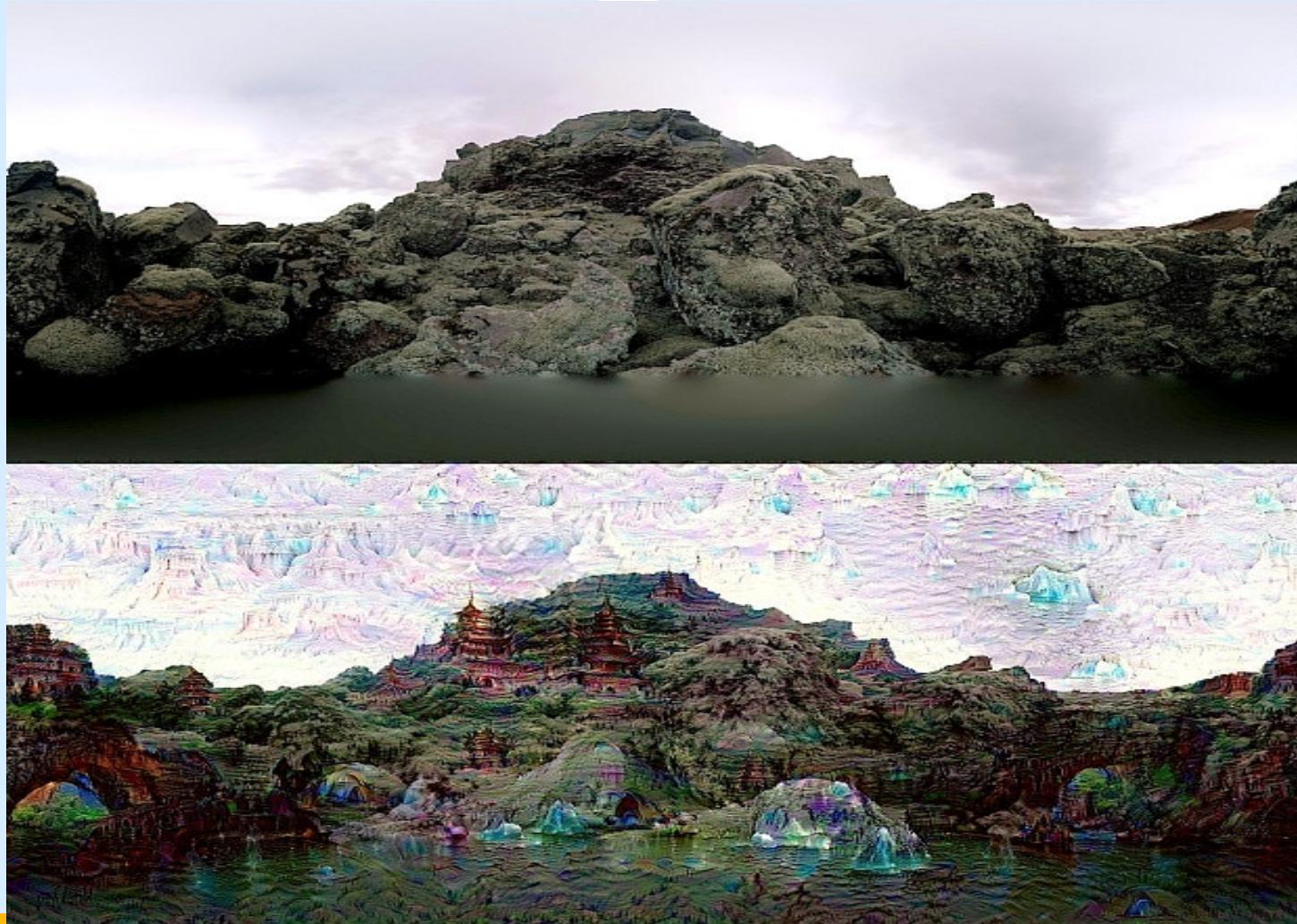


# So what are the buzzwords these days?



- Machine Learning?
  - Nice
- Neural Networks?
  - Good
- Deep Learning?
  - Great
- Machine Dreaming??????????
  - 😞

# Detour: Machine Dreaming



# The concept of learning in a ML system



- Learning = Improving with experience at some task
  - Improve over task  $T$ ,
  - With respect to performance measure,  $P$
  - Based on experience,  $E$ .

# Example: Spam Filter

## **Example:** Spam Filtering

Spam - is all email the user does not want to receive and has not asked to receive

*T:* Identify Spam Emails

*P:*

% of spam emails that were filtered

% of ham/ (non-spam) emails that  
were incorrectly filtered-out

*E:* a database of emails that were labelled  
by users



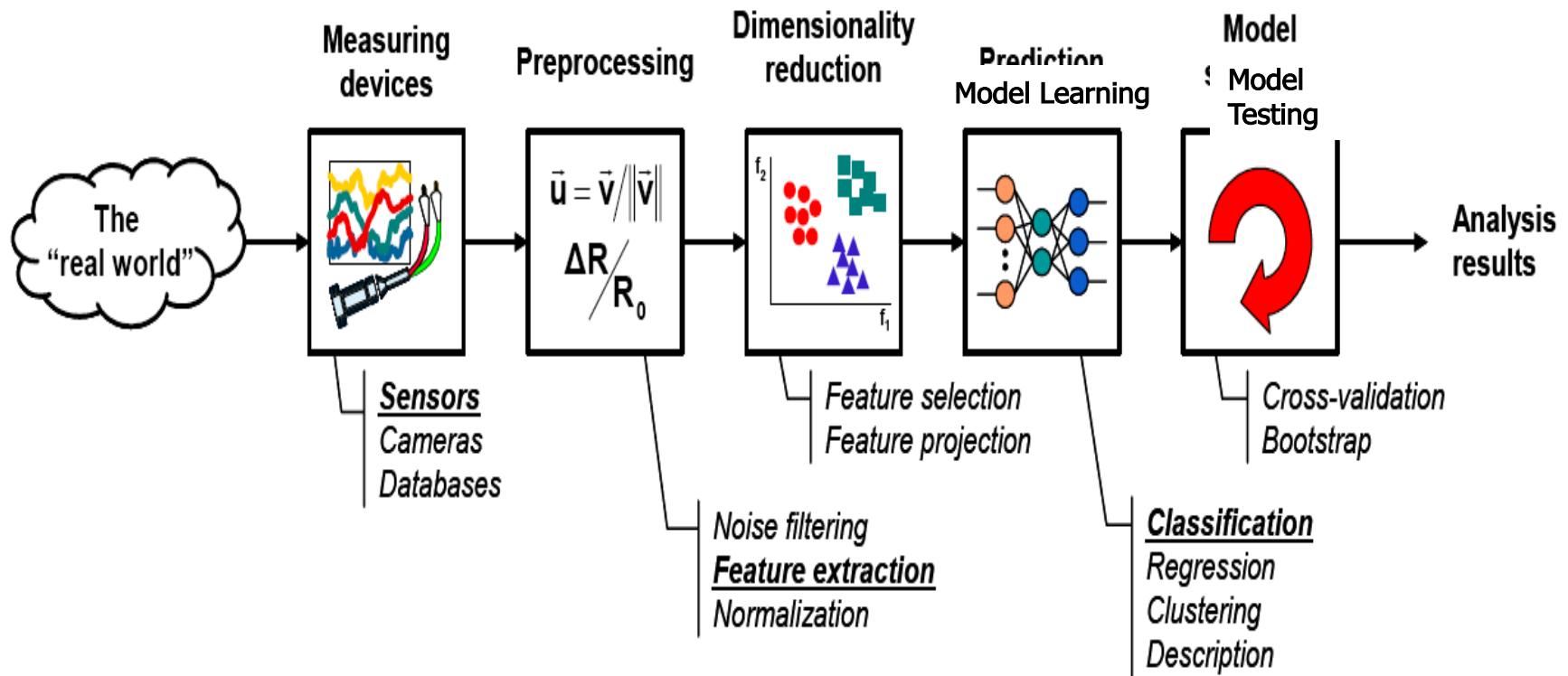
# Language of ML



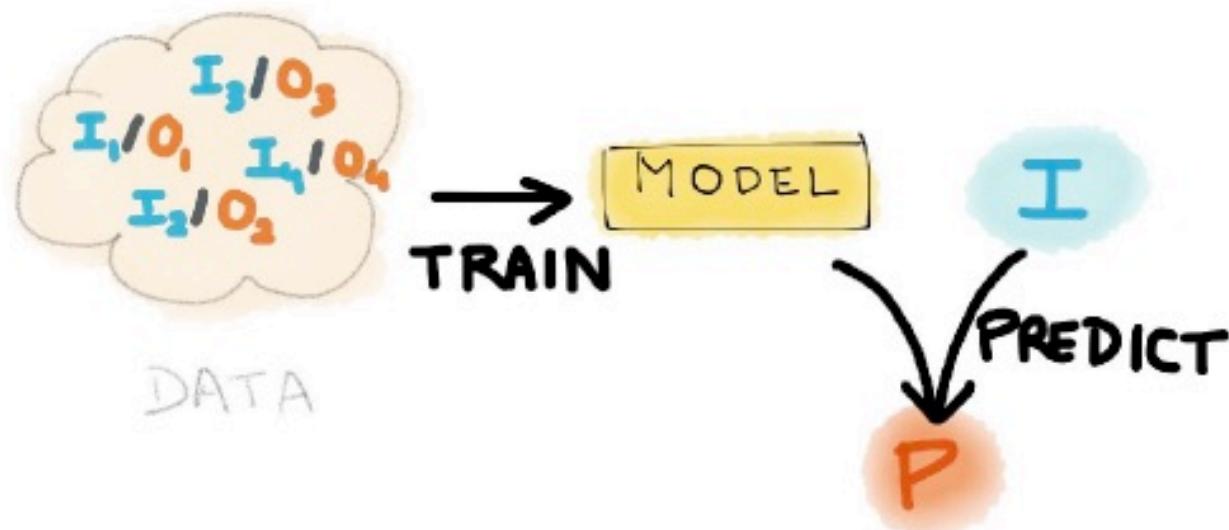
## Preprocessing



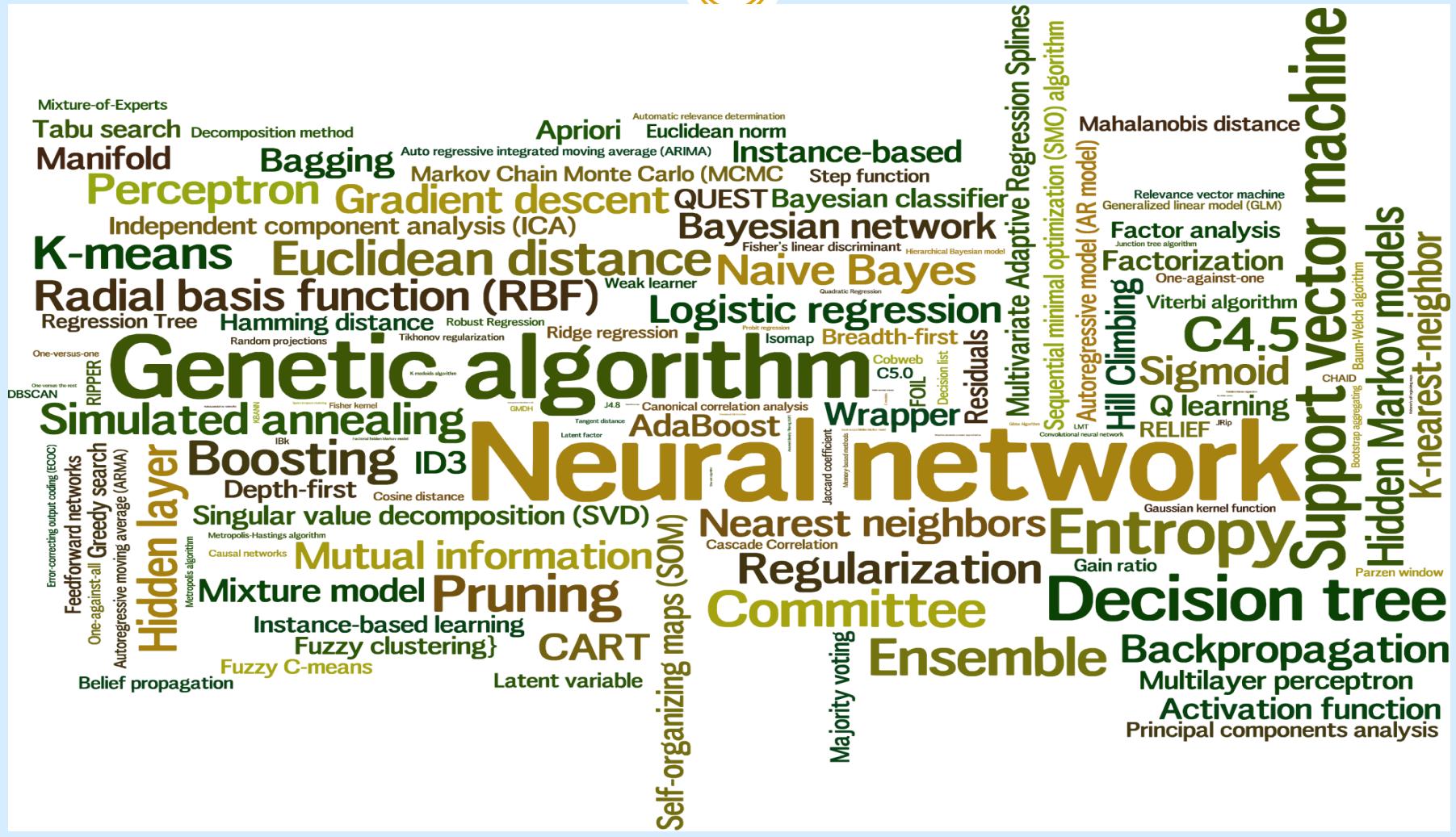
# The Learning Process



# Simplify....



# World of ML Algos



# Connect between ML and Statistics/Probability



- ML algorithms are based on a bunch of assumptions
- Eg. Data distribution: Normal, linear/non-linear
- They are a simplified representation of reality
- ML algos tell us how likely a given result is?
  - This means results come to us not as Yes/No but probability of an observation being Yes/No
  - This probability output is result of assumptions and actual algorithm
- No ML algo can exist without probability and statistical concepts

# Terminology



- **Feature**
  - Attribute or a variable
  - Table – row is observation, and column a feature
  - Dimensionality
  - Categorical or numerical
- **Label**
  - Dependent variable
  - To be predicted
  - Numerical or categorical
- **Model**
  - Mathematical representation of relationship betn variables

# Machine Learning



- Machine learning algo fits model to data
- Computationally intensive
- Train and test data
- Applications
  - Classification
  - Regression
  - Clustering
  - Recommendation

# (Broad) Categories of ML



- **Supervised (Classification or Prediction)**
  - First, train the algorithm to create a model.
  - What is training?
  - Data comes along with output as well.
  - New data is labelled based on model created from training data.
- **Unsupervised (Clustering)**
  - We do not know output, just have a bunch of observations
  - What we try to do is keep them in buckets according to certain similarities
- **Reinforcement based (think ANN, Deep Learning)**
  - Learning based on feedback or reward

# What if only this was given?



Age	Spectacle prescription	Astigmatism	Tear production rate
Young	Myope	No	Reduced
Young	Myope	No	Normal
Young	Myope	Yes	Reduced
Young	Myope	Yes	Normal
Young	Hypermetrope	No	Reduced
Young	Hypermetrope	No	Normal
Young	Hypermetrope	Yes	Reduced
Young	Hypermetrope	Yes	Normal
Pre-presbyopic	Myope	No	Reduced
Pre-presbyopic	Myope	No	Normal
Pre-presbyopic	Myope	Yes	Reduced
Pre-presbyopic	Myope	Yes	Normal
Pre-presbyopic	Hypermetrope	No	Reduced
Pre-presbyopic	Hypermetrope	No	Normal
Pre-presbyopic	Hypermetrope	Yes	Reduced
Pre-presbyopic	Hypermetrope	Yes	Normal
Presbyopic	Myope	No	Reduced
Presbyopic	Myope	No	Normal
Presbyopic	Myope	Yes	Reduced
Presbyopic	Myope	Yes	Normal
Presbyopic	Hypermetrope	No	Reduced
Presbyopic	Hypermetrope	No	Normal
Presbyopic	Hypermetrope	Yes	Reduced
Presbyopic	Hypermetrope	Yes	Normal
Presbyopic	Hypermetrope	Yes	Reduced

# Example



Age	Spectacle prescription	Astigmatism	Tear production rate	Recommended lenses
Young	Myope	No	Reduced	None
Young	Myope	No	Normal	Soft
Young	Myope	Yes	Reduced	None
Young	Myope	Yes	Normal	Hard
Young	Hypermetrope	No	Reduced	None
Young	Hypermetrope	No	Normal	Soft
Young	Hypermetrope	Yes	Reduced	None
Young	Hypermetrope	Yes	Normal	hard
Pre-presbyopic	Myope	No	Reduced	None
Pre-presbyopic	Myope	No	Normal	Soft
Pre-presbyopic	Myope	Yes	Reduced	None
Pre-presbyopic	Myope	Yes	Normal	Hard
Pre-presbyopic	Hypermetrope	No	Reduced	None
Pre-presbyopic	Hypermetrope	No	Normal	Soft
Pre-presbyopic	Hypermetrope	Yes	Reduced	None
Pre-presbyopic	Hypermetrope	Yes	Normal	None
Presbyopic	Myope	No	Reduced	None
Presbyopic	Myope	No	Normal	None
Presbyopic	Myope	Yes	Reduced	None
Presbyopic	Myope	Yes	Normal	Hard
Presbyopic	Hypermetrope	No	Reduced	None
Presbyopic	Hypermetrope	No	Normal	Soft
Presbyopic	Hypermetrope	Yes	Reduced	None
Presbyopic	Hypermetrope	Yes	Normal	None

# Supervised



```
If tear production rate = reduced then recommendation = none
If age = young and astigmatic = no
    and tear production rate = normal then recommendation = soft
If age = pre-presbyopic and astigmatic = no
    and tear production rate = normal then recommendation = soft
If age = presbyopic and spectacle prescription = myope
    and astigmatic = no then recommendation = none
If spectacle prescription = hypermetrope and astigmatic = no
    and tear production rate = normal then recommendation = soft
If spectacle prescription = myope and astigmatic = yes
    and tear production rate = normal then recommendation = hard
If age young and astigmatic = yes
    and tear production rate = normal then recommendation = hard
If age = pre-presbyopic
    and spectacle prescription = hypermetrope
    and astigmatic = yes then recommendation = none
If age = presbyopic and spectacle prescription = hypermetrope
    and astigmatic = yes then recommendation = none
```

# Supervised Learning



- Idea is to learn like humans from past experience
- But machines do not have any experience
- Experience here comes from data
- So what do we learn
  - Attempt to solve a target function
  - In the process we would come up with a formula/decision rules

# Supervised and Unsupervised



- Supervised
  - Data comes with a y “output” variable.
  - Algorithm is learning under the supervision
- Unsupervised
  - There is no y variable
  - Task is try and find some sort of clusters or classes in the data

# Supervised Algorithms



- Trains model with labeled dataset
- Create model based on training data and use it to classify new data
- Largely, classification and prediction
- Prediction
  - Predict numerical variable (generally continuous valued)
  - Linear regression, decision tree, ensemble
- Classification
  - Predict categorical variables (discrete or nominal)
  - Logistic regression, Neural network, Naive Bayes
- Recommendations

# Classification

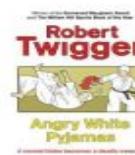


- “Which category of products is most interesting to this customer?”
- “Is this movie a romantic comedy, documentary, or thriller?”
- “Is this review written by a customer or a robot?”
- “Will the customer buy this product?”
- “Is this email spam or not spam?”

# Recommendation



## Books



### Readers



3	2	?	5	?	?
4	?	3	0	?	5
1	4	?	?	3	2
?	3	?	2	?	4

# Recommendation



- “Which movies should be recommended to a user?”
- “If the user just listened to a song, which song would he like now?”
- “Which news articles are relevant for a user in a particular context?”
- “Which advertisements should be displayed for a user on a mobile app?”
- “Which products are frequently bought together?”

# Two step process



- It is a two-step process:
  - First, generate model using training dataset
  - Second, once the model is generated then predict new observations using that model

# Before you run an algo, you do feature engg



- Derive new features from the initial data set:
  1. Aggregations: Count, Sum, Min, Max, Avg, Std
  2. Temporal: Elapsed time, Trends
  3. Continuous to Categorical: Converting real values to enumerations.
  4. Categorical to Binary: Converting enumerations to binary features.
  5. Domain-specific derived features.

# You have a huge menu to choose from



## Machine Learning Algorithms *(sample)*

**Continuous**

### Unsupervised

- Clustering & Dimensionality Reduction
  - SVD
  - PCA
  - K-means

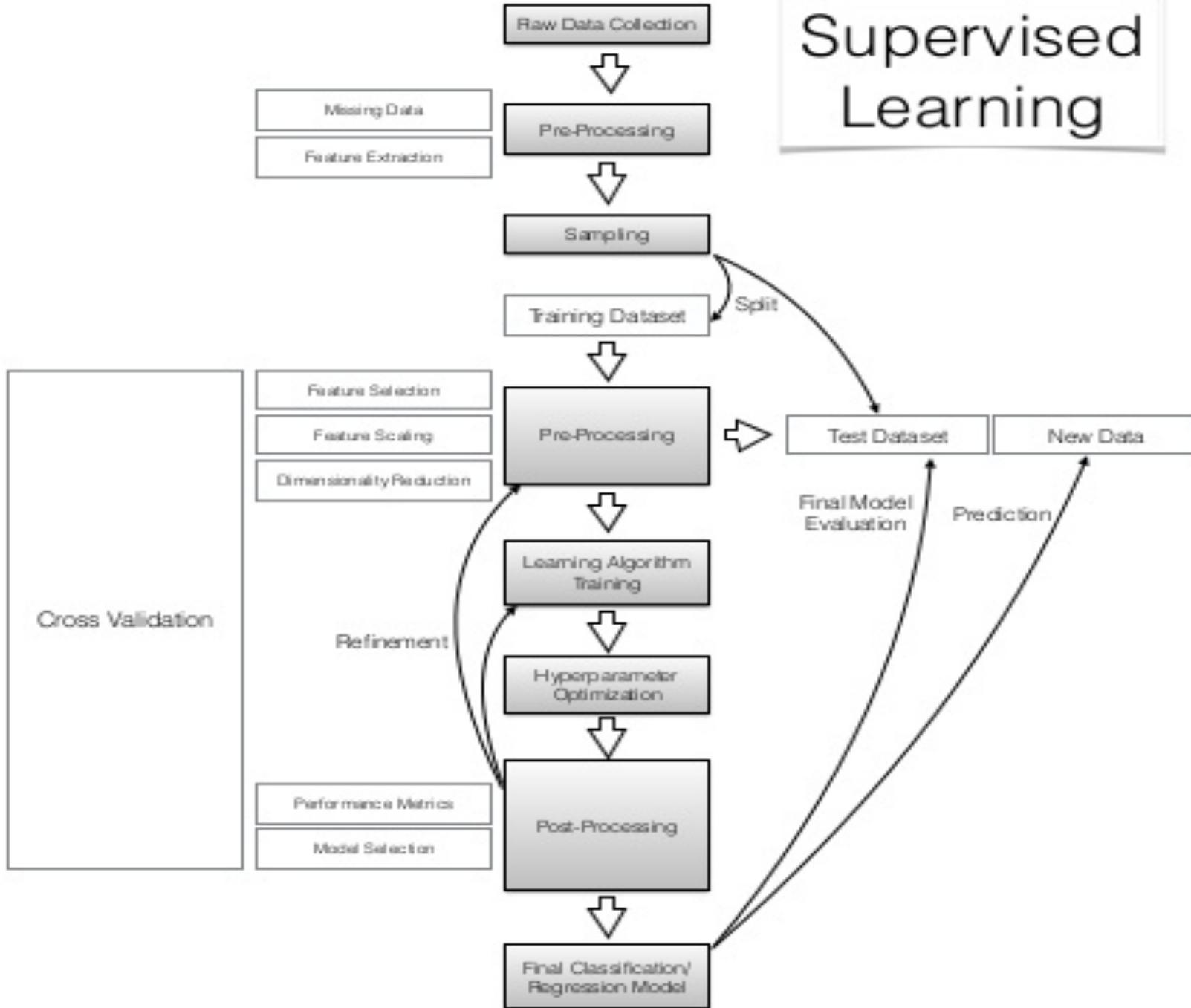
**Categorical**

- Association Analysis
  - Apriori
  - FP-Growth
- Hidden Markov Model

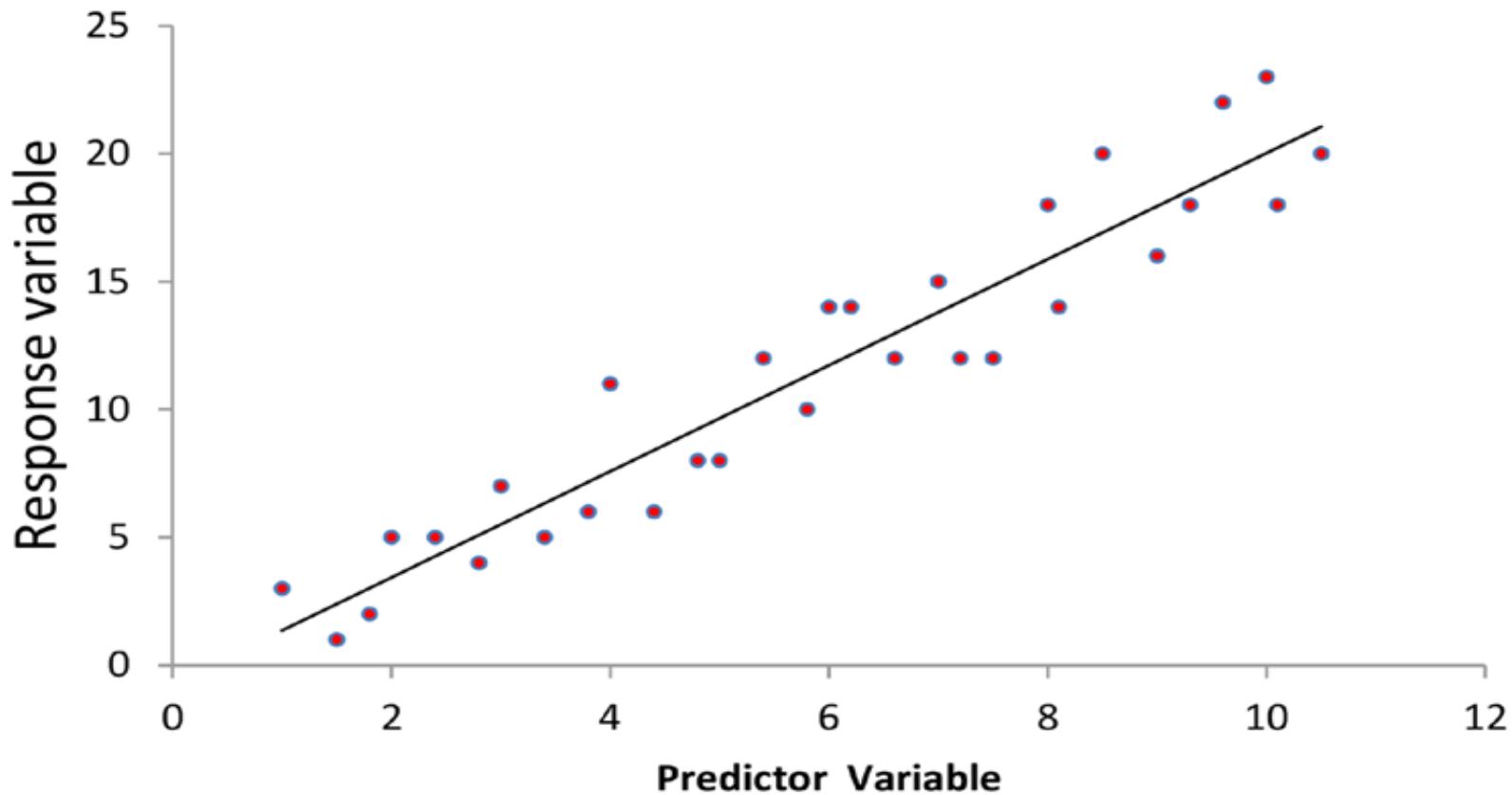
### Supervised

- Regression
  - Linear
  - Polynomial
- Decision Trees
- Random Forests
- Classification
  - KNN
  - Trees
  - Logistic Regression
  - Naive-Bayes
  - SVM

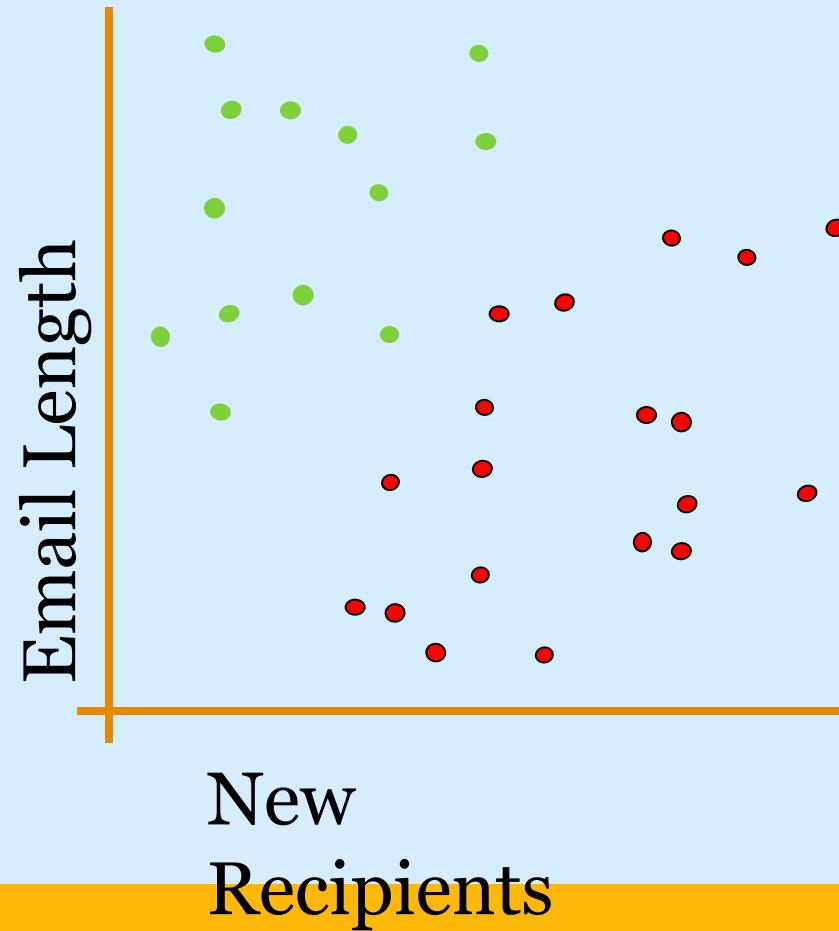
# Supervised Learning



# Linear Regression

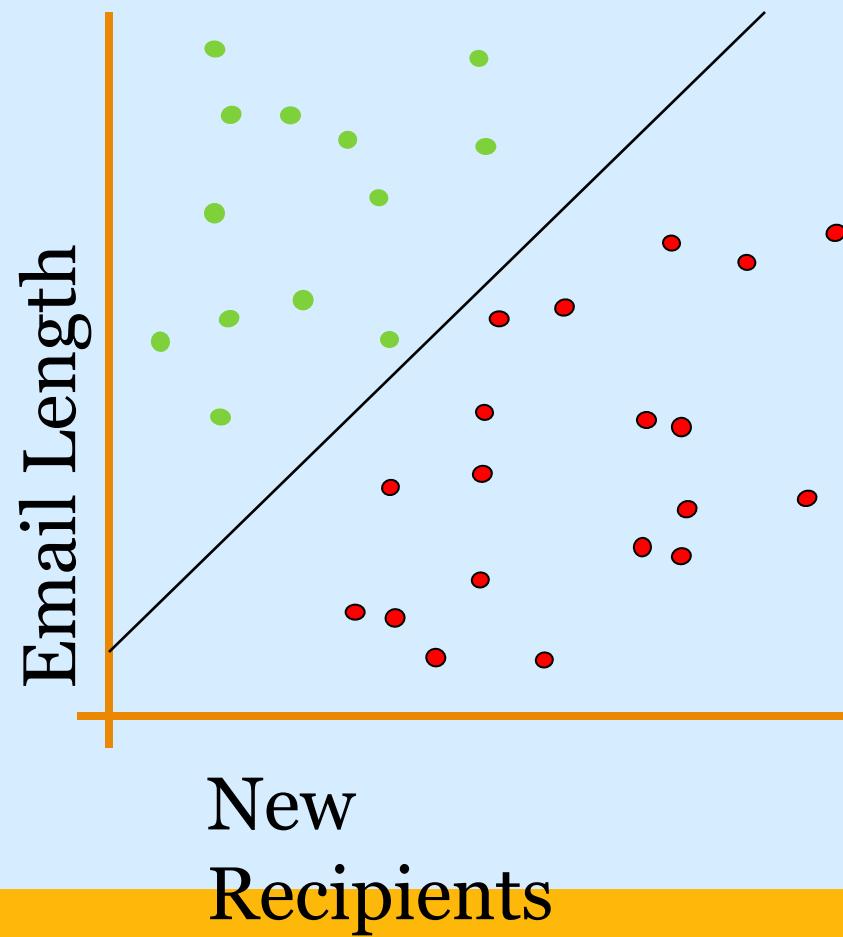


# Classifiers



How would you  
classify this data?

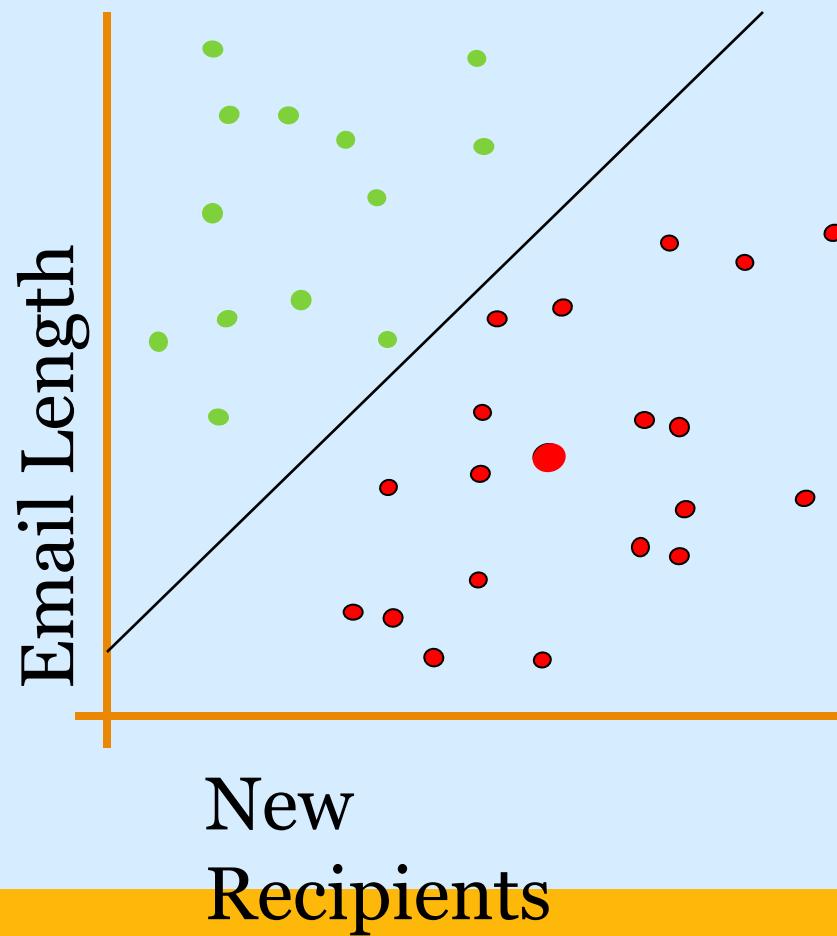
# Classifiers



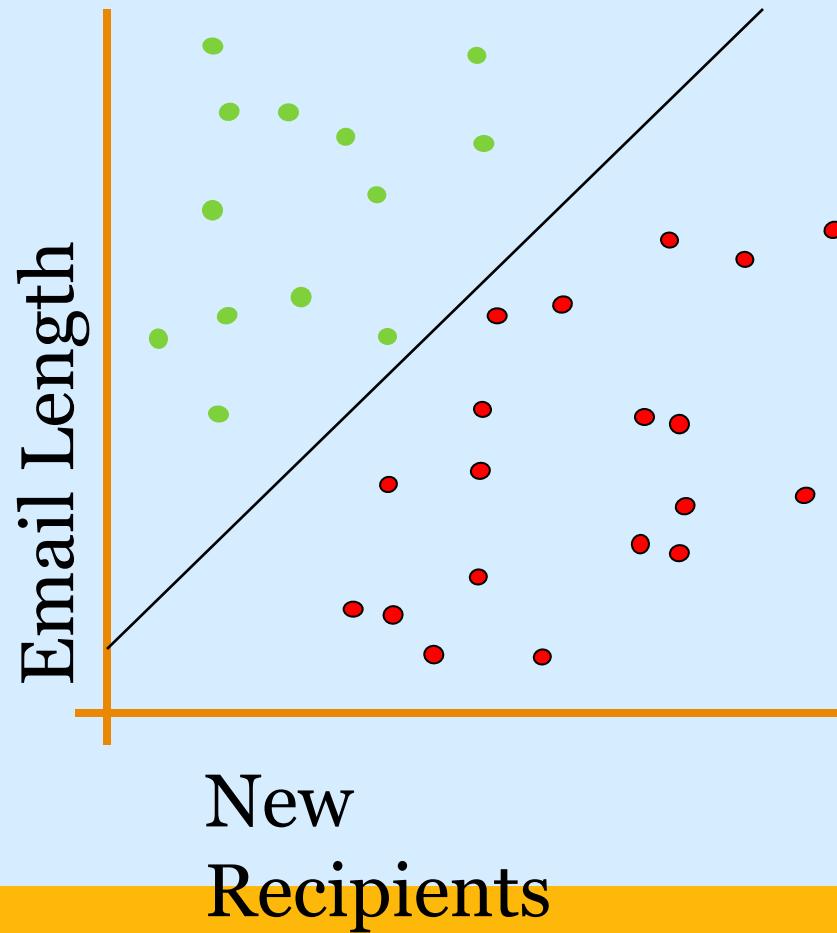
How would you  
classify this data?

# When a new email is sent

1. We first place the new email in the space
2. Classify it according to the subspace in which it resides

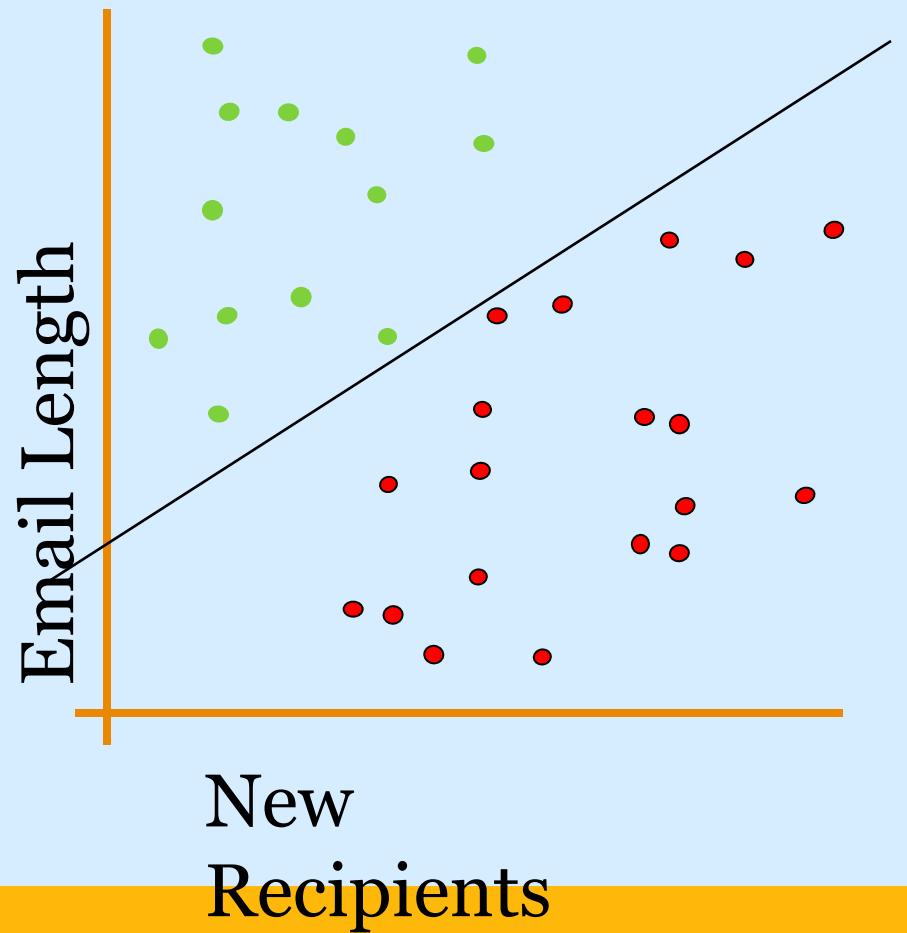


# Linear Classifiers



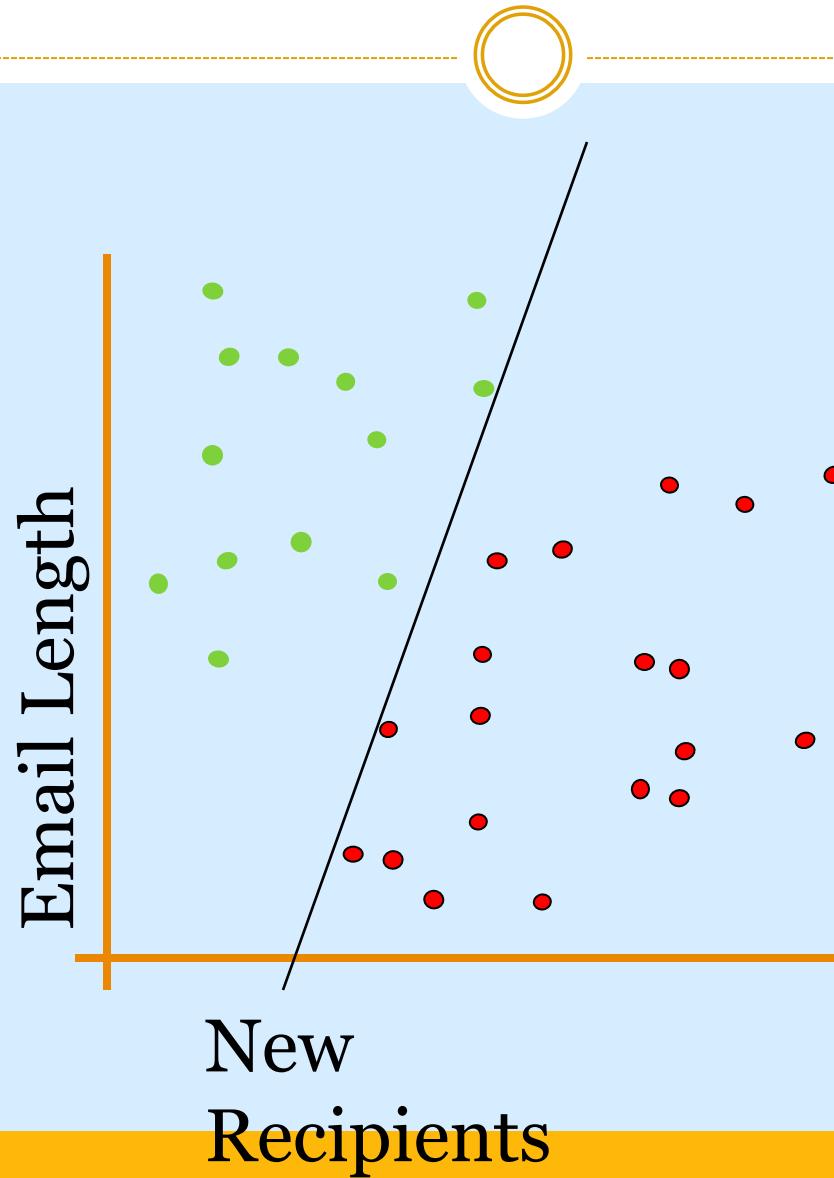
How would you  
classify this data?

# Linear Classifiers



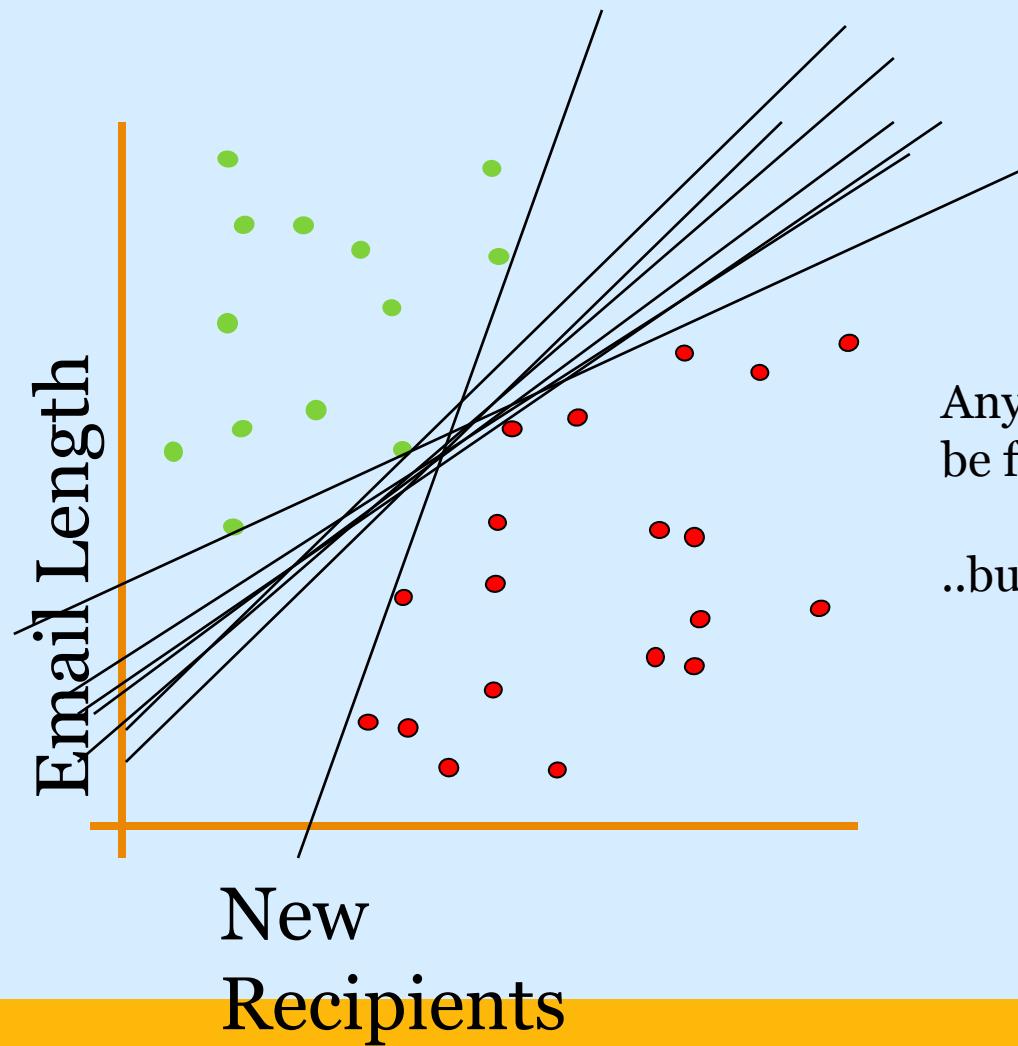
How would you  
classify this data?

# Linear Classifiers

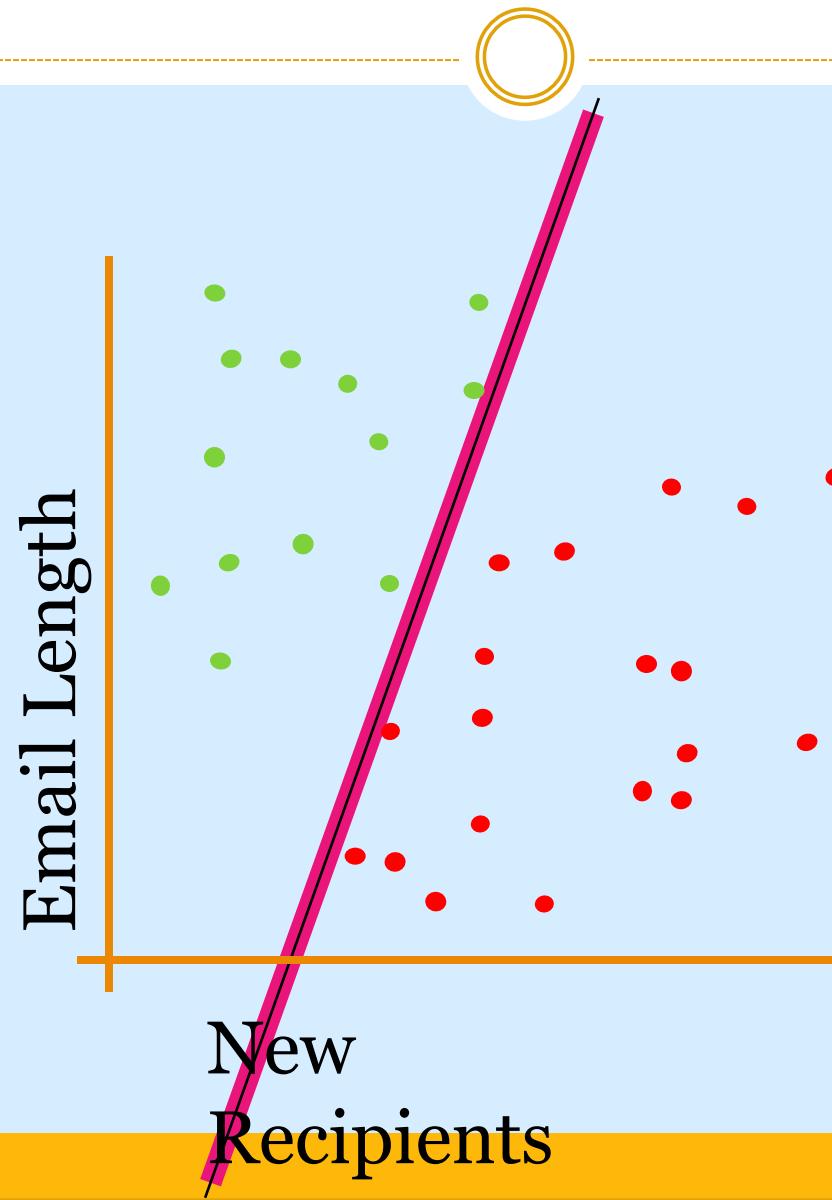


How would you  
classify this data?

# Linear Classifiers

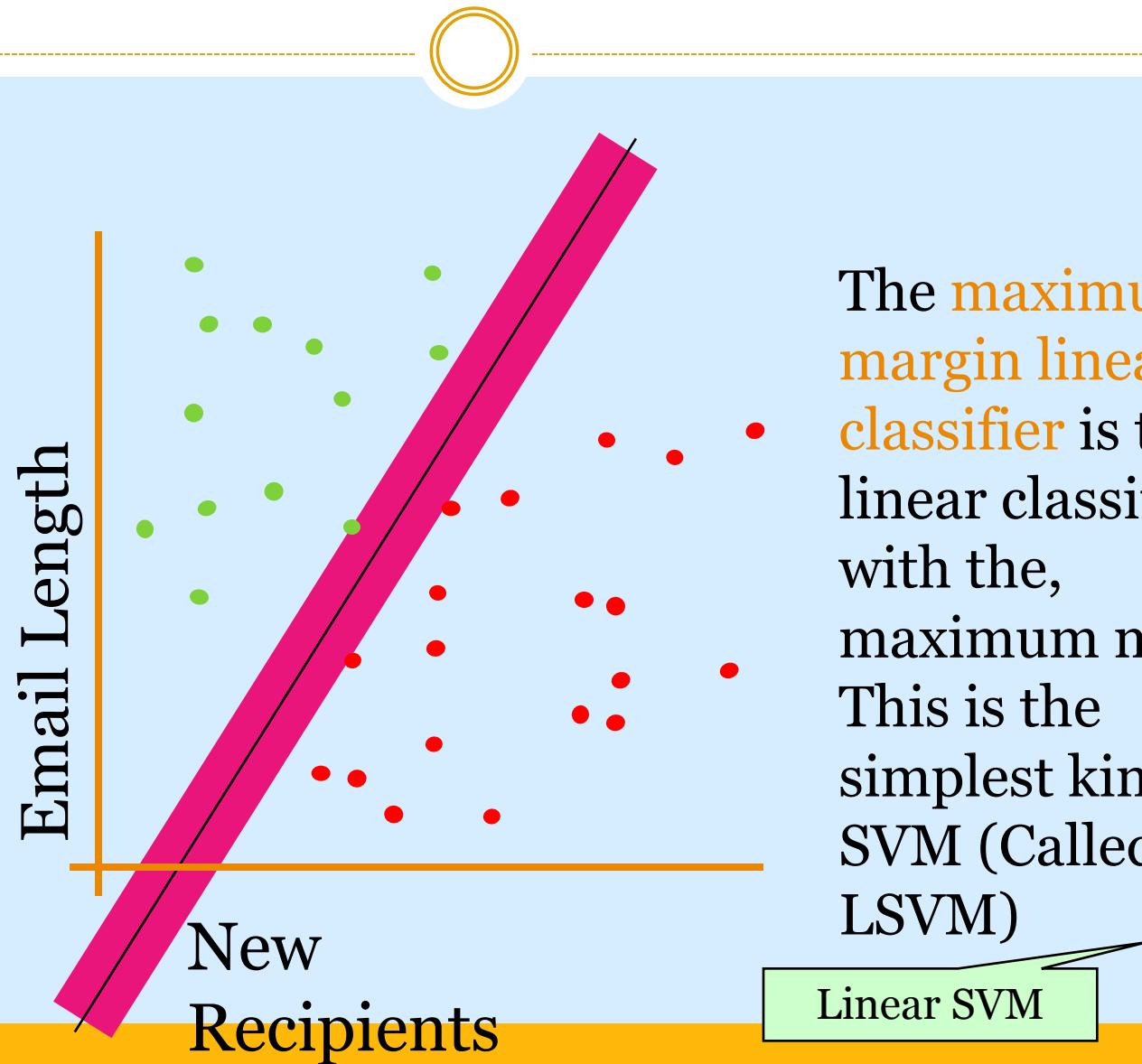


# Classifier Margin

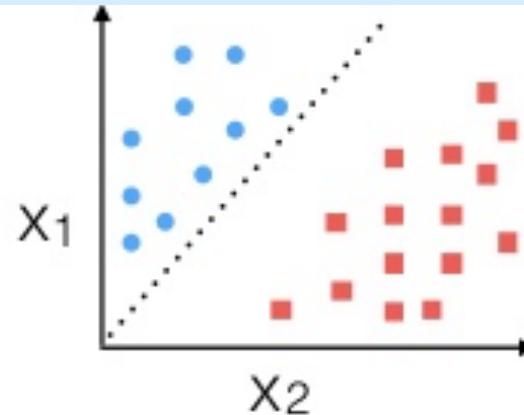
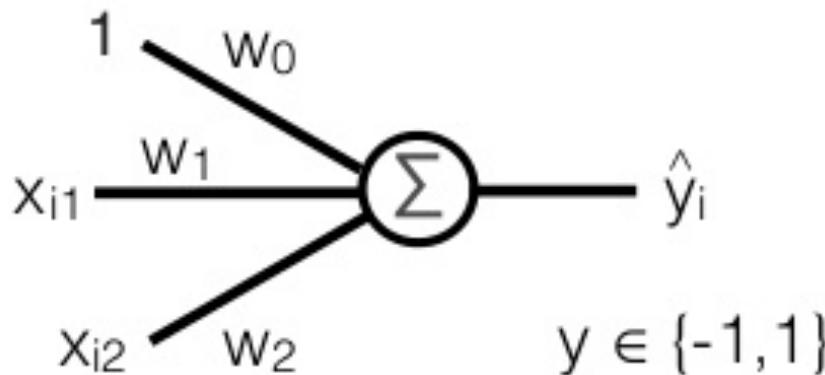


Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

# Maximum Margin



# General Algorithms for this class



$$\hat{y} = \mathbf{w}^T \mathbf{x} = w_0 + w_1 x_1 + w_2 x_2$$

$$\hat{y}_i \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x}_i \geq \theta \\ -1 & \text{otherwise} \end{cases}$$

$w_j$  = weight

$x_i$  = training sample

$y_i$  = desired output

$\hat{y}_i$  = actual output

$t$  = iteration step

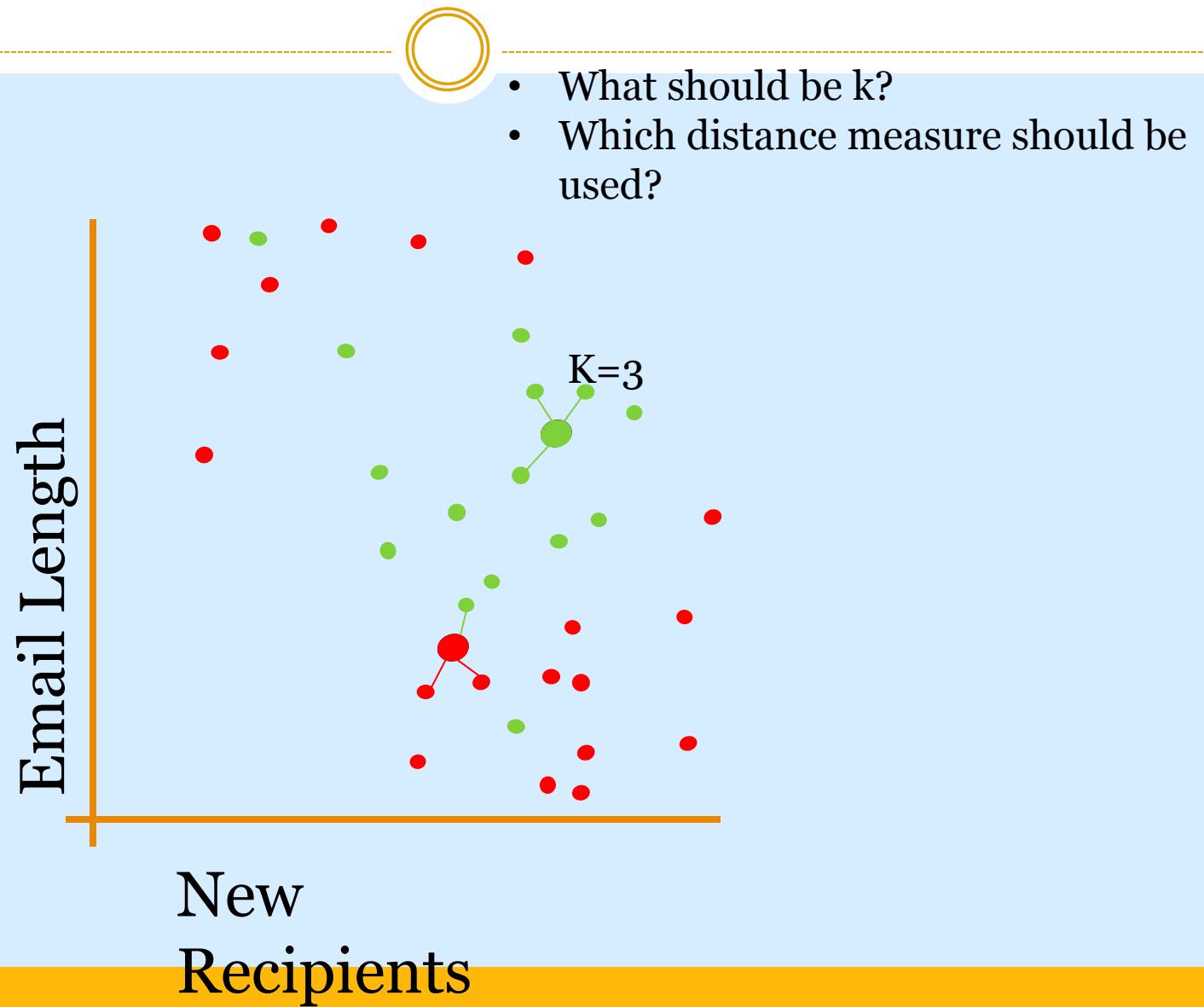
$\eta$  = learning rate

update rule:

$$w_j(t+1) = w_j(t) + \eta(y_i - \hat{y}_i)x_i$$

until  
 $t+1 = \text{max iter}$   
or error = 0

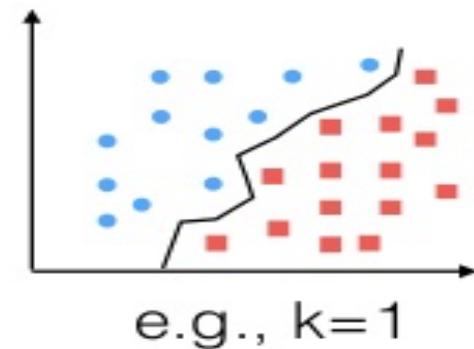
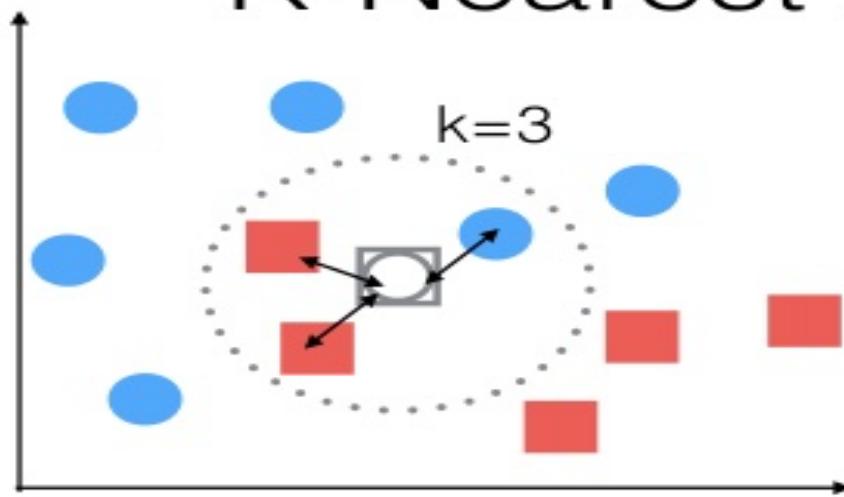
# Another classifier: k-Nearest Neighbors



# Non-Parametric Classifiers

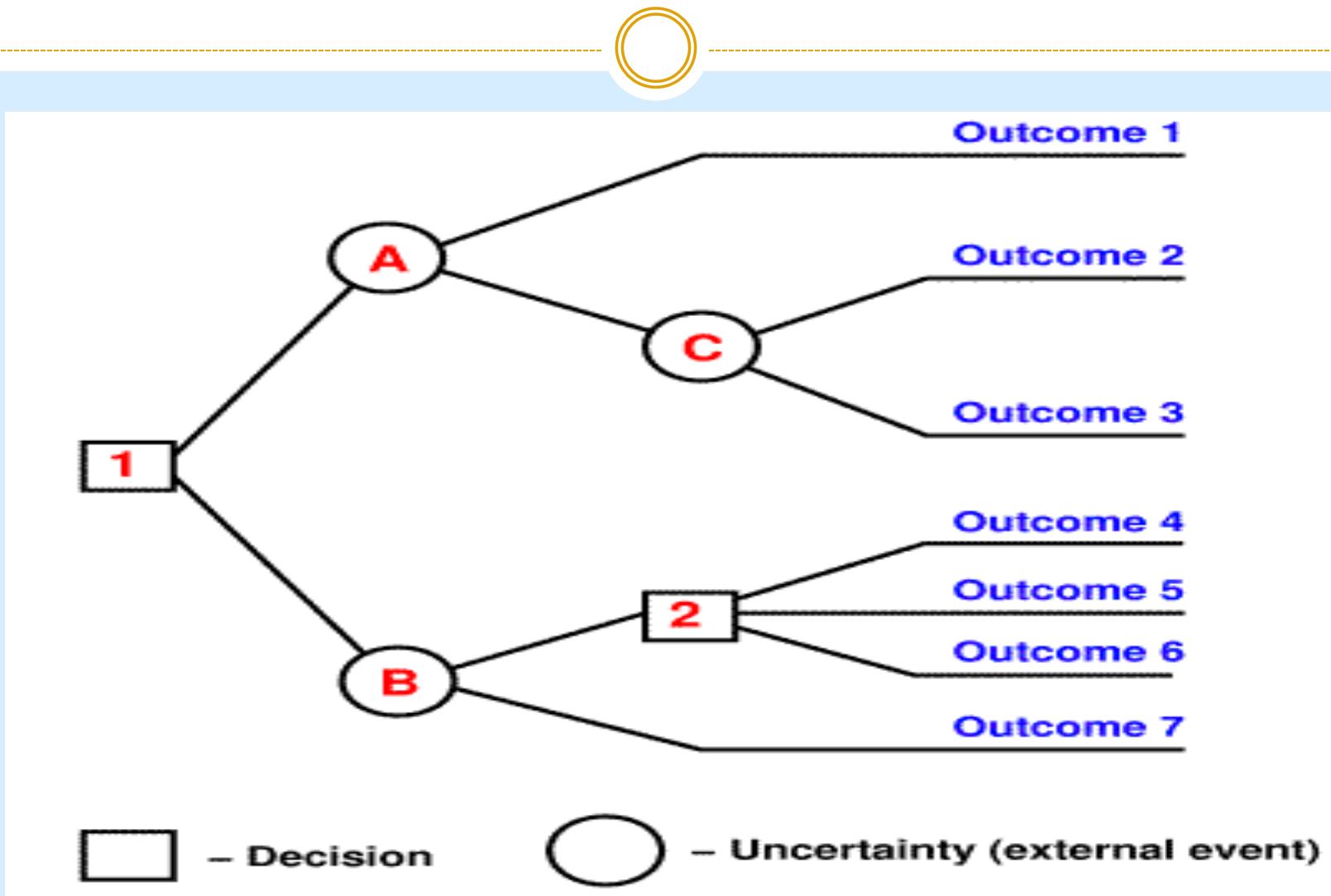


## Non-Parametric Classifiers: K-Nearest Neighbor

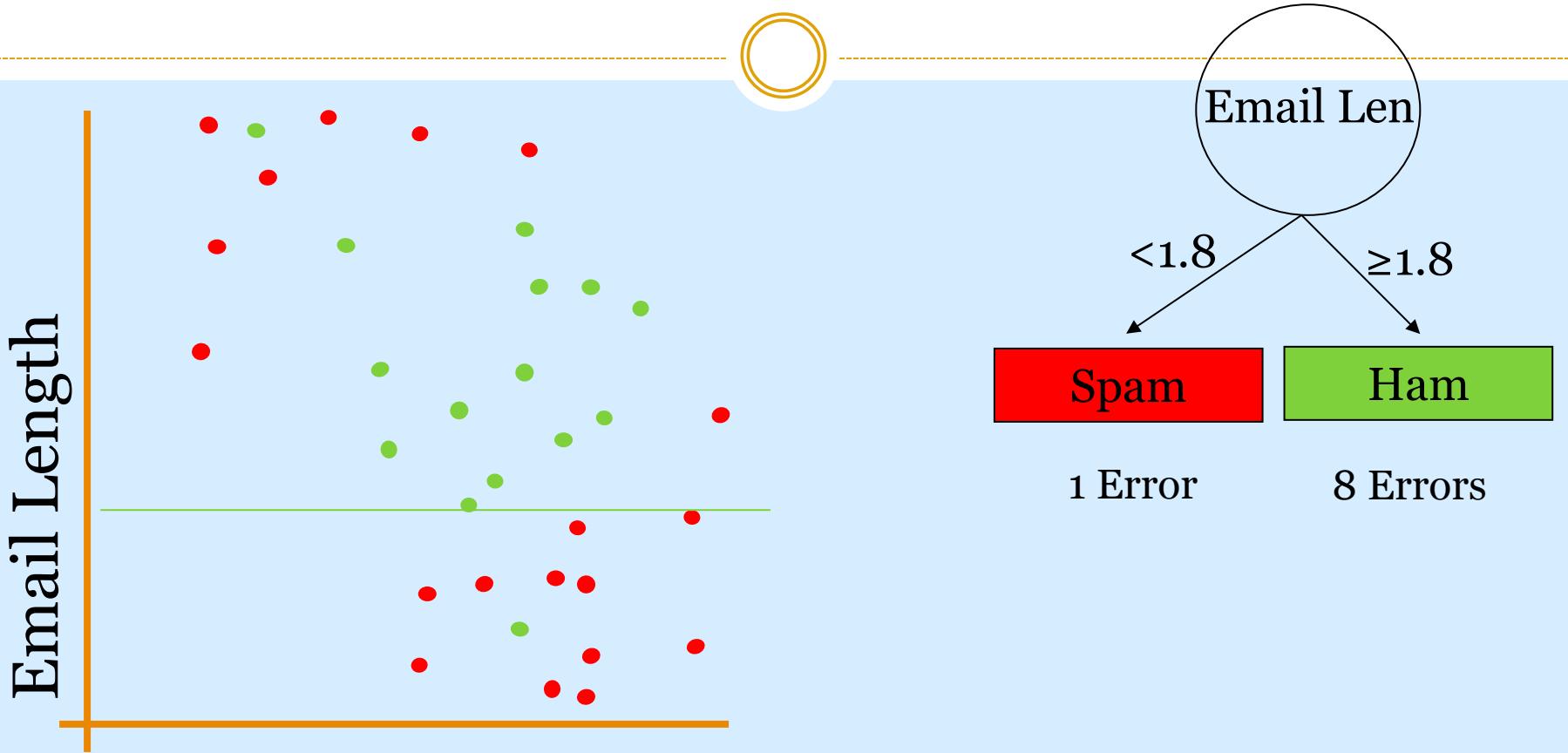


- Simple!
- Lazy learner
- Very susceptible to curse of dimensionality

# Decision Tree



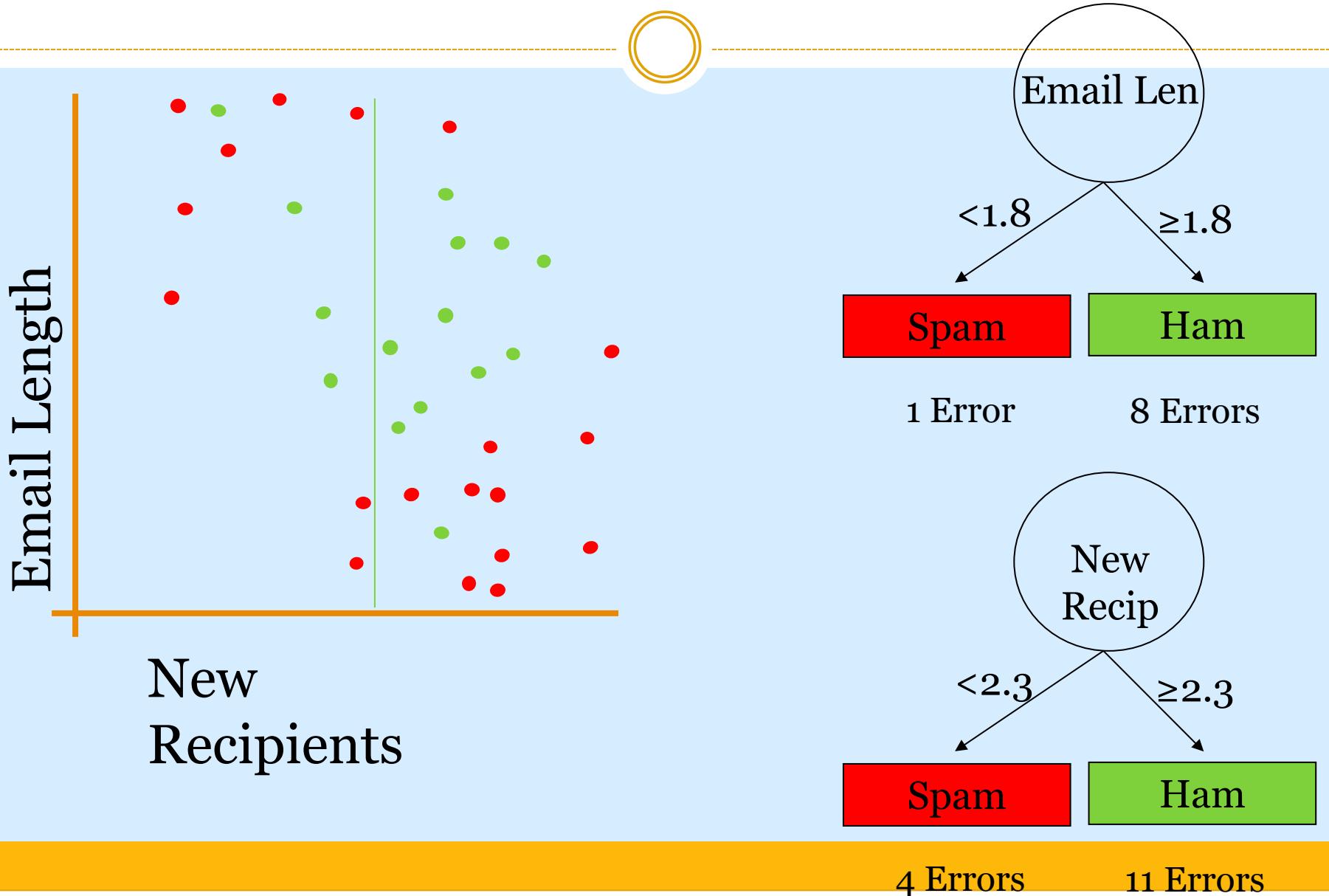
# Top Down Induction of Decision Trees



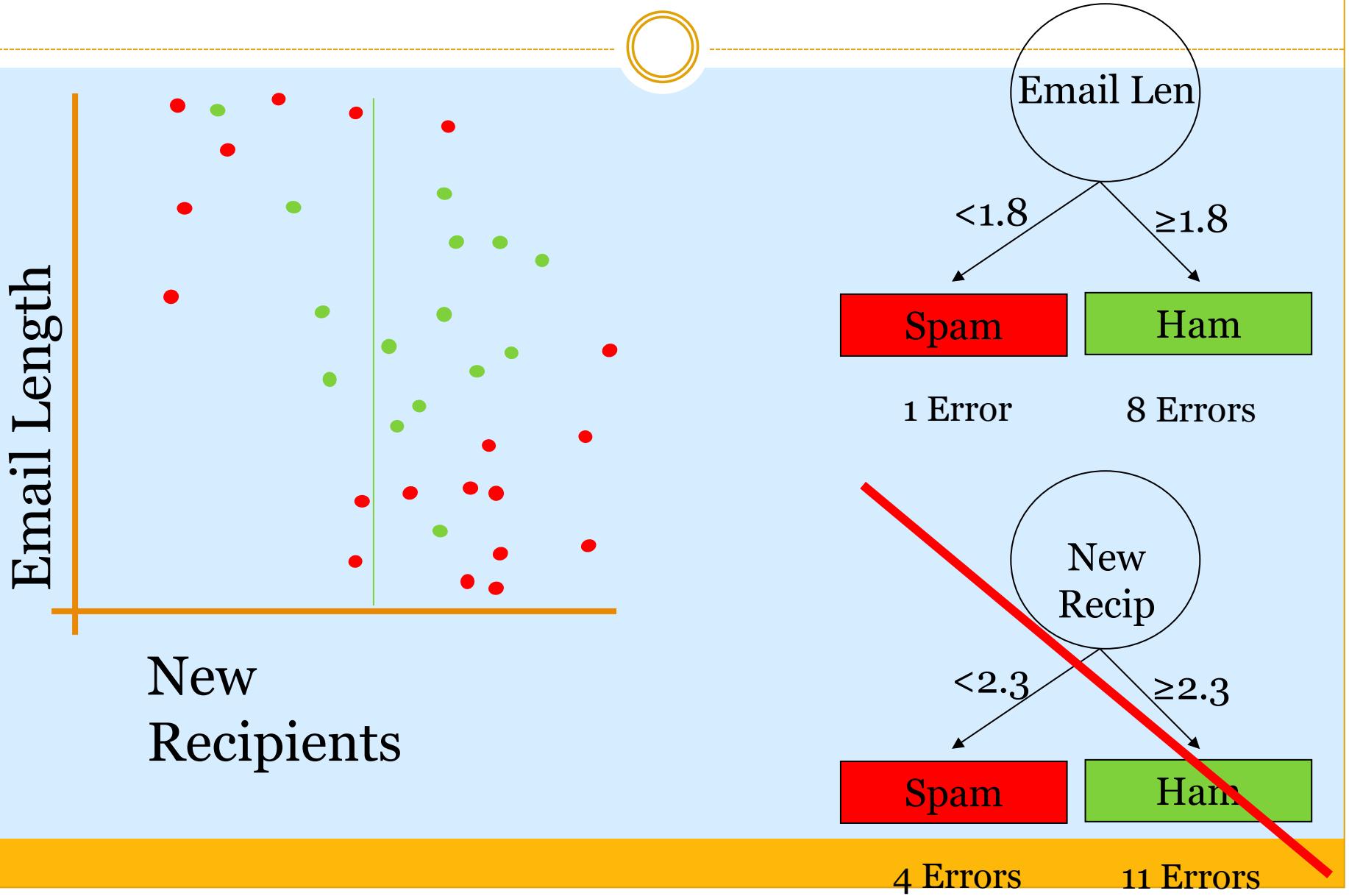
New  
Recipients

A single level decision tree is also known as Decision Stump

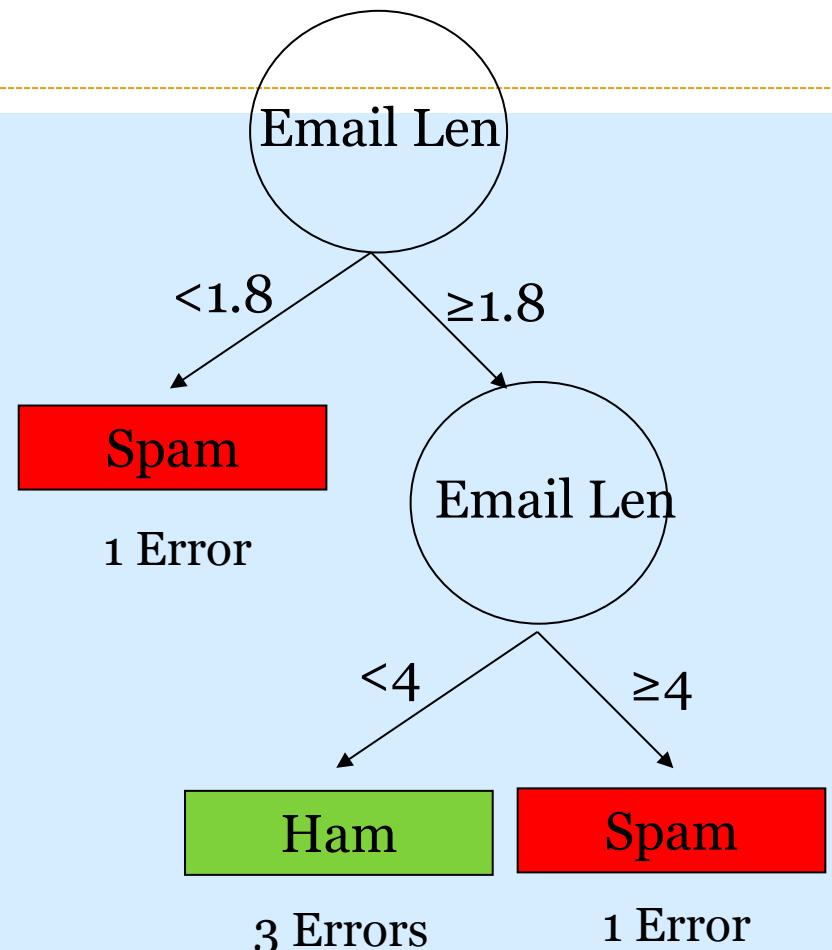
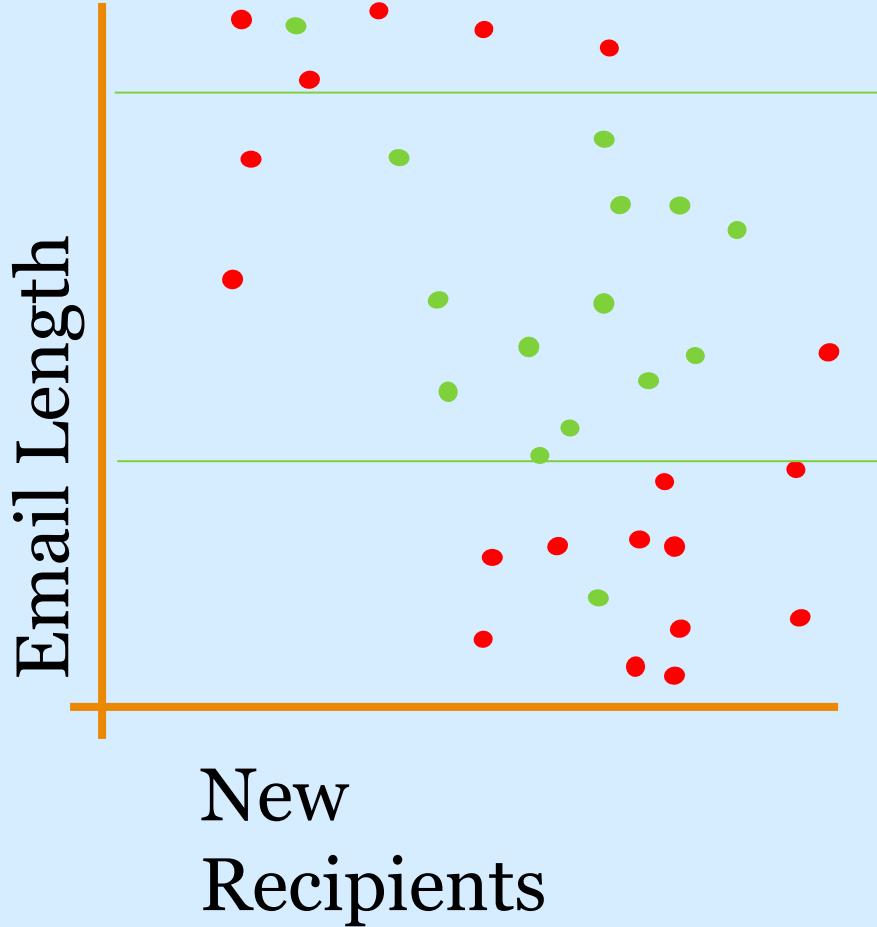
# Top Down Induction of Decision Trees



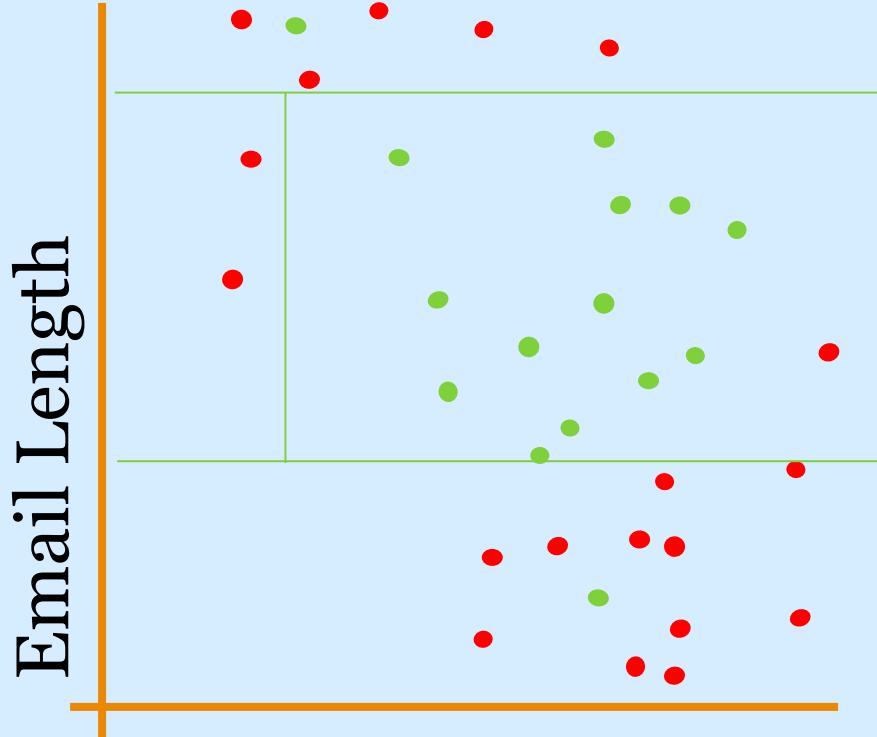
# Top Down Induction of Decision Trees



# Top Down Induction of Decision Trees



# Top Down Induction of Decision Trees



New  
Recipients



Email Len

<1.8      ≥1.8

Spam

1 Error

Email Len

<4      ≥4

Spam

1 Error

New  
Recip

<1

≥1

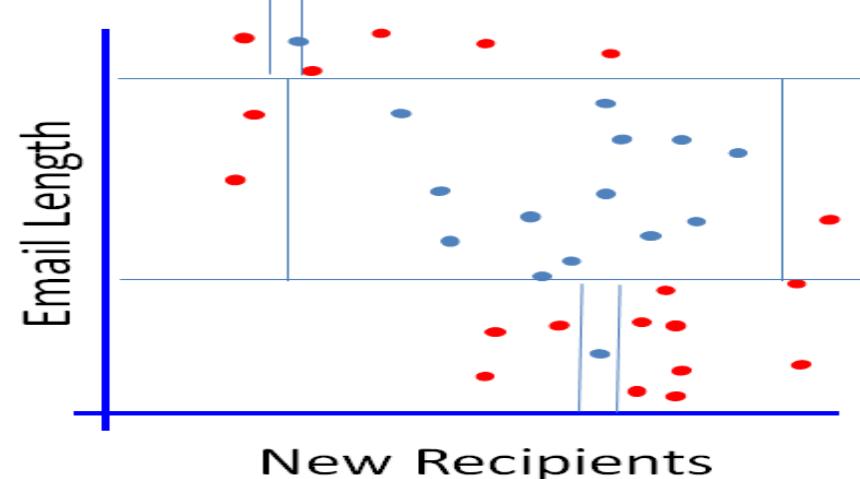
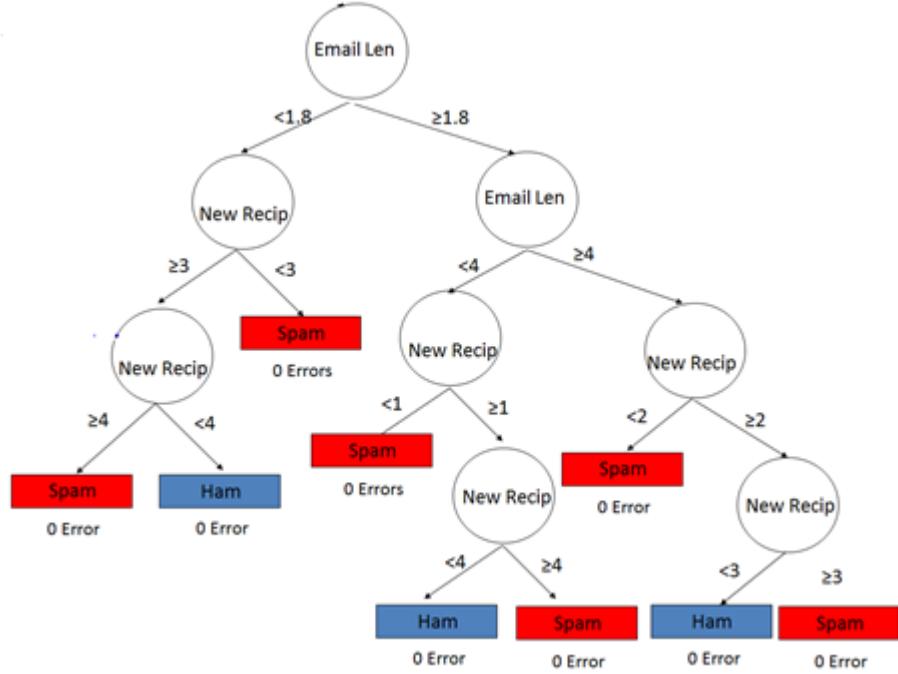
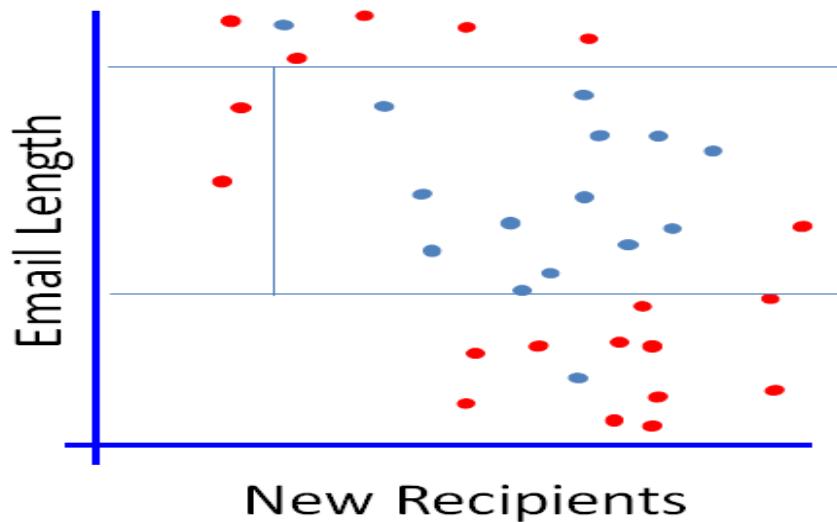
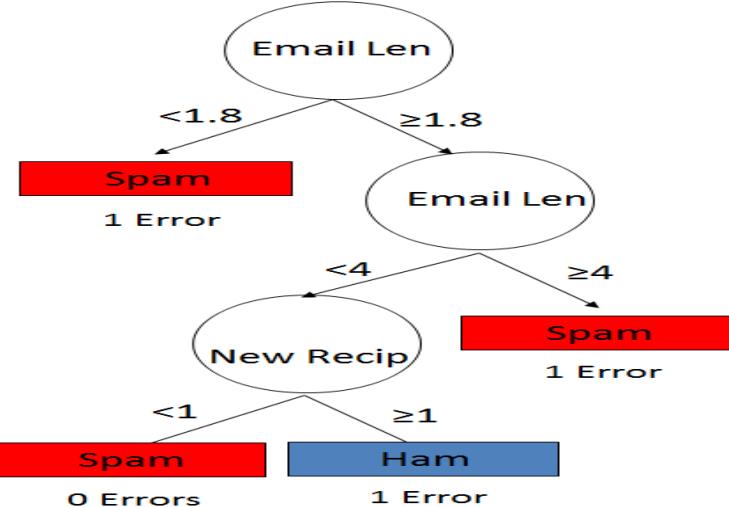
Spam

Ham

0 Errors

1 Error

# Which One?



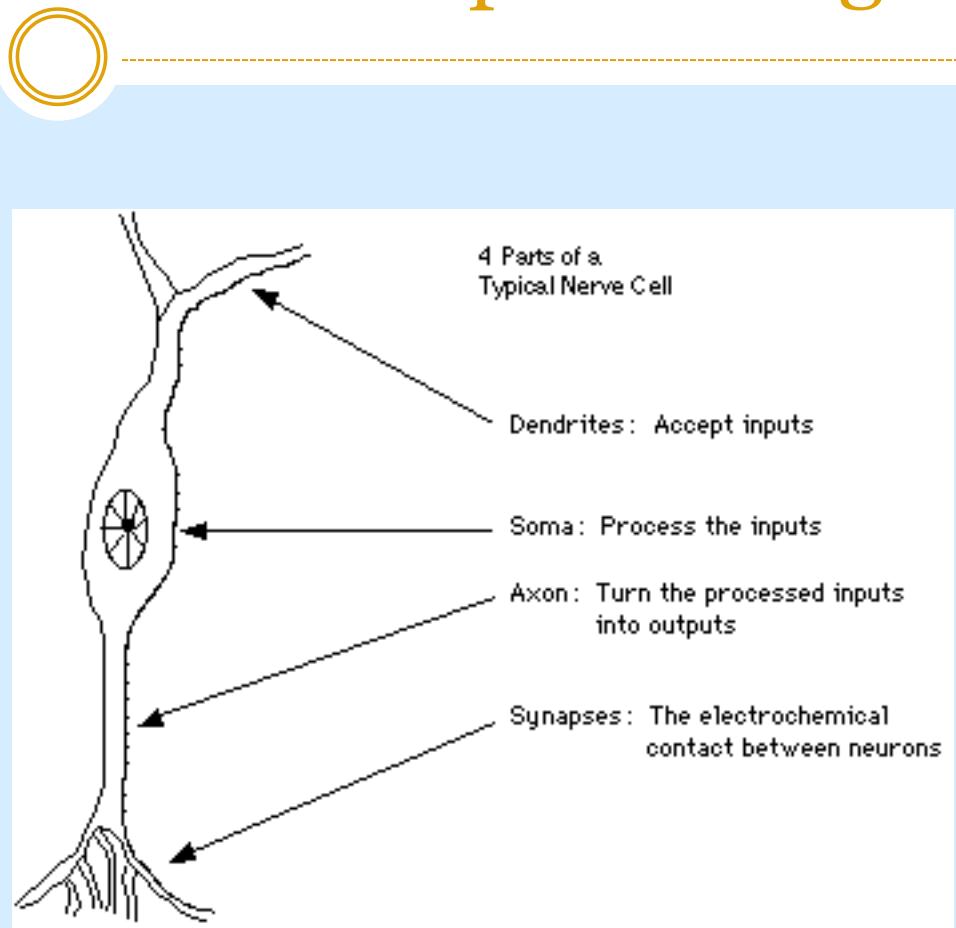
# Moving Forward



- Earlier algorithms discussed form general branch of Machine Learning
- ML is practice of using algorithms to parse data, learn from it, then determine or make predictions
- Another approach is mimicking human brain
- Entire branch of Aritifical Intelligence based on this approach is known as Artificial Neural Networks and Deep Learning
- It is composed of a large number of highly interconnected processing elements called neurons

# Let us simply understand brain processing

- Four parts of a nerve cell:
- Dendrite: Accepts the input
- Soma: Processes the input
- Axon: Turns processed input into output
- Synapses: Connection between neurons



# Our brain is a huge neural network



# Moving Forward



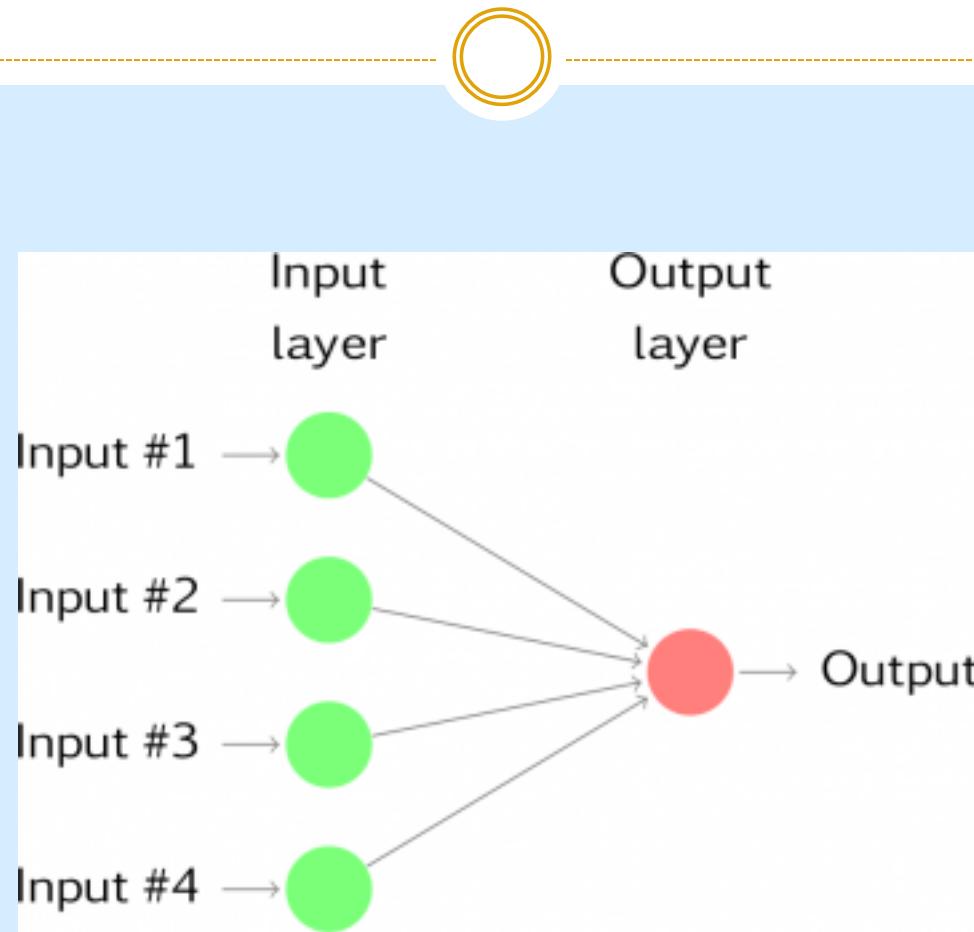
- Artificial Neural Network (ANNs) are programs designed to solve any problem by trying to mimic the structure and the function of our nervous system.
- Neural networks are based on simulated neurons, Which are joined together in a variety of ways to form networks.
- Neural network resembles the human brain in the following two ways: -
  - A neural network acquires knowledge through learning.
  - A neural network's knowledge is stored within the interconnection strengths known as synaptic weight.

# Artificial Neural Network (ANN)

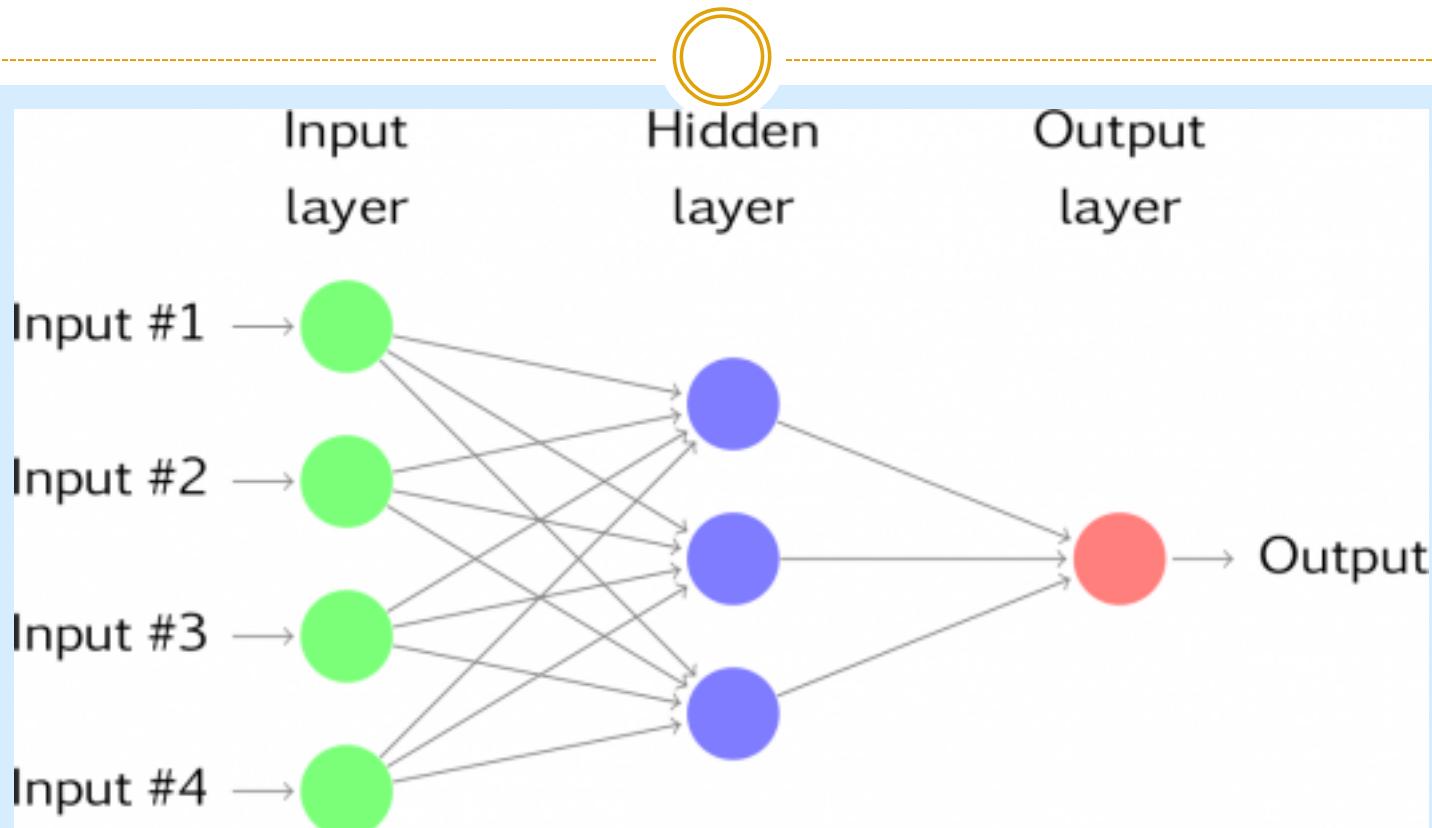


- An ANN is based on the same idea
- There are inputs to the network ( $x_1, x_2 \dots$ )
- Each of these inputs is multiplied by a connection weight ( $w_1, w_2 \dots$ )
- The products are then summed and fed to a function to generate output
- $\text{Sum} = w_1x_1 + w_2x_2 + ..$

# Neural Network



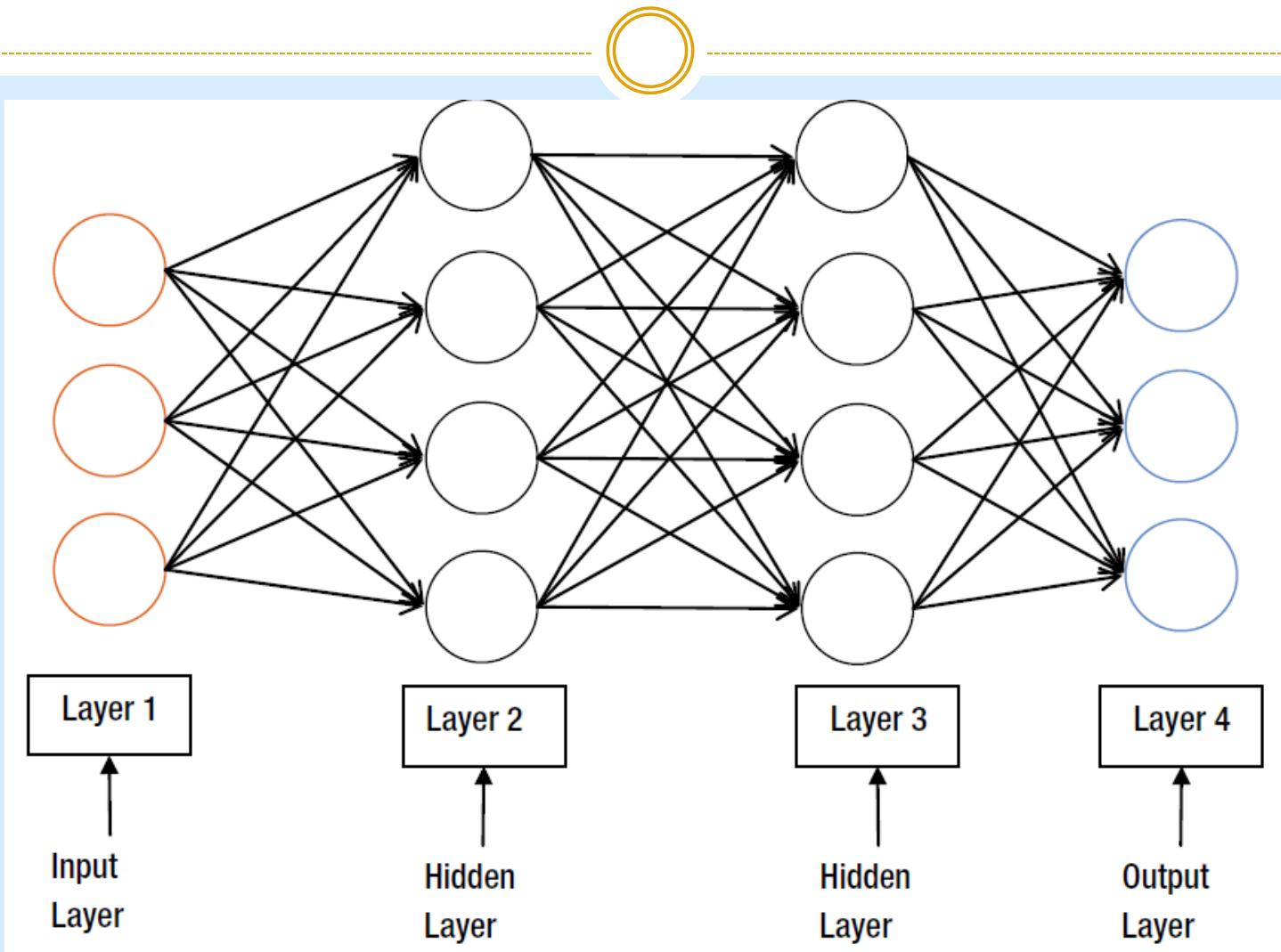
# Let us add complexity



$$z_j = b_j + \sum_{i=1}^4 w_{i,j}x_i$$

$$s(z) = \frac{1}{1 + e^{-z}}$$

# Neural Network



# Illustrative Example



- For example, a 3-input neuron is taught to output 1 when the input ( $X_1, X_2$  and  $X_3$ ) is 111 or 101 and to output 0 when the input is 000 or 001

X 1:		0	0	0	0	1	1	1	1
X 2:		0	0	1	1	0	0	1	1
X 3:		0	1	0	1	0	1	0	1
O U T:		0	0	0/ 1	0/ 1	0/ 1	1	0/ 1	1

# Illustrative Example



- Take the pattern 010. It differs from 000 in 1 element, from 001 in 2 elements, from 101 in 3 elements and from 111 in 2 elements. Therefore, the 'nearest' pattern is 000 which belongs in the 0-taught set. Thus the firing rule requires that the neuron should not fire when the input is 001. On the other hand, 011 is equally distant from two taught patterns that have different outputs and thus the output stays undefined (0/1).

X 1:		0	0	0	0	1	1	1	1
X 2:		0	0	1	1	0	0	1	1
X 3:		0	1	0	1	0	1	0	1
O U T:		0	0	0	0/ 1	0/ 1	1	1	1

# Typical Structure of an ANN



- The algorithm is a learning rule which suggests a way of modifying weights to represent a function from input to output
- The network architecture is a feedforward network where computational units are structured in a multi-layered network: an input layer, one or more hidden layer(s), and an output layer
- The units on a layer have full connections to units on the adjacent layers, but no connection to units on the same layer. This leads us to the idea of back propagation

# Backpropagation



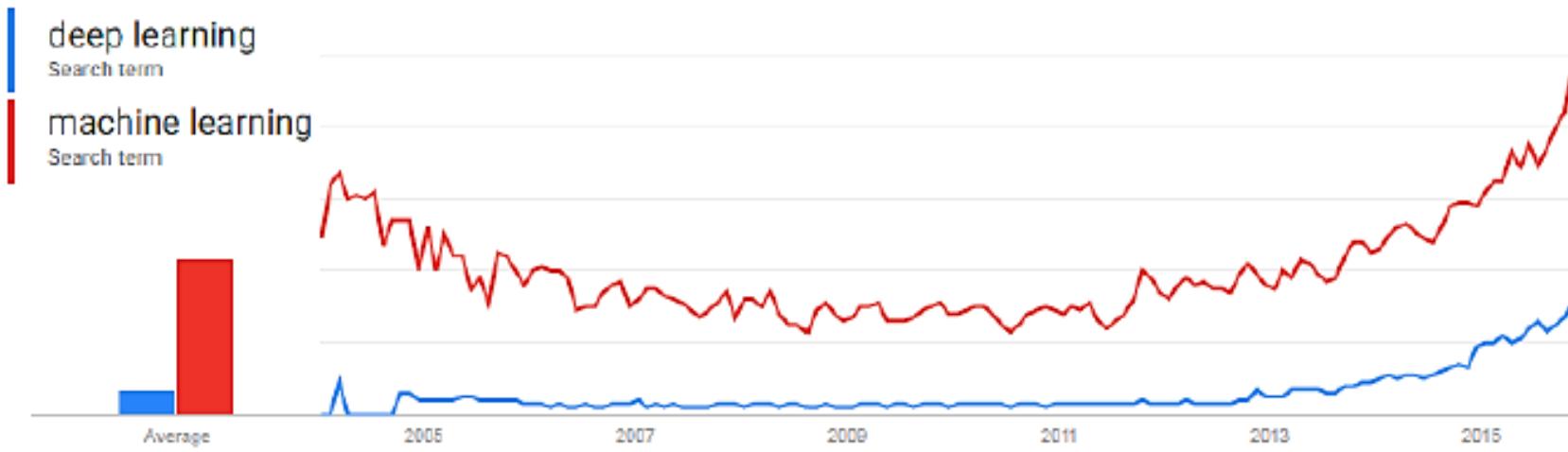
- Calculate the difference (error) between the expected and actual output value
- Adjust the weights in order to minimize the error
- Minimize the error by performing a gradient decent on the error surface
- The amount of the weight change for each input pattern in an epoch is proportional to the error

# Why use Neural Networks?



- ability to derive meaning from complicated or imprecise data
- extract patterns and detect trends that are too complex to be noticed by either humans or other computer techniques
- Adaptive learning
- Real Time Operation

# Deep Learning



# Deep Learning



## Artificial Intelligence Takes Off at Google

Number of software projects within Google that uses a key AI technology, called Deep Learning.



Source: Google

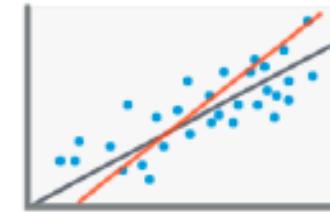
# Why this hype about Deep Learning?



- To understand this let us understand our learning till now:
- Supervised Learning: Learning with a labeled training set. *Example: email spam detector with training set of already labeled emails*
- Unsupervised Learning: Discovering patterns in unlabeled data *Example: cluster similar documents based on the text content*
- Reinforcement Learning: learning based on feedback or reward *Example: learn to play chess by winning or losing*



Classification  
(supervised – predictive)



Regression  
(supervised – predictive)

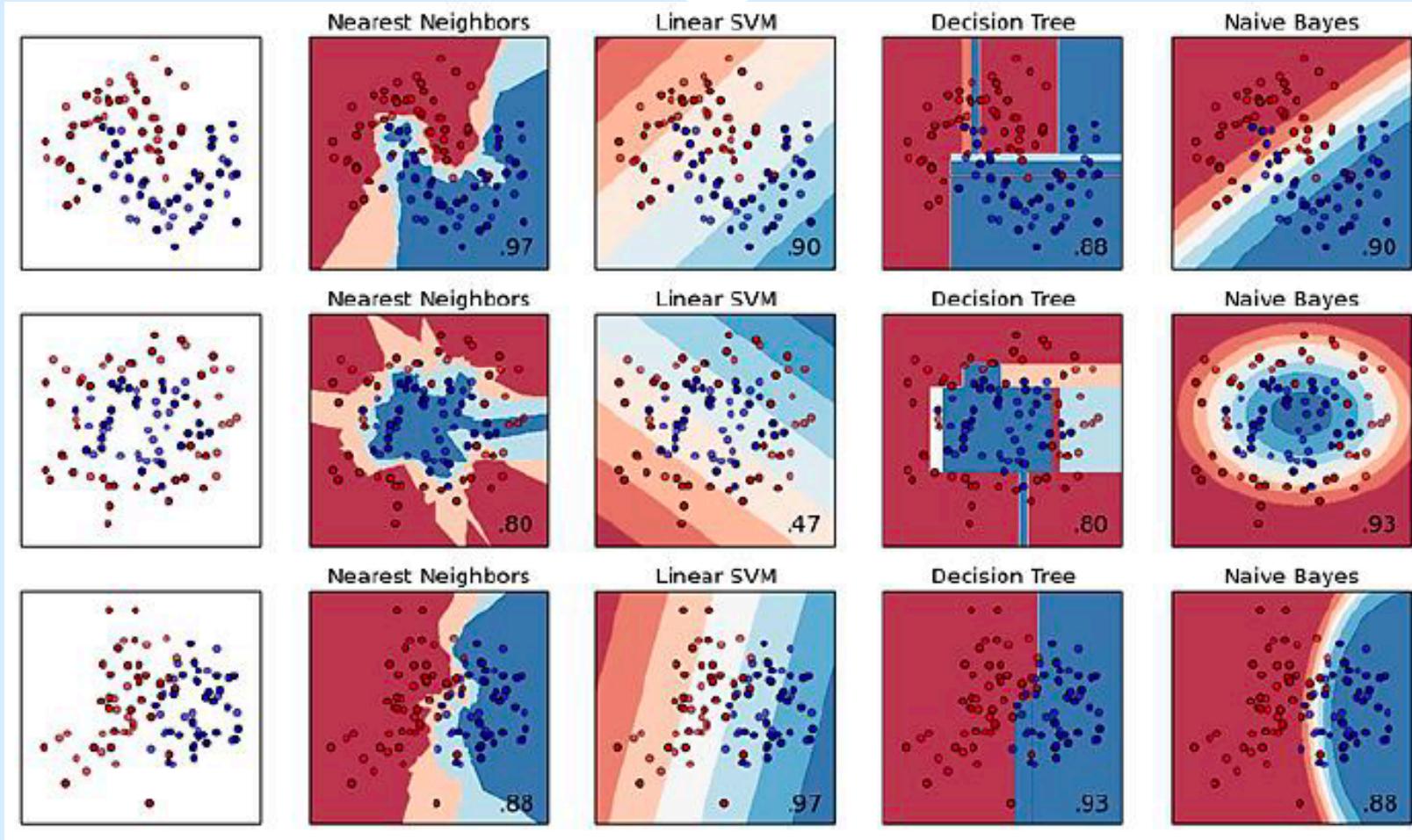


Clustering  
(unsupervised – descriptive)



Anomaly Detection  
(unsupervised – descriptive)

# Comparison of algorithmic accuracy



# Deep Learning



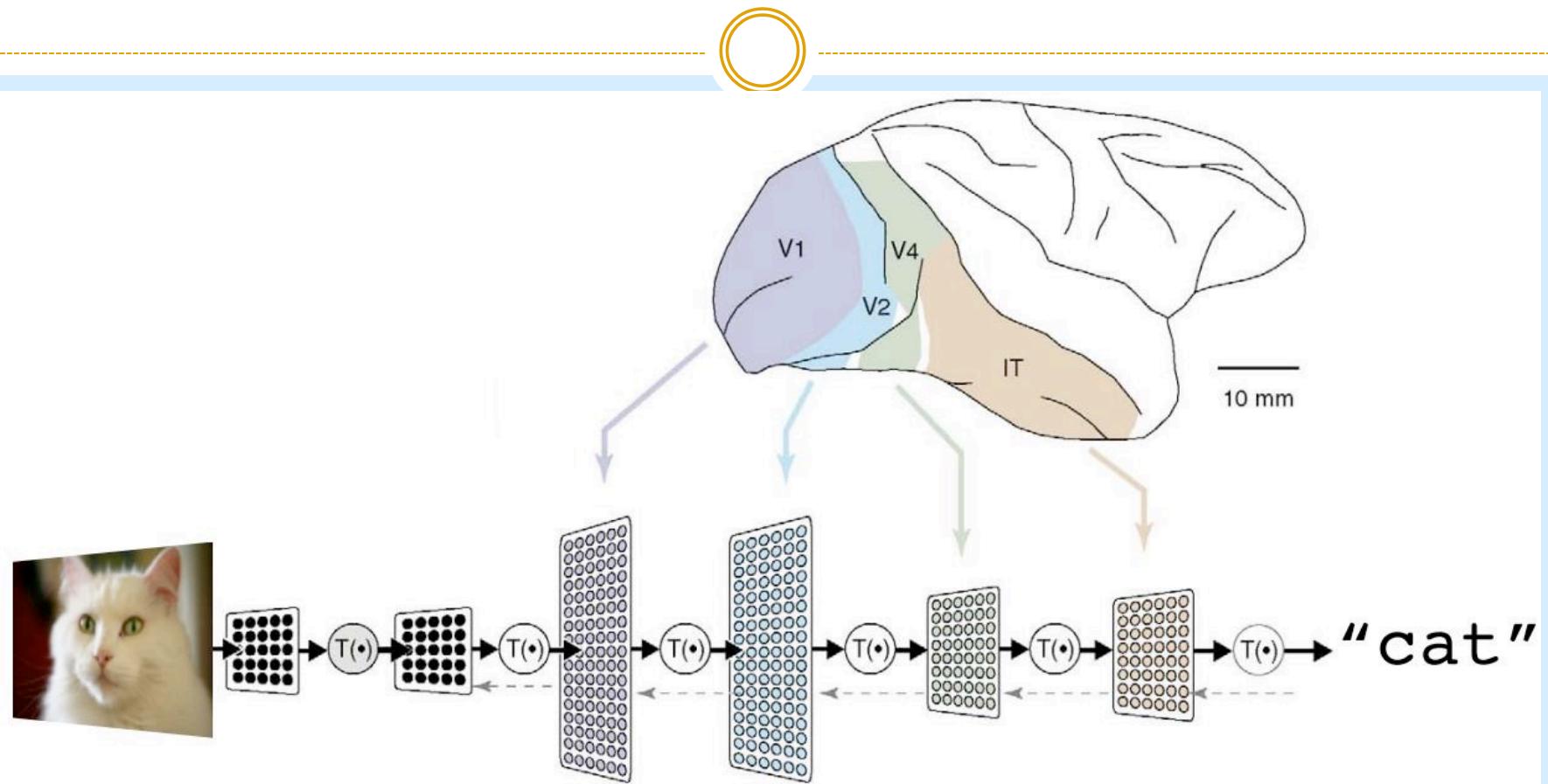
- So what is the idea behind Deep Learning?
- Part of the machine learning field of learning representations of data. Exceptionally effective at learning patterns.
- Utilizes learning algorithms that derive meaning out of data by using a hierarchy of multiple layers that mimic the neural networks of our brain.
- If you provide the system tons of information, it begins to understand it and respond in useful ways.

# Deep Learning

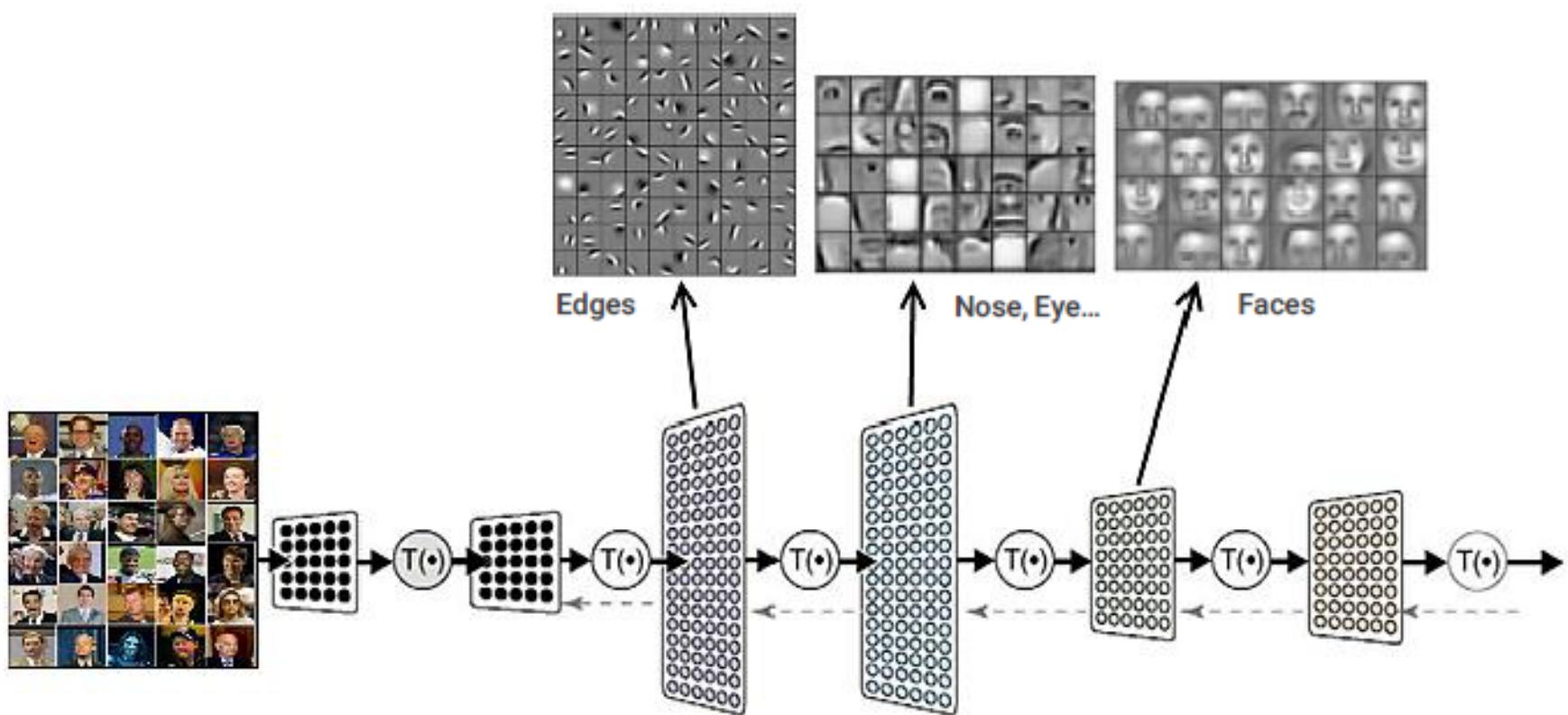


- Conceptually, similar to ANN
- Hierarchy of layers (think layers in ANN)
- Each layer transforms image/input to abstract representations
- Output layer combines these abstract layers to form predictions
- Popular algos: Convolutional Neural Networks, TensorFlow, Theano, DGCAN

# Deep Learning



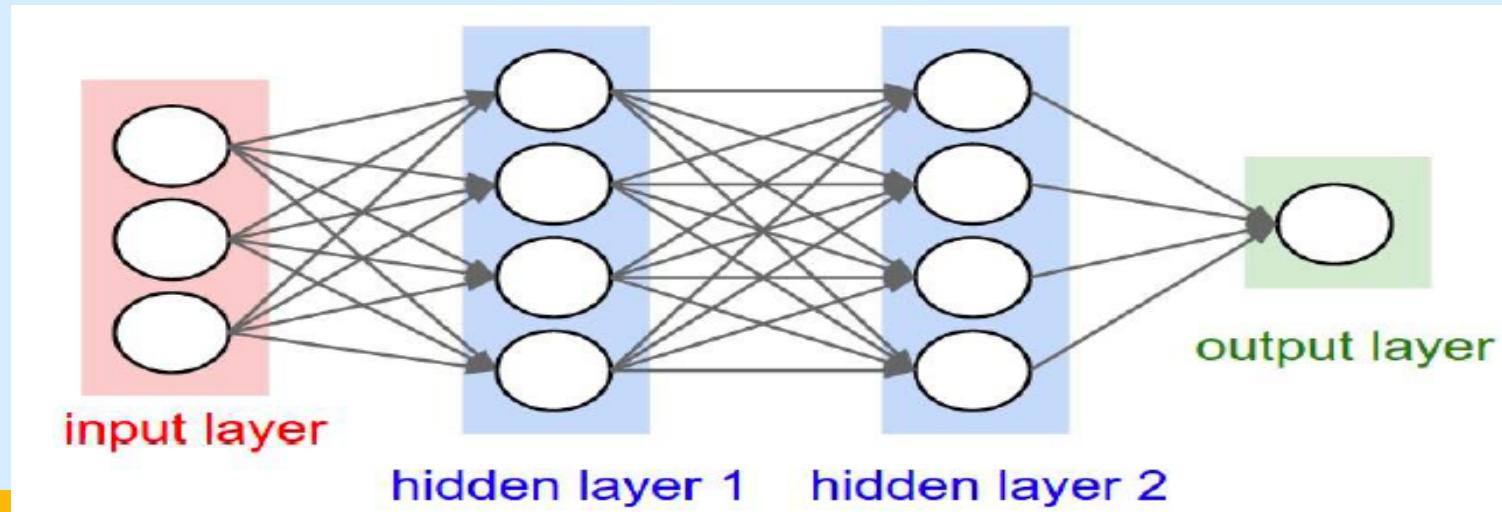
# Deep Learning



# Still very similar to ANN



- Consists of one input, one output and multiple fully-connected hidden layers in-between. Each layer is represented as a series of neurons and progressively extracts higher and higher-level features of the input until the final layer essentially makes a decision about what the input shows.



# Training Process



- Learns by generating an error signal that measures the difference between the predictions of the network and the desired values and then using this error signal to change the weights(or parameters) so that predictions get more accurate.



# Deep Learning



Human captions from the training set



Automatically captioned



# Deep Learning



Older      Mouth Open      Eyes Open      Smiling      Moustache      Glasses



# Applications of Deep Learning



- Speech recognition, Computer Vision, NLP

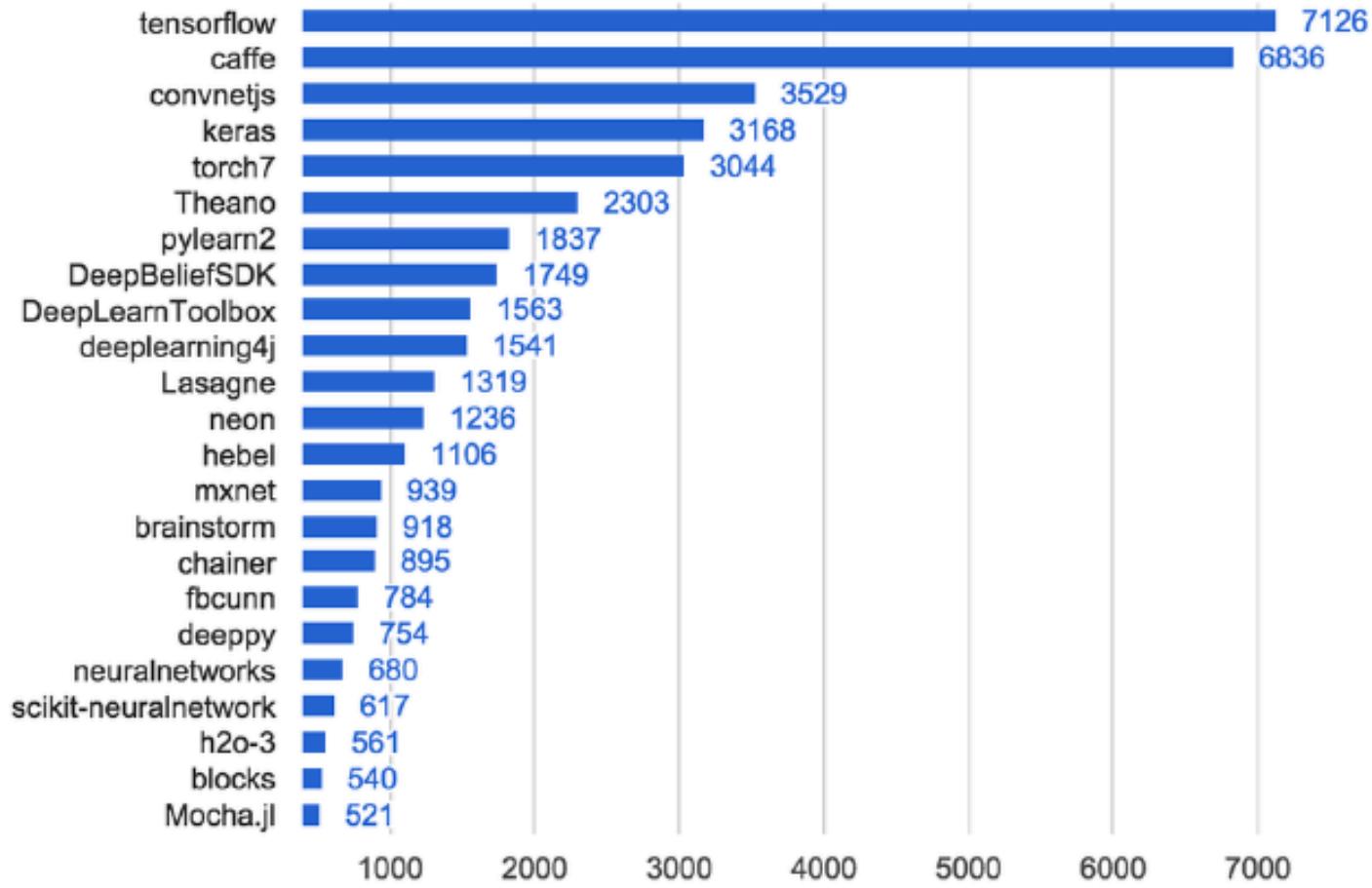


# Usage Requirements



- Large dataset with good quality input-output mappings
- Measurable and describable goals
- Enough computing power – generally GPUs are needed
- Excels in tasks where the basic unit (*pixel, word*) has very little meaning in itself, but the combination of such units has a useful meaning

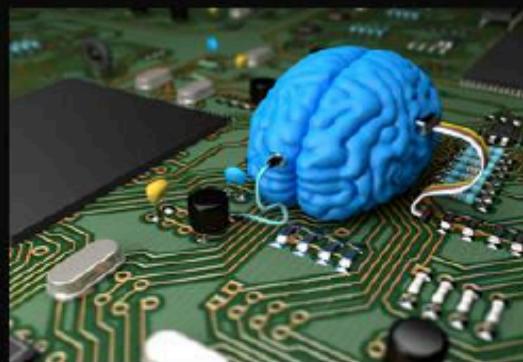
# Some Popular Packages



# Deep Learning



What society thinks I do



What my friends think I do



What other computer  
scientists think I do



What mathematicians think I do



What I think I do

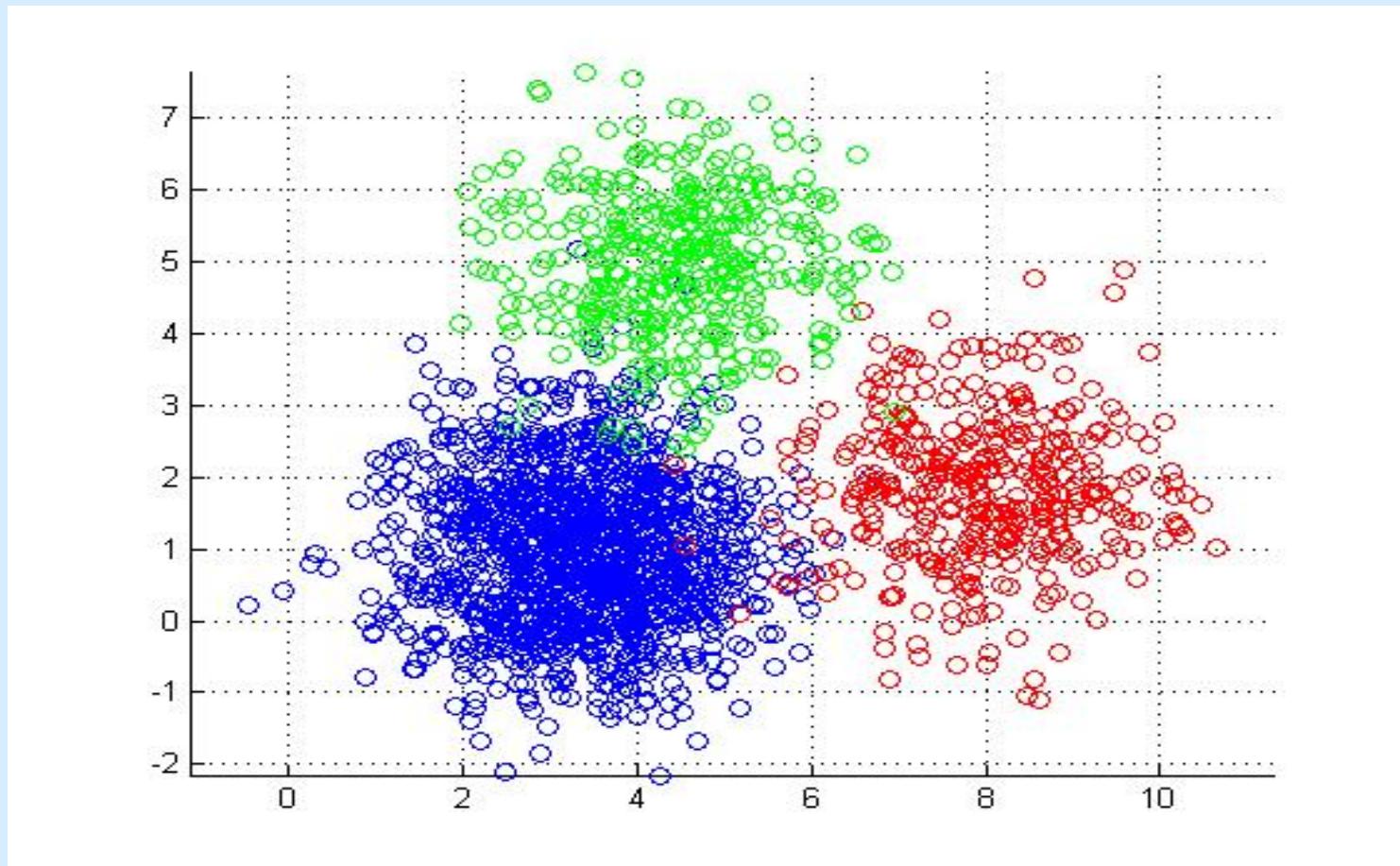
```
from theano import *
```

What I actually do

# Unsupervised Algorithm



- Think Clustering
- Used when dataset is unlabelled
- Goal is to cluster observations in meaningful groups
- Try to place observations in different buckets
  - Learning patterns inherent in data
  - PCA, SVD, k-means



# Applications

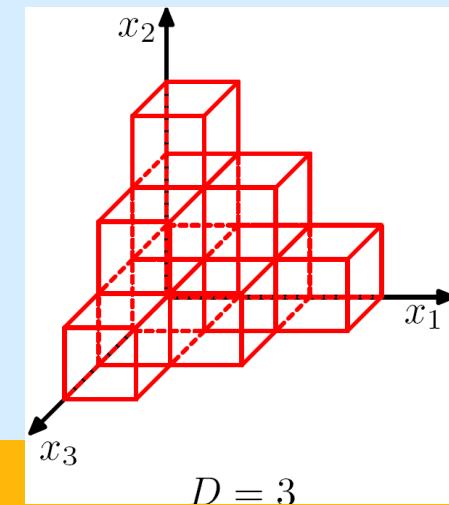
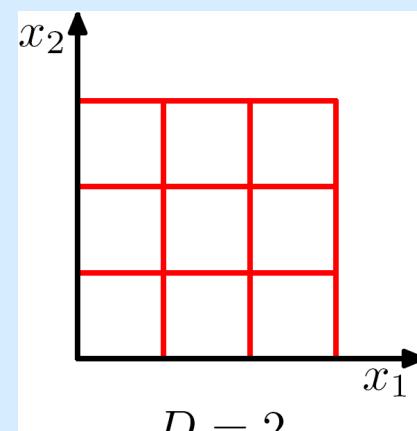
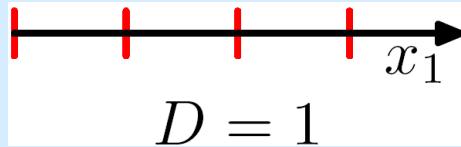
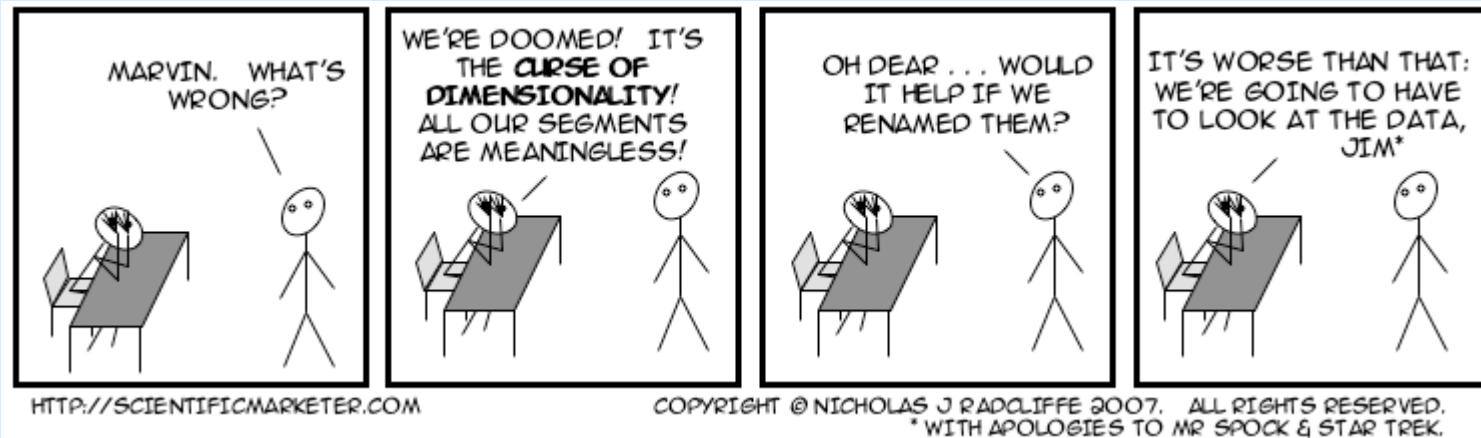


- “How to group customers for targeted marketing purposes?”
- “Which neighborhoods in a country are most similar to each other?”
- “What groups of insurance policy holders have high claim costs?”
- “How to group the products in a store based on their attributes?”
- “How to group pictures based on their description?”

# Dimensionality



- What is that?



# PCA



- Basically, reduce number of dimensions
- If you have 50 columns, are all of them needed?
  - Or can we squeeze info of those columns into 3 columns?
  - The three columns would be some combination of 50 columns
- This tremendously helps training of models
- Generally used as a step before actual regressions/classification



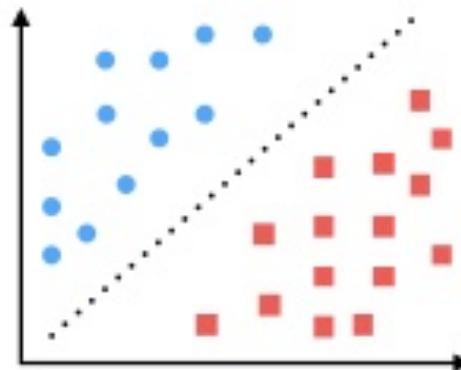
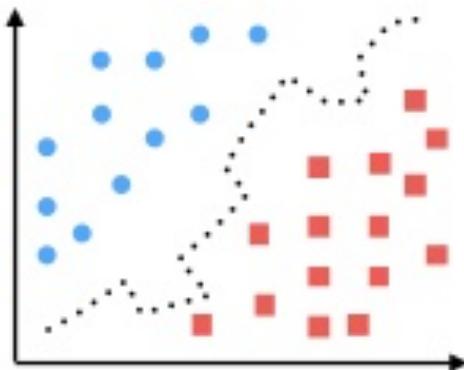
While **dimensionality reduction** is an important tool in machine learning/data mining, we must always be aware that it can distort the data in misleading ways.

Above is a two dimensional projection of an intrinsically three dimensional world....



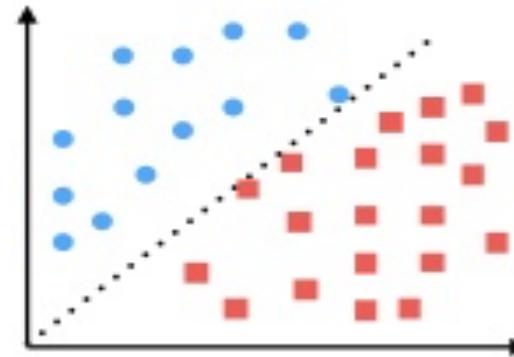
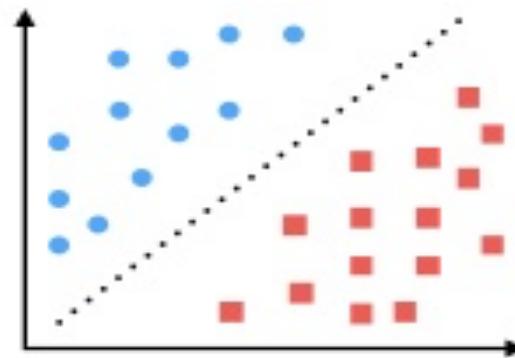
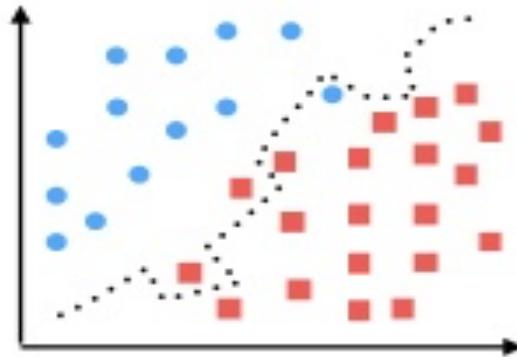
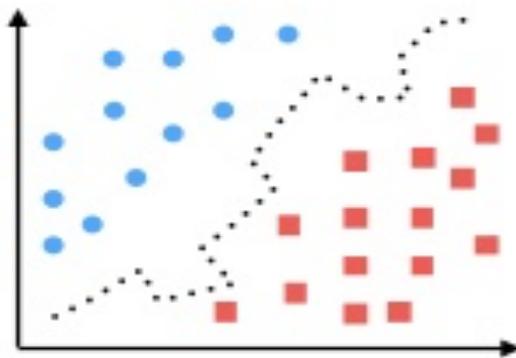
*Original photographer unknown*

# Problem of over-fitting



How well does the model perform on unseen data?

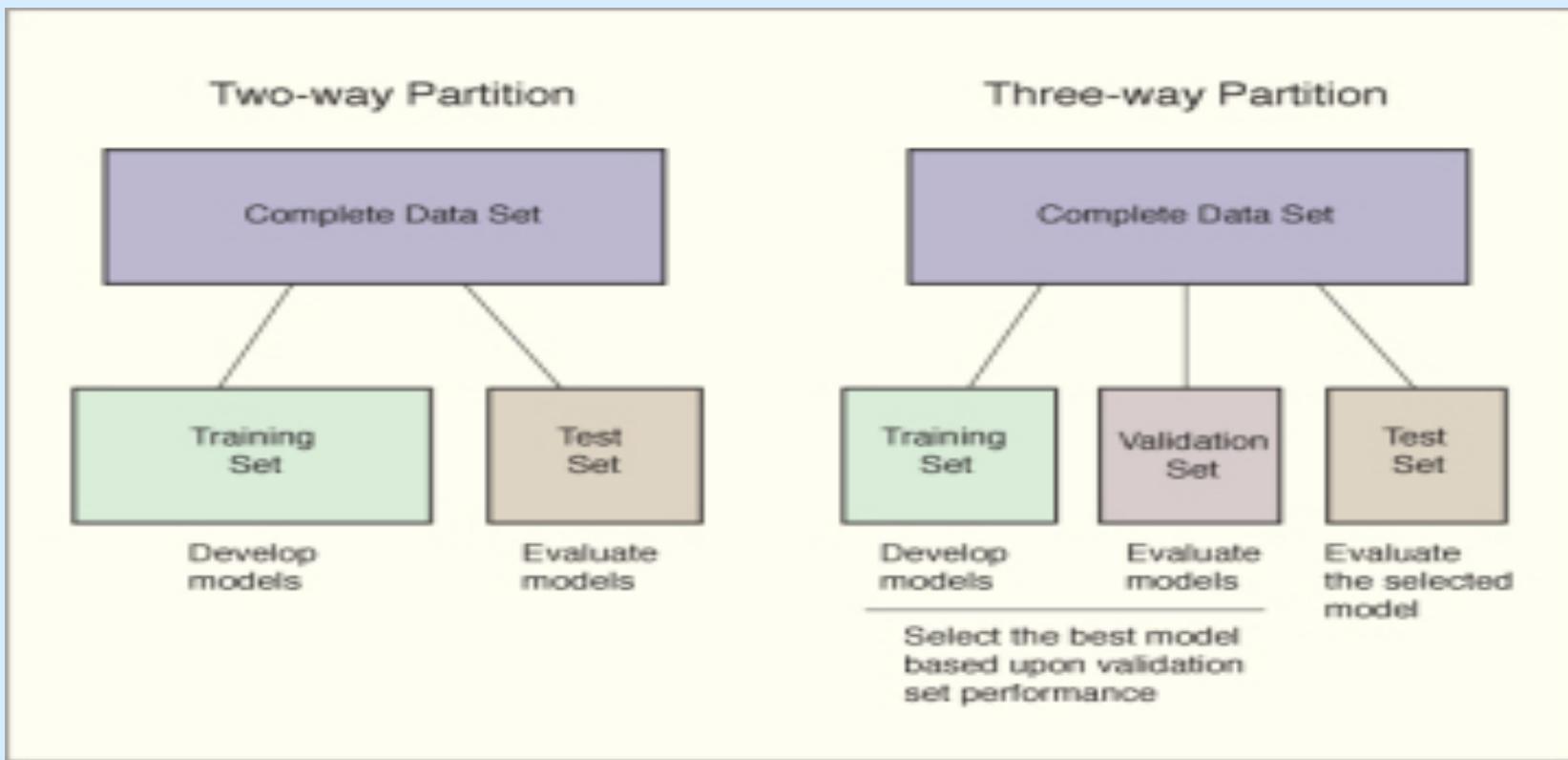
# Overfitting



# Evaluation of Accuracy



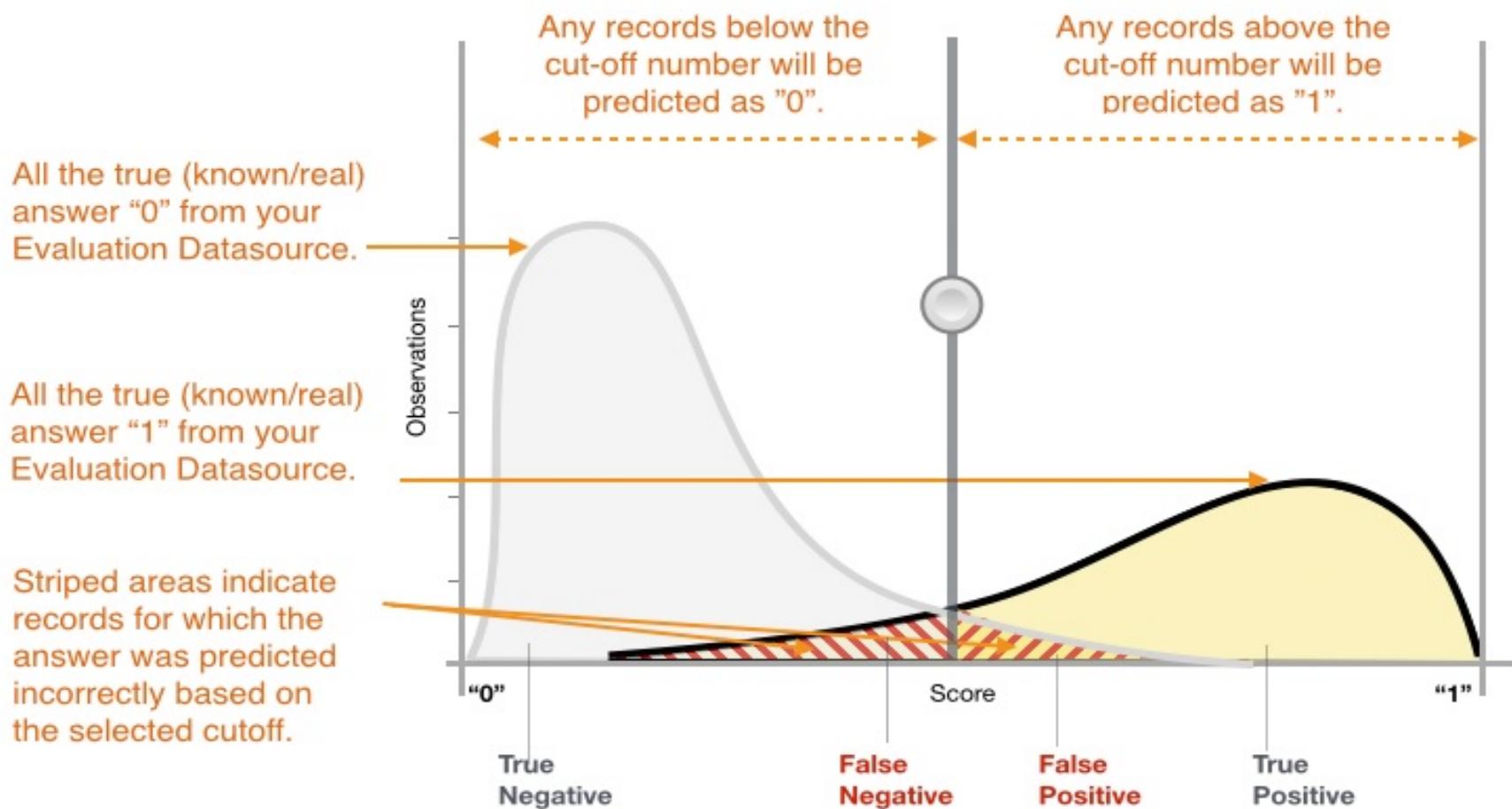
- Question: How to evaluate the predictive accuracy of your model?
- Answer: Partition the data set into Train and Test sets.
- Train is like the “past” you learn from, and Test is like the “future” you predict.



# Some ways of evaluating – Binary Classification



		Actual Value (as confirmed by experiment)	
		positives	negatives
Predicted Value (predicted by the test)	positives	<b>TP</b> True Positive	<b>FP</b> False Positive
	negatives	<b>FN</b> False Negative	<b>TN</b> True Negative



# Evaluating Classification Accuracy



- Simply put, accuracy = (no of correct classifications)/total observations
- There are other things also to consider
- Efficiency
  - Time to construct the model
  - Time to use the model
- Scalability
- Interpretable
  - Can you take the model to a non-technical business audience and convince them to deploy it in business?
- Often times, a succinct model is a good model

# Evaluating Continuous Predictions

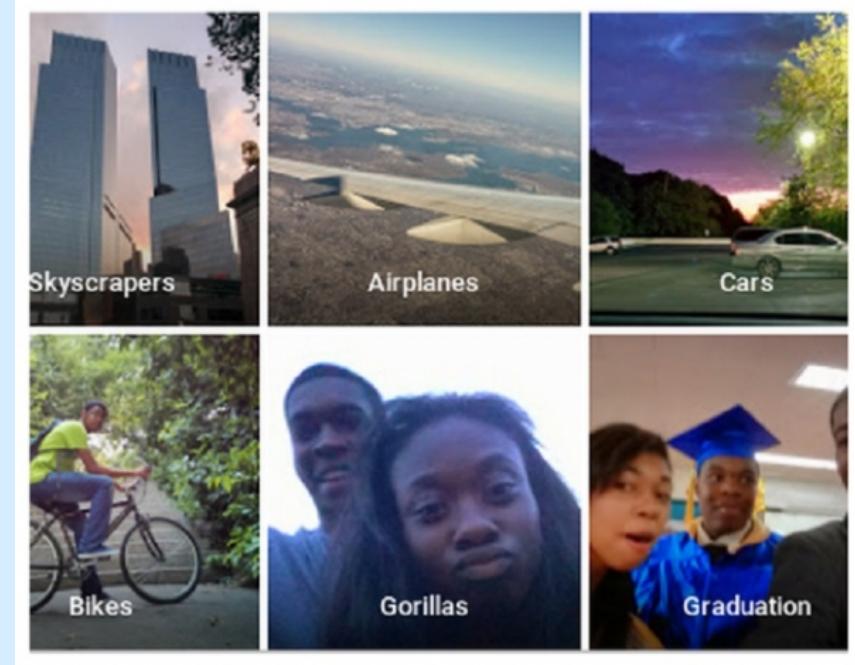


- One popular metric is Root Mean Squared Error (RMSE)
- Idea is very simple. We want to know quality of prediction
  - For each observation: actual-predicted and then square it
  - Sounds similar to deviation?
  - Well it actually is
  - Deviation is error from the actuals
  - That's what we care for

# Is ML all great?



- In May 2015 Flickr released an automatic image tagging capability that mistakenly labeled a black man for an ape.
- Soon afterwards, Google came up with a photo labeling tool similar to Flickr, which made similar mistakes. Black men were tagged as gorillas.
- A recent Carnegie Mellon University study showed that Google displayed ads in a way that discriminated based on the gender of the user.



# Fundamental Assumptions of Learning



- Training data distribution is similar to test data distribution
- This does not work out in real life
  - Data is representation of reality
  - As reality changes, data also changes
- If the violations are more severe, this would imply poor prediction results even if your model works very well on training data

# So whom to blame



## "No Free Lunch" :(

D. H. Wolpert. The supervised learning no-free-lunch theorems. In Soft Computing and Industry, pages 25–42. Springer, 2002.

Our model is a simplification of reality



Simplification is based on assumptions (model bias)



Assumptions fail in certain situations

Roughly speaking:

***"No one model works best for all possible situations."***

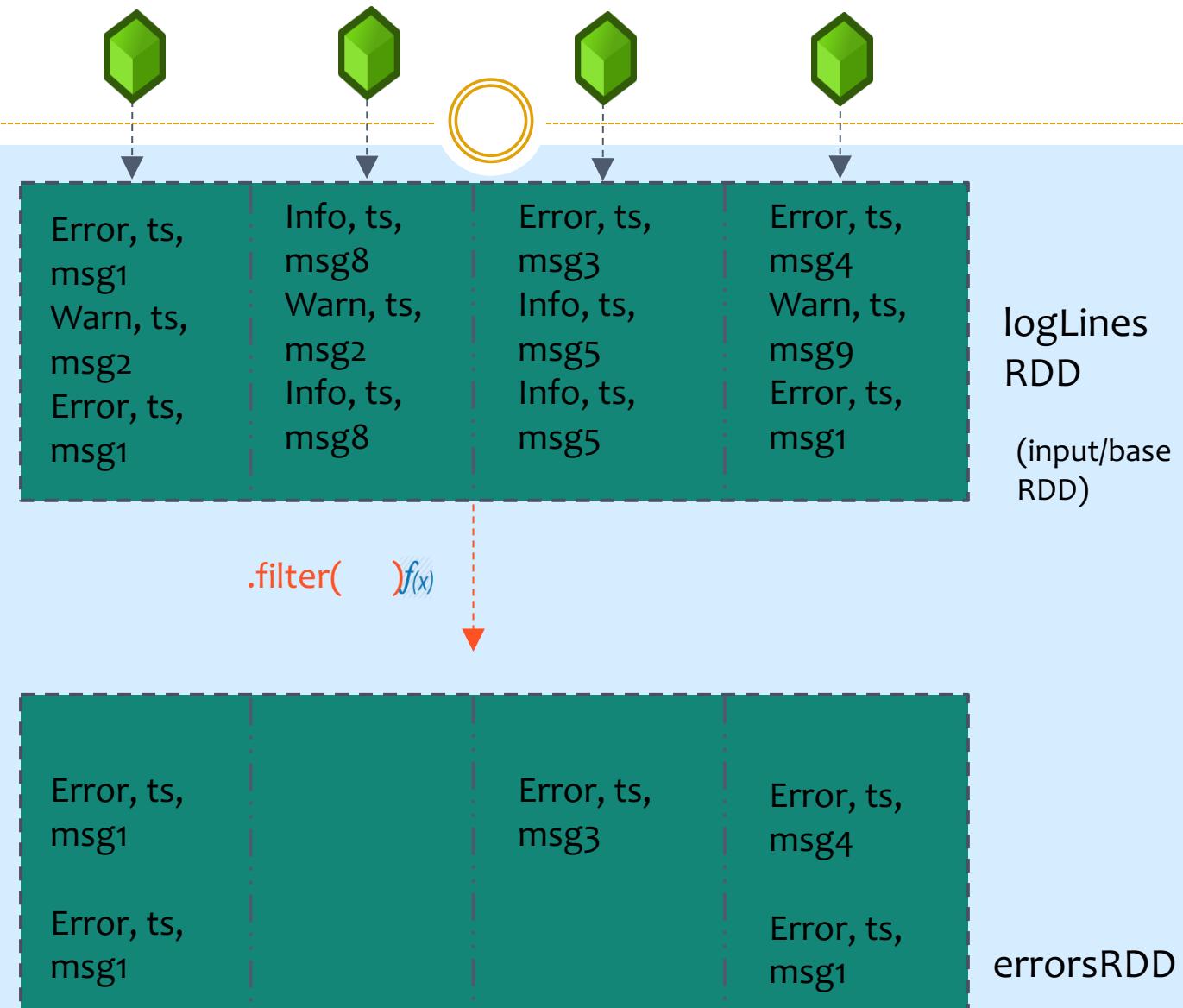
# Quick Review of Spark

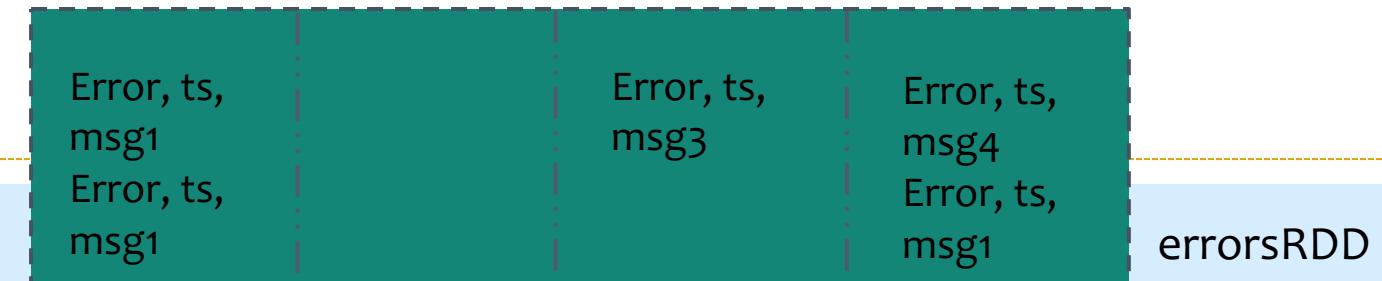


# Spark Core Components

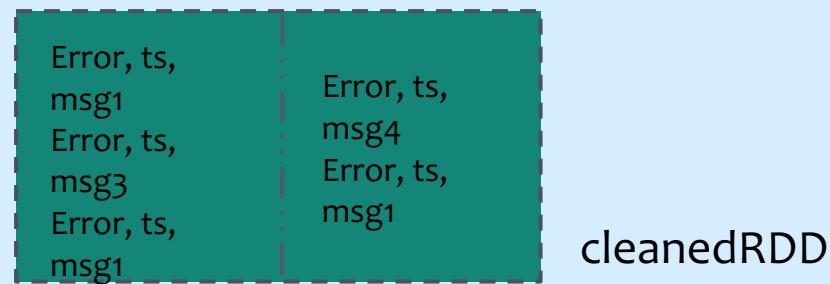


- Three building blocks
- **SparkContext**
  - Entry point to the Spark program
  - This is what you call from pyspark or from your own application
- **RDD**
  - Resilient Distributed Dataset
  - Container for your data
- Operations – Actions and Transformations





.coalesce( 2 )



.collect( )

```

ec2-user@ip-10-0-12-60:~$ spark
Welcome to
version 1.1.0

Using Scala version 2.10.4 (Java HotSpot(TM) 64-Bit Server VM, Java 1.7.0_71)
Type in expressions to have them evaluated.
Type help for more information.
Creating SparkContext...
SparkContext available as sc.
Spark context available as sc.
Type in expressions to have them evaluated.
Type help for more information.

scala> val keyvalueRDD = sc.cassandraTable("tinykeyspace", "keyvaluestable")
keyvalueRDD: com.datastax.spark.connector.rdd.CassandraRDD[com.datastax.spark.connector.CassandraRow] = CassandraRDD[0] at RDD at CassandraRDD.scala:49

scala> keyvalueRDD.count()
res1: Long = 4
scala>

```

Driver

# Technical Stuff Starts Now



# Machine Learning in Spark



- Two libraries
  - MLLib
  - Spark ML (Pipelines API)
- Scalable (Overcome the problem of data fitting one machine only)
- MLLib is more mature
- Commonly used ML algos and statistical analysis tools
- Machine Learning on large datasets
- Scalable, faster
- Both batch and streaming data

# MLLib & Spark ML



Spark SQL

Spark  
Streaming

MLlib

GraphX

Spark Core

# Statistical Utilities



- Summary stats
- Correlation
- Hypothesis Testing
- Sampling
- Random data

# Algorithms



- Regression
  - Linear and Logistic, SVM, Naive Bayes, Decision Tree, Random Forest
- Clustering
  - K-means, Gaussian, LDA, PCA, SVD
- Data Mining and Feature Extraction
  - TF-IDF, Word2Vec, Normalization, Item set mining, Association Rules, Collaborative Filtering
- APIs for Deep Learning and ANN

# Algorithms

- 
- ```
graph LR; A[• Alternating Least Squares  
• Lasso  
• Ridge Regression  
• Logistic Regression  
• Decision Trees  
• Naïve Bayes  
• Support Vector Machines  
• K-Means  
• Gradient descent  
• L-BFGS  
• Random data generation  
• Linear algebra  
• Feature transformations  
• Statistics: testing, correlation  
• Evaluation metrics] --> B[Collaborative Filtering for Recommendation]; A --> C[Prediction]; A --> D[Clustering]; A --> E[Optimization]; A --> F[Many Utilities]
```
- Alternating Least Squares
  - Lasso
  - Ridge Regression
  - Logistic Regression
  - Decision Trees
  - Naïve Bayes
  - Support Vector Machines
  - K-Means
  - Gradient descent
  - L-BFGS
  - Random data generation
  - Linear algebra
  - Feature transformations
  - Statistics: testing, correlation
  - Evaluation metrics

# API



- Data Types – MLLib operates on these data types
- Vector
- LabeledPoint
- Rating

# Vector



- Represents one observation in a dataset
- Represents element in n-dimensional space
- Used for representing features
- Vector of length n represents observation with n features
- Two types of vectors
  - DenseVector
  - SparseVector

# Data Type - Vector



- **Dense**
  - Stores a value at each position in vector
  - Use if your dataset does not have too many zeros
    - `import org.apache.spark.mllib.linalg._`
    - `val denseVector = Vectors.dense(1.0, 0.0, 3.0)`
- **Sparse**
  - Stores only non-zero values
    - `import org.apache.spark.mllib.linalg._`
    - `val sparseVector = Vectors.sparse(10, Array(3, 6), Array(100.0, 200.0))`

# LabeledPoint



- An observation in a labeled dataset
- Contains both label and features (Vector)
- Primary RDD abstraction
- You should transform dataset into RDD of LabeledPoints to do ML on Spark
- Can represent both categorical and numerical values

# LabeledPoint Example

- import org.apache.spark.mllib.linalg.Vectors
- import  
org.apache.spark.mllib.regression.LabeledPoint
- val positive = LabeledPoint(1.0, Vectors.dense(10.0,  
30.0, 20.0))
- val negative = LabeledPoint(0.0, Vectors.sparse(3,  
Array(0, 2), Array(200.0, 300.0)))

# Rating



- Used with recommendation algorithms
- Represents user rating
- Must transform dataset into an RDD of Ratings to run a recommendation algorithm
- Has three fields
  - User
  - Product
  - Rating
- import org.apache.spark.mllib.recommendation.\_
- val rating = Rating(100, 10, 3.0)

# Feature Extraction and Basic Statistics



- Several functions for basic operations
- Scaling, normalization, summary statistics, correlations, sampling etc
- Random Number Generators
- TF-IDF

# Algorithms and Models



- Model represented by a class
- ML algorithm also represented by a class
- Two key methods
  - Train
  - Predict

# Regression Algorithm



- Many different regression methods available
  - LinearRegression, IsotonicRegression, RidgeRegression
- Train method fits a linear regression model to data
  - Returns an object of type LinearRegressionModel
- Takes RDD of LabeledPoints and returns regression model



- Let us move to hands-on