

OVERVIEW

MACHINE

LEARNING

INTRODUCTION TO DATA SCIENCE

TIM KRASKA



TOKENIZATION AND STEMMING

WORKING WITH TEXT

TOKENIZATION

Input: “*Friends, Romans and Countrymen*”

Output: Tokens

- *Friends*
- *Romans*
- *and*
- *Countrymen*

A **token** is an instance of a sequence of characters

COMMON STEPS

- **Remove Stop Words** (a, an, the, to, be, ...)
- **Normalization to terms**
 - **deleting periods:** U.S.A. → USA
 - **deleting hyphens:** *anti-discriminatory* → *antidiscriminatory*
 - **Abbreviations:** Massachusetts Institute of Technology → MIT
 - **Case-folding:** Meal → meal, Brown → brown
 - **Language-issues:** *Tuebingen, Tübingen* → *Tubingen*
 - **asymmetric expansion:** *windows* → *window*
 - ...
 - *Why is this a form of entity resolution? What examples above are problematic?*
- **Thesauri and soundex**
 - *car = automobile* *color = colour*
- **Stemming**

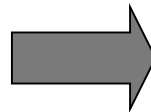
STEMMING

Reduce terms to their “roots” before indexing

“Stemming” suggest crude affix chopping

- language dependent
- e.g., *automate(s)*, *automatic*, *automation* all reduced to *automat*.

for example compressed
and compression are both
accepted as equivalent to
compress.



for exampl compress and
compress ar both accept
as equival to compress

PORTER'S ALGORITHM

Commonest algorithm for stemming English

- Results suggest it's at least as good as other stemming options

Conventions + 5 phases of reductions

- phases applied sequentially
- each phase consists of a set of commands
- sample convention: *Of the rules in a compound command, select the one that applies to the longest suffix.*

TYPICAL RULES IN PORTER

sses* → *ss

ies* → *i

ational* → *ate

tional* → *tion

Weight of word sensitive rules

***(m>1) EMENT* →**

- *replacement* → *replac*
- *cement* → *cement*

OTHER STEMMERS

Other stemmers exist, e.g., Lovins stemmer

- <http://www.comp.lancs.ac.uk/computing/research/stemming/general/lovins.htm>
- Single-pass, longest suffix removal (about 250 rules)

Full morphological analysis – at most modest benefits for retrieval

Do stemming and other normalizations help?

- English: very mixed results. Helps recall for some queries but harms precision on others
 - E.g., operative (dentistry) \Rightarrow oper
- Definitely useful for Spanish, German, Finnish, ...
 - 30% performance gains for Finnish!

MACHINE LEARNING

INTRODUCTION TO DATA SCIENCE

TIM KRASKA

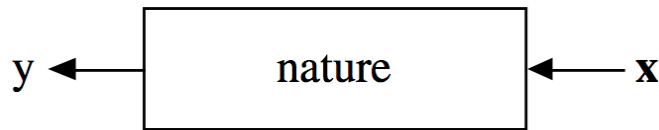


TITANIC DATASET

Label	Features							
survived	pclass	sex	age	sibsp	parch	fare	cabin	embarked
0	3	male	22	1	0	7.25		S
1	1	female	38	1	0	71.2833	C85	C
1	3	female	26	0	0	7.925		S
1	1	female	35	1	0	53.1	C123	S
0	3	male	35	0	0	8.05		S
0	3	male		0	0	8.4583		Q
0	1	male	54	0	0	51.8625	E46	S
0	3	male	2	3	1	21.075		S
1	3	female	27	0	2	11.1333		S
1	2	female	14	1	0	30.0708		C
1	3	female	4	1	1	16.7	G6	S
1	1	female	58	0	0	26.55	C103	S
0	3	male	20	0	0	8.05		S

DIFFERENCE BETWEEN STATISTICS AND MACHINE LEARNING

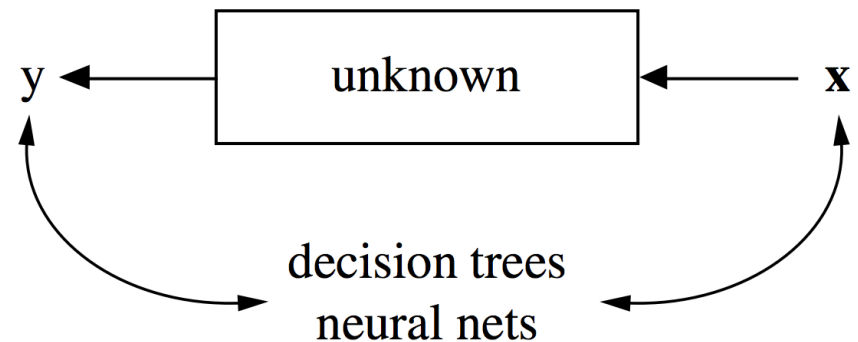
One view:



Emphasis on stochastic models of nature:

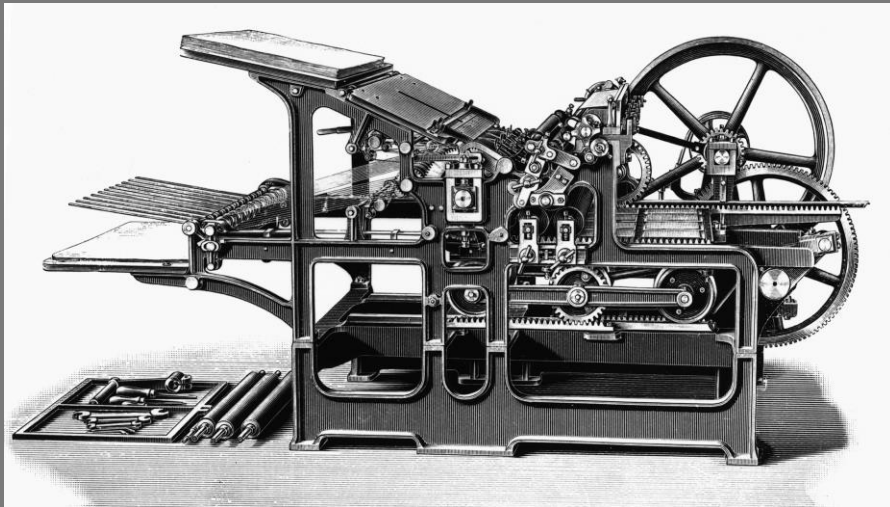


Find a function that predicts y from x :
no model of nature implied or needed





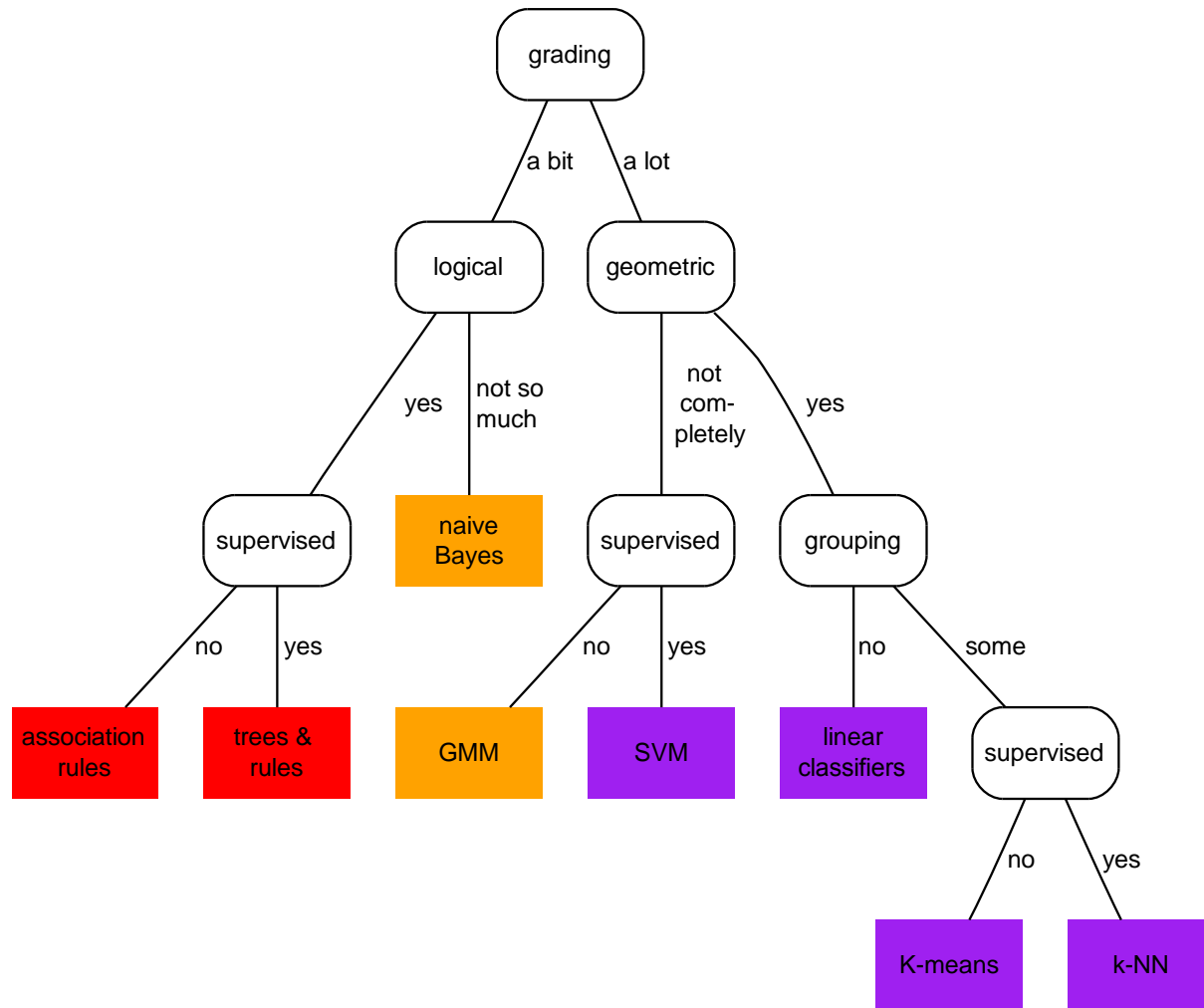
VS



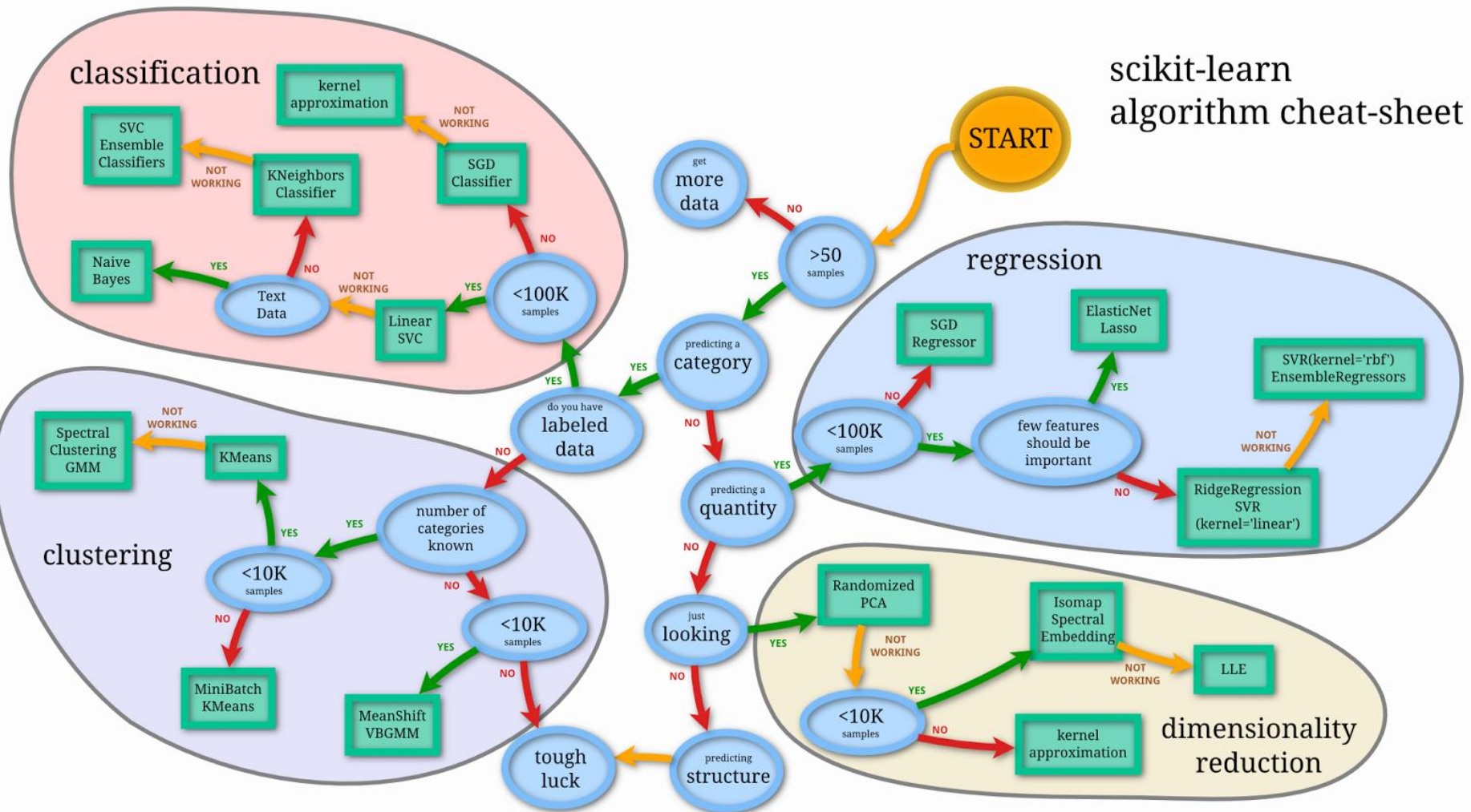
MACHINE LEARNING PROBLEMS

	Supervised Learning	Unsupervised Learning
Discrete	classification or categorization	clustering
Continuous	regression	dimensionality reduction

ML on ML



(BE CAREFULLY)

scikit-learn
algorithm cheat-sheet

MACHINE LEARNING PROBLEMS

	Supervised Learning	Unsupervised Learning
Discrete	classification or categorization	clustering
Continuous	regression	dimensionality reduction

CLUSTERING STRATEGIES

K-means

- Iteratively re-assign points to the nearest cluster center

Agglomerative clustering

- Start with each point as its own cluster and iteratively merge the closest clusters

Mean-shift clustering

- Estimate modes of PDF (i.e., the value x at which its probability mass function takes its maximum value)

Spectral clustering

- Split the nodes in a graph based on assigned links with similarity weights

As we go down this chart, the clustering strategies have more tendency to transitively group points even if they are not nearby in feature space

K-MEANS

Clusters based on **centroids** (aka the **center of gravity** or mean) of points in a cluster, c :

$$\vec{\mu}(c) = \frac{1}{|c|} \sum_{\vec{x} \in c} \vec{x}$$

Reassignment of instances to clusters is based on distance to the current cluster centroids.

- (Or one can equivalently phrase it in terms of similarities)

K-MEANS ALGORITHM

```
Select  $K$  random data points  $\{s_1, s_2, \dots, s_K\}$  as centroids  $c_j$ .  
Until clustering converges or other stopping criterion {  
  For each data point  $x_i$ :  
    Assign  $x_i$  to the closes centroid such that  
       $dist(x_i, c_j)$  is minimal.  
  For each cluster  $c_j$ , update the centroids  
     $c_j = \mu(c_j)$   
}
```

MORE FORMALLY

Given a set of observations (x_1, x_2, \dots, x_n) , where each observation is a d -dimensional real vector, k -means clustering aims to partition the n observations into k sets ($k \leq n$) $S = \{S_1, S_2, \dots, S_k\}$:

$$\arg \min_S \sum_{i=1 \dots k} \sum_{x_j \in S_i} \text{dist}(x_j, c_i)$$

$$\text{dist} = \left(x_j - c_i \right)^2$$

$$c_i = m(S_i) = \frac{1}{|S_i|} \sum_{x_i \in S_i} x_i$$

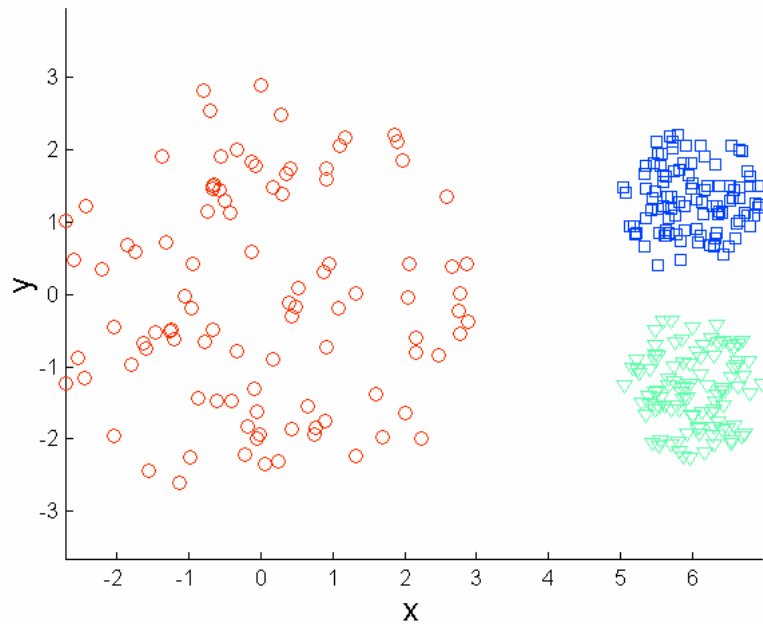


TERMINATION CONDITIONS

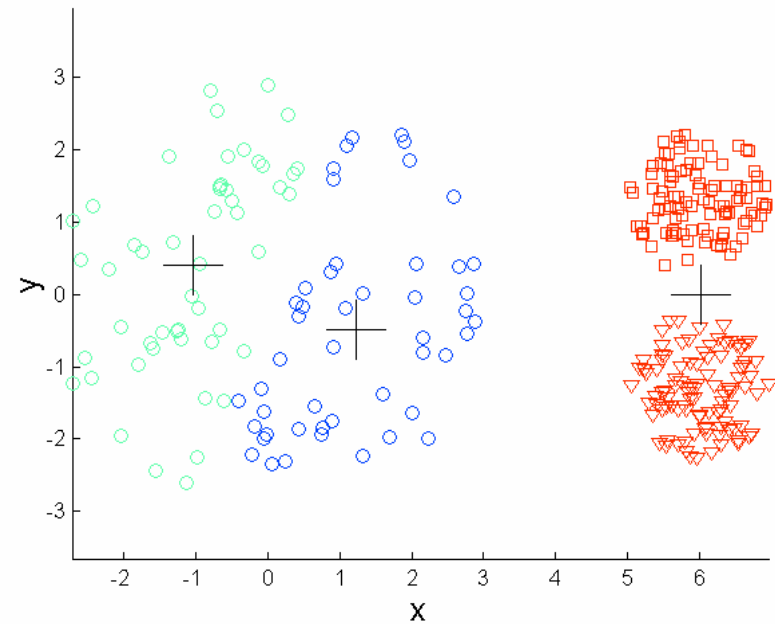
Several possibilities, e.g.,

- A fixed number of iterations.
- Partition unchanged.
- Centroid positions don't change.

LIMITATIONS OF K-MEANS: DIFFERING DENSITY



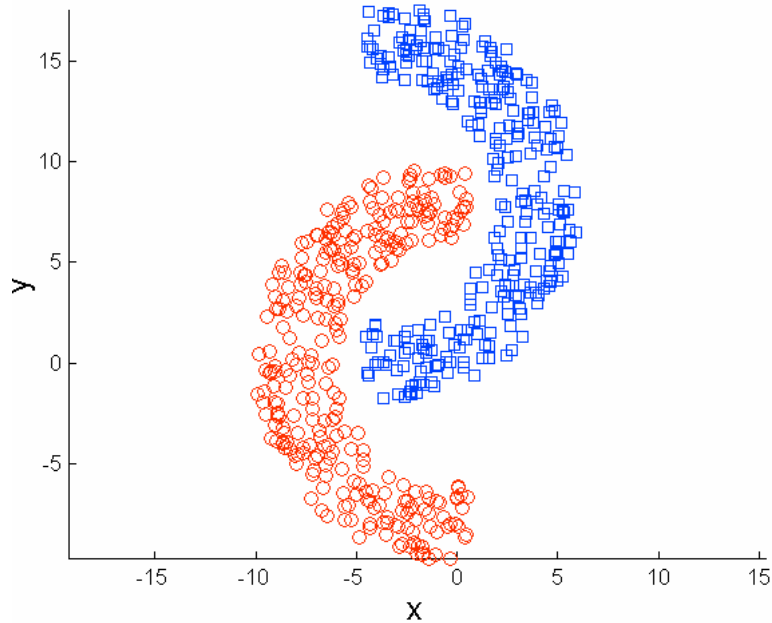
Original Points



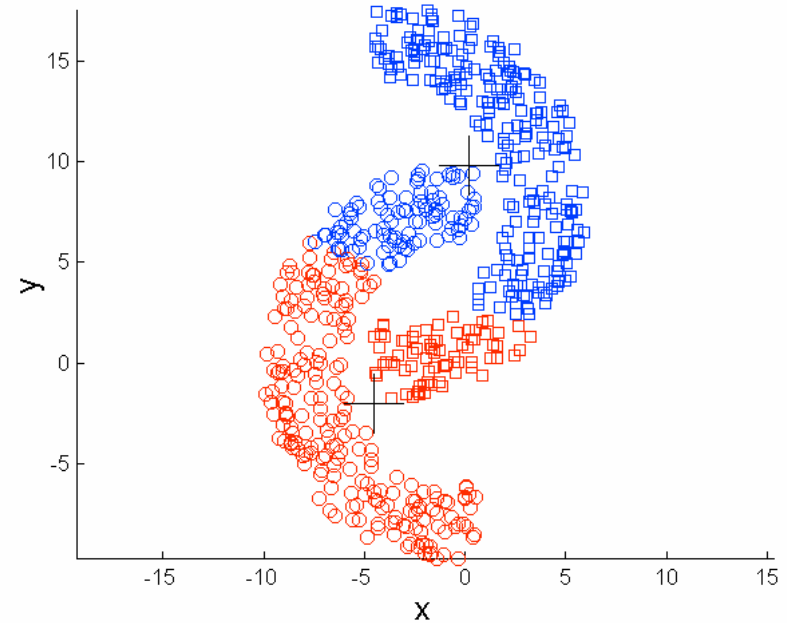
K-means (3 Clusters)

LIMITATIONS OF K-MEANS

Non-globular Shapes

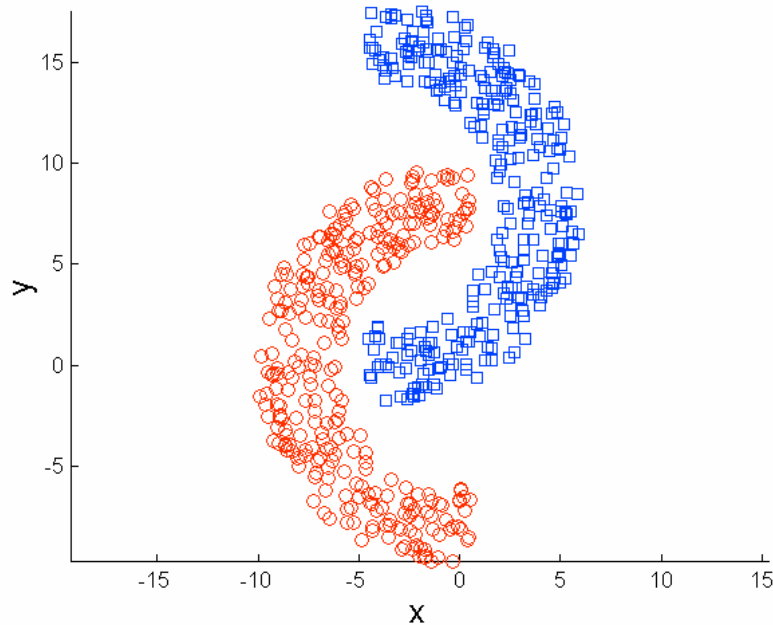


Original Points

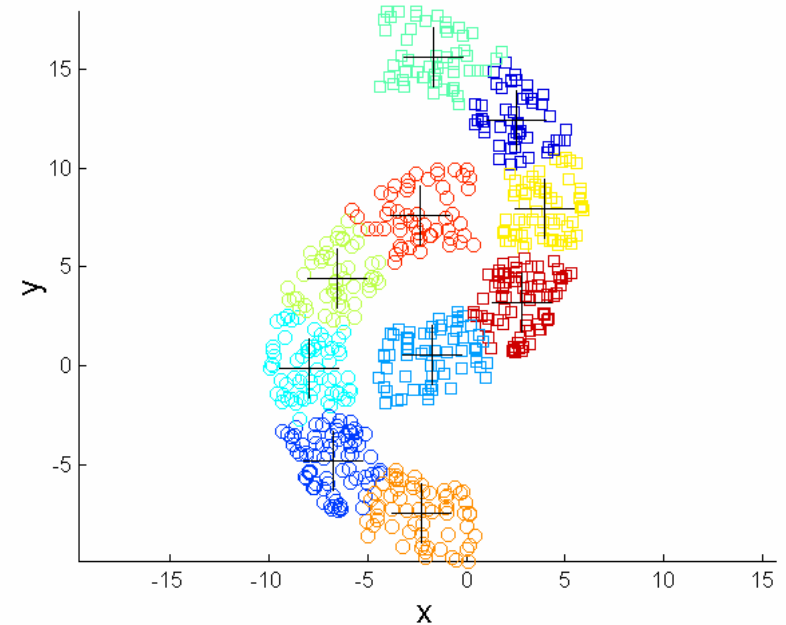


K-means (2 Clusters)

OVERCOMING K-MEANS LIMITATIONS



Original Points



K-means Clusters

Can you think of other ways to overcome the limitations?

CLUSTERING STRATEGIES

K-means

- Iteratively re-assign points to the nearest cluster center

Agglomerative clustering

- Start with each point as its own cluster and iteratively merge the closest clusters

Mean-shift clustering

- Estimate modes of PDF (i.e., the value x at which its probability mass function takes its maximum value)

Spectral clustering

- Split the nodes in a graph based on assigned links with similarity weights

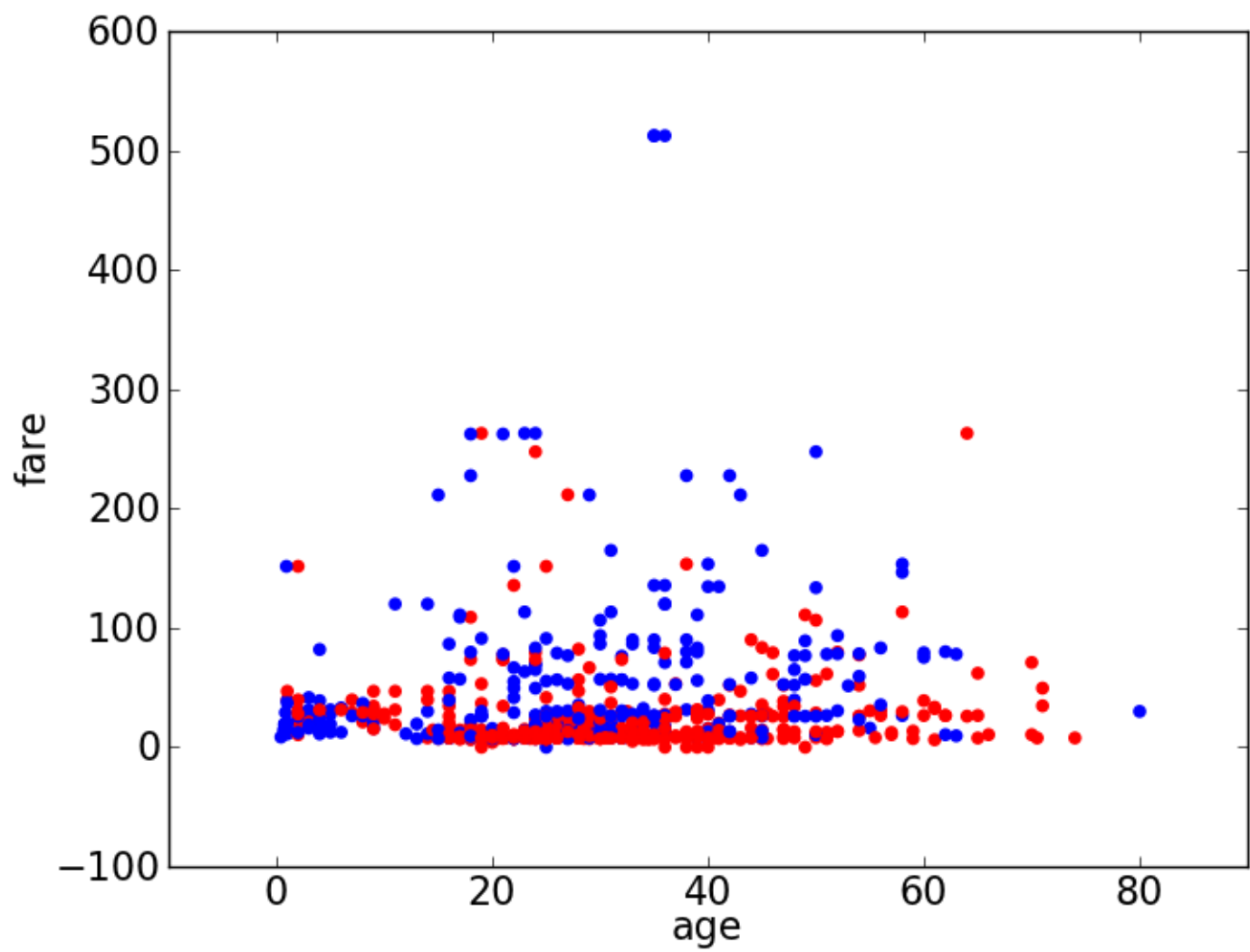
As we go down this chart, the clustering strategies have more tendency to transitively group points even if they are not nearby in feature space

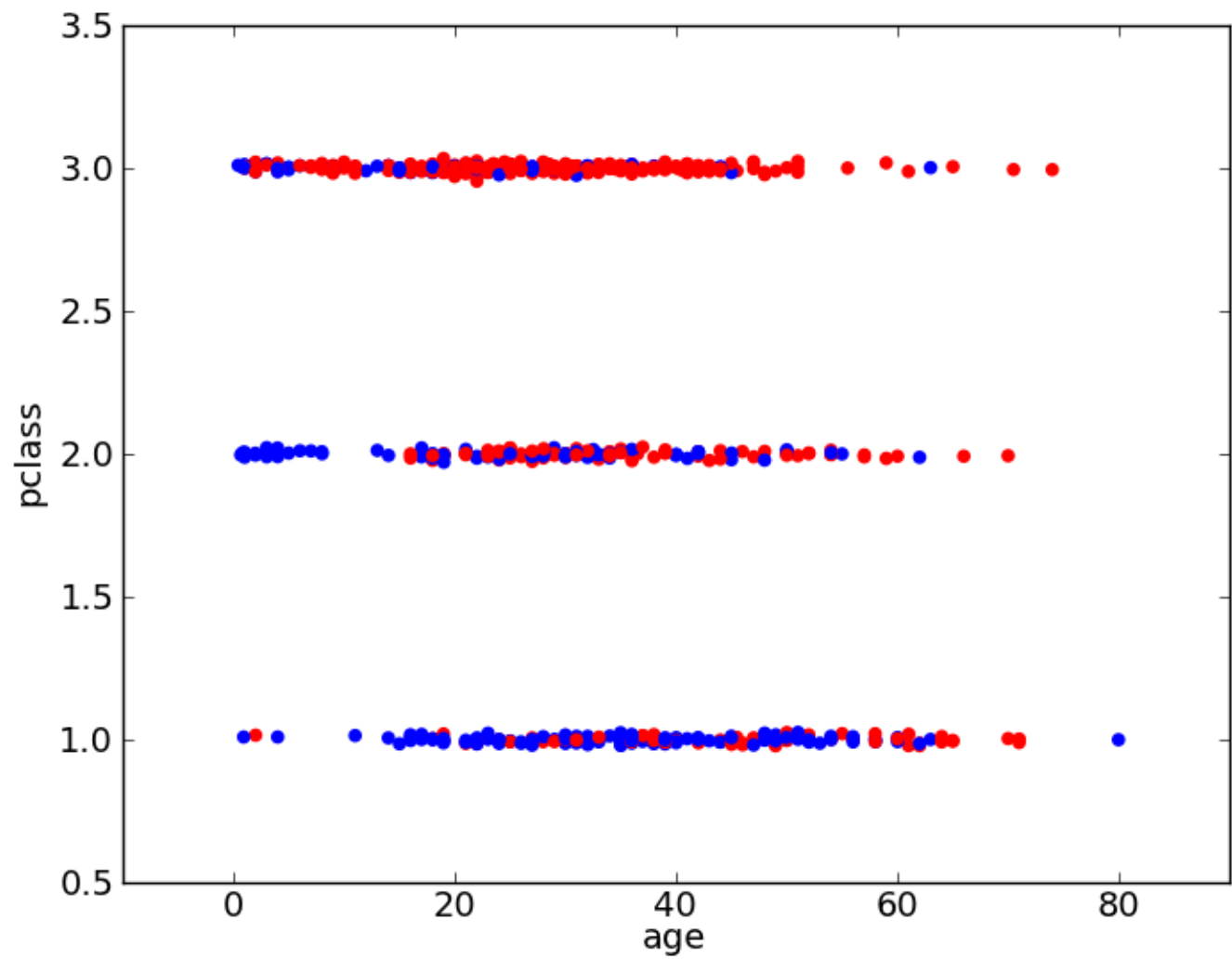
MACHINE LEARNING PROBLEMS

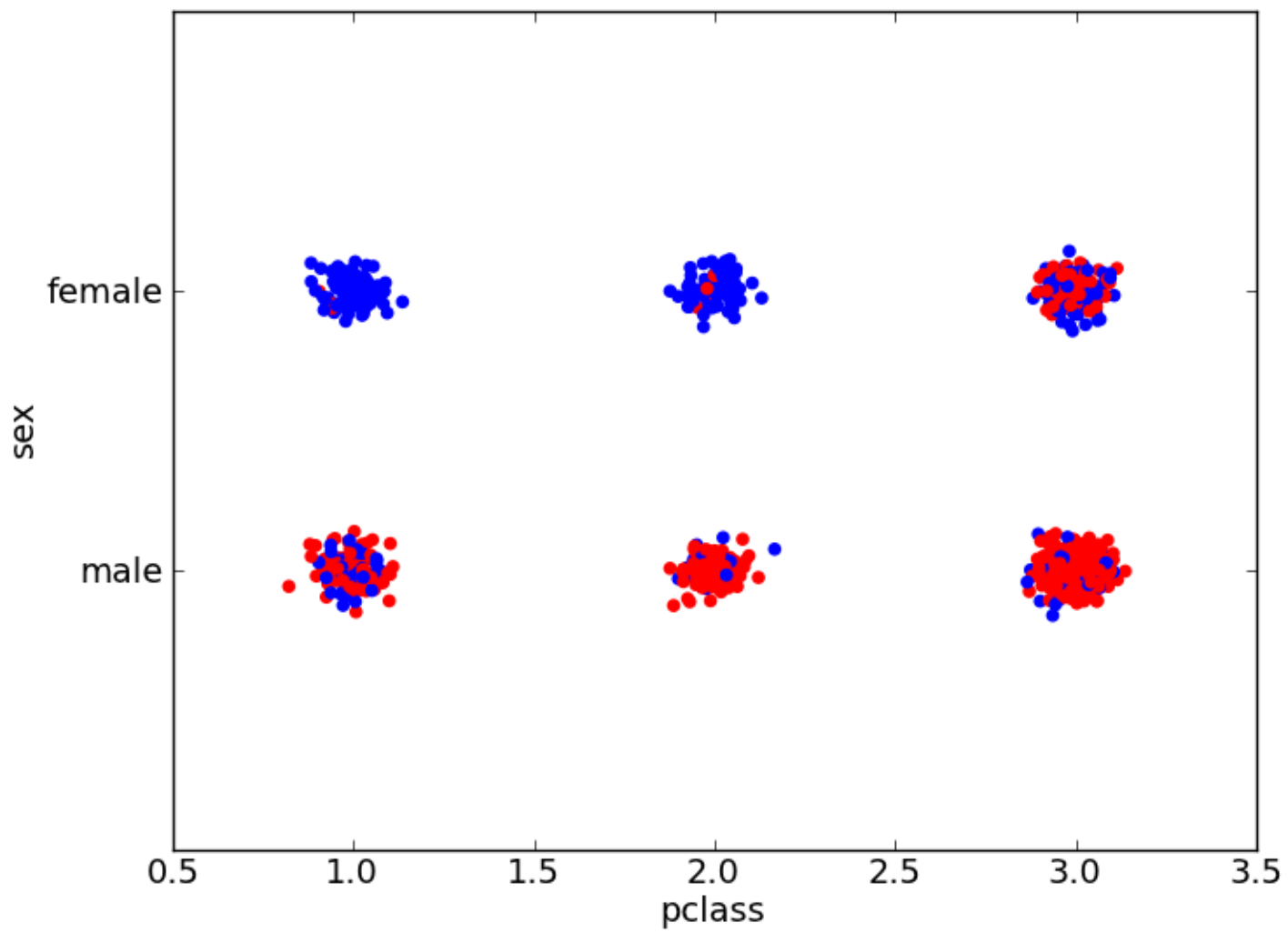
	Supervised Learning	Unsupervised Learning
Discrete	classification or categorization	clustering
Continuous	regression	dimensionality reduction

TITANIC DATASET

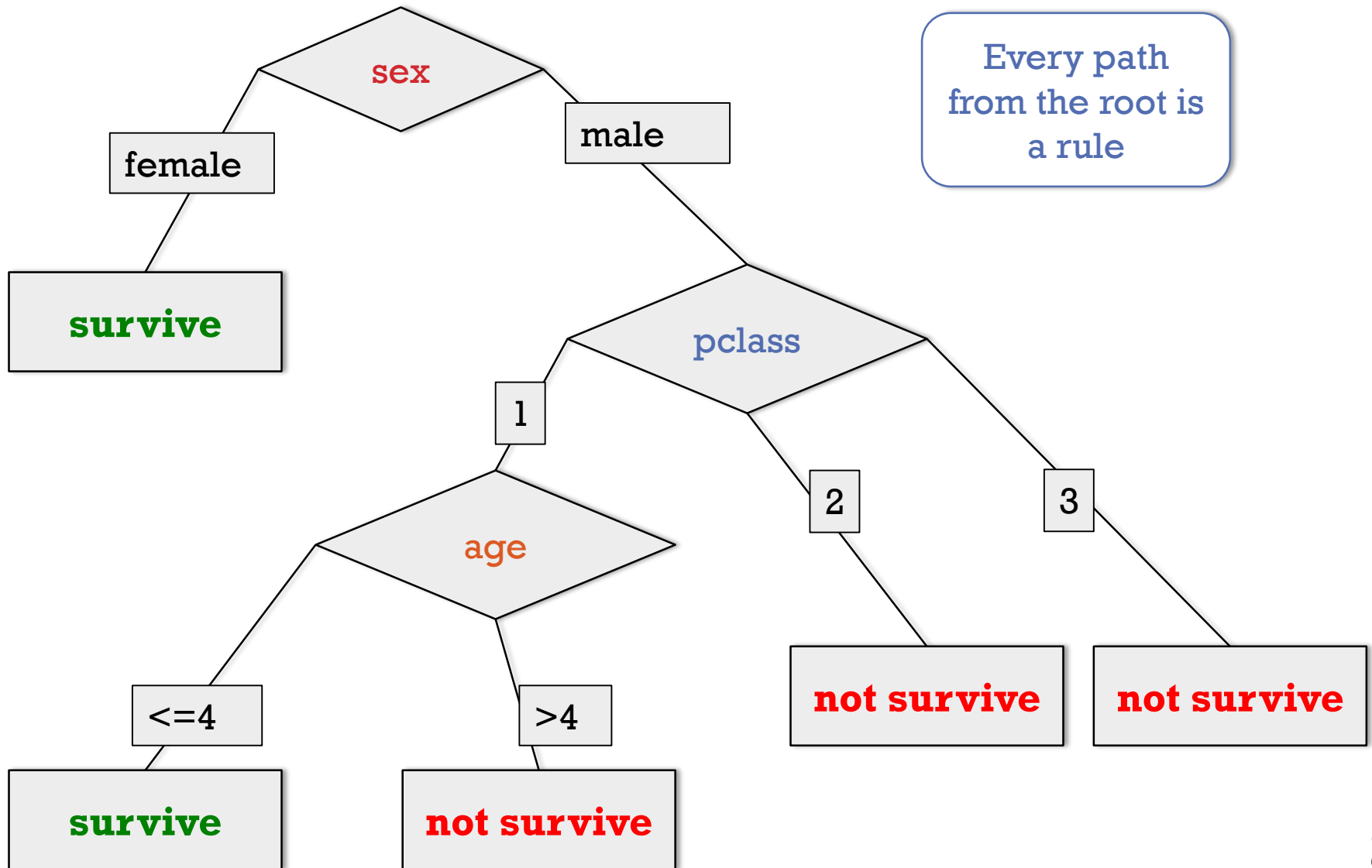
survived	pclass	sex	age	sibsp	parch	fare	cabin	embarked
0	3	male	22	1	0	7.25		S
1	1	female	38	1	0	71.2833	C85	C
1	3	female	26	0	0	7.925		S
1	1	female	35	1	0	53.1	C123	S
0	3	male	35	0	0	8.05		S
0	3	male		0	0	8.4583		Q
0	1	male	54	0	0	51.8625	E46	S
0	3	male	2	3	1	21.075		S
1	3	female	27	0	2	11.1333		S
1	2	female	14	1	0	30.0708		C
1	3	female	4	1	1	16.7	G6	S
1	1	female	58	0	0	26.55	C103	S
0	3	male	20	0	0	8.05		S







DECISION TREE



WHAT IS A CLASSIFIER

Apply a prediction function to a feature representation of an image/data-set to get the desired output:

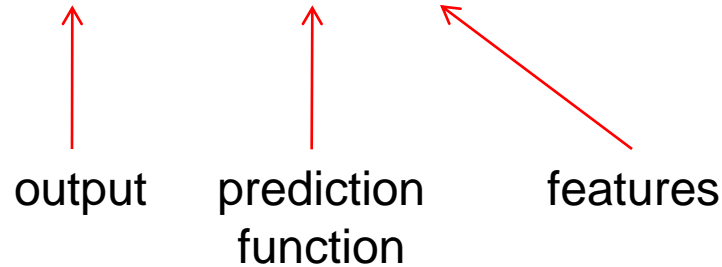
$$f(\text{apple image}) = \text{"apple"}$$

$$f(\text{tomato image}) = \text{"tomato"}$$

$$f(\text{cow image}) = \text{"cow"}$$

THE MACHINE LEARNING FRAMEWORK

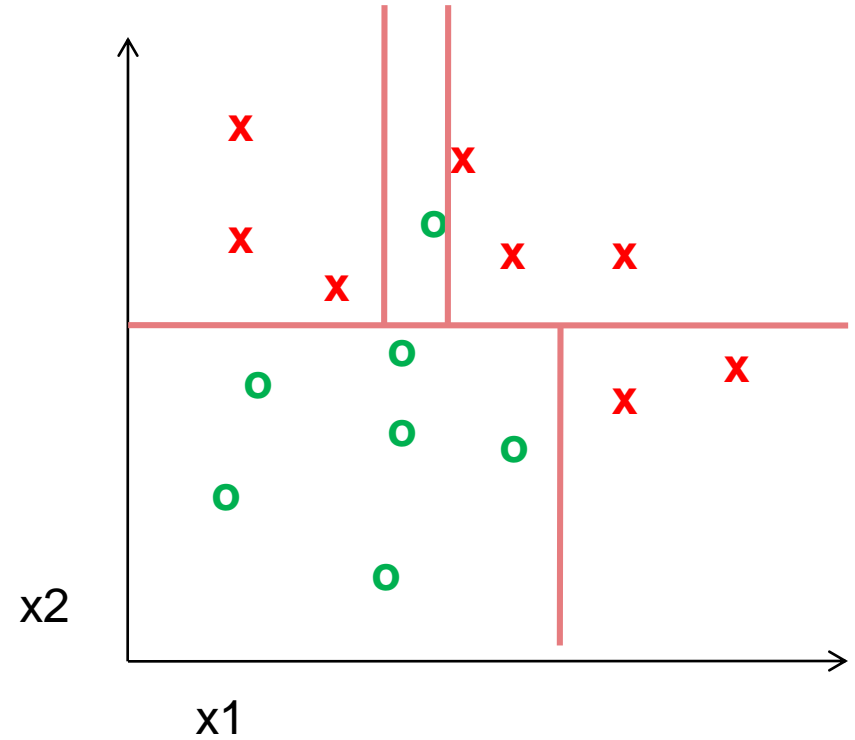
$$\mathbf{y} = \mathbf{f}(\mathbf{x})$$



Training: given a *training set* of labeled examples $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$, estimate the prediction function \mathbf{f} by minimizing the prediction error on the training set

Testing: apply \mathbf{f} to a never before seen *test example* \mathbf{x} and output the predicted value $\mathbf{y} = \mathbf{f}(\mathbf{x})$

DECISION BOUNDARIES: DECISION TREES



CLASSIFIER OVERVIEW

Representation	Evaluation	Optimization
Instances	Accuracy/Error rate	Combinatorial optimization
<i>K</i> -nearest neighbor	Precision and recall	Greedy search
Support vector machines	Squared error	Beam search
Hyperplanes	Likelihood	Branch-and-bound
Naive Bayes	Posterior probability	Continuous optimization
Logistic regression	Information gain	Unconstrained
Decision trees	K-L divergence	Gradient descent
Sets of rules	Cost/Utility	Conjugate gradient
Propositional rules	Margin	Quasi-Newton methods
Logic programs		Constrained
Neural networks		Linear programming
Graphical models		Quadratic programming
Bayesian networks		
Conditional random fields		

RECAP: DECISION TREES

Representation

- A set of rules: IF...THEN conditions

Evaluation

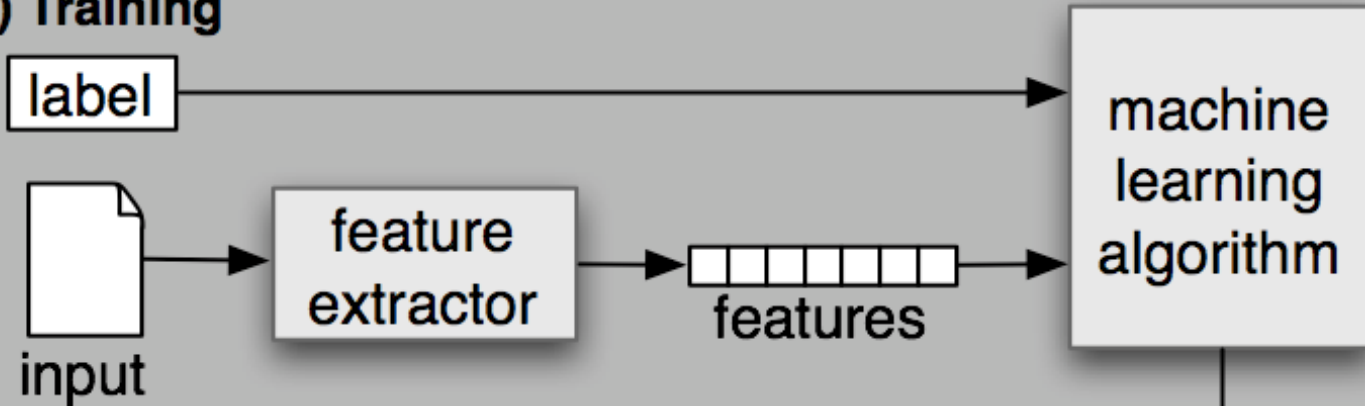
- coverage: # of data points that satisfy conditions
- accuracy = # of correct predictions / coverage

Optimization

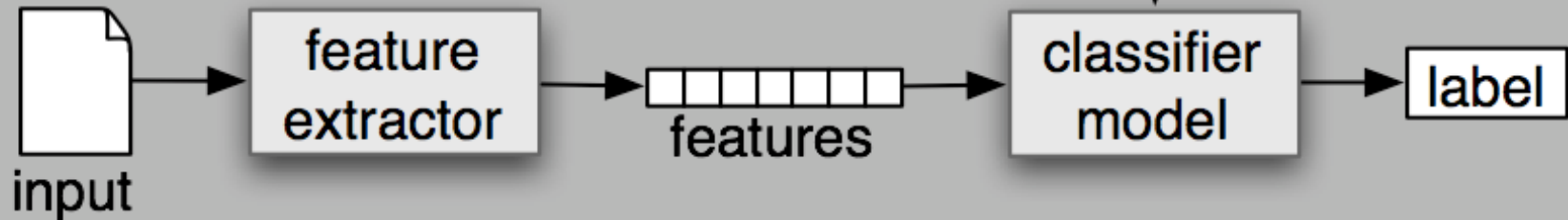
- Build decision tree that maximize accuracy

ML PIPELINE (SUPERVISED)

(a) Training



(b) Prediction



FEATURES

Fact Table
- <u>Shop ID</u>
- <u>Customer ID</u>
- <u>Date ID</u>
- <u>Product ID</u>
- Amount
- Volume
- Profit
- ...

Fact Table
- <u>Shop ID</u>
- <u>Customer ID</u>
- <u>Date ID</u>
- <u>Product ID</u>
- Amount
- Volume
- Profit
- Delivery Time
- ...

Product
- <u>Product ID</u>
- Type_ID
- Brand_ID
- Length
- Height
- Depth
- Weight
- ...

Product_Type
- <u>Type ID</u>
- Name
- Description
- ...

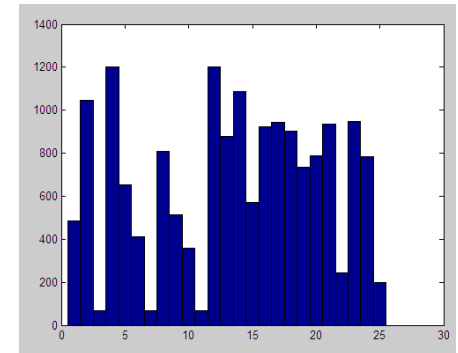
Brand
- <u>Brand ID</u>
- Name
- ...

Customer State	Product Type	Product Weight	Volume (L*H*D)	Month	Delivery Time

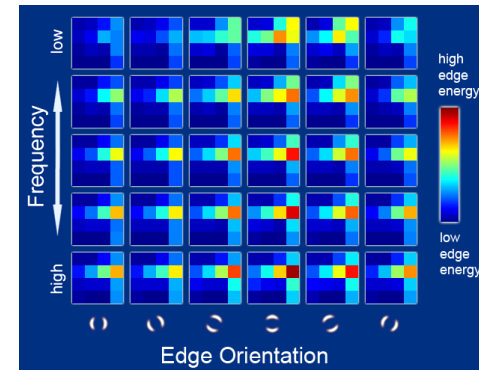
IMAGE FEATURES

Raw pixels

Histograms

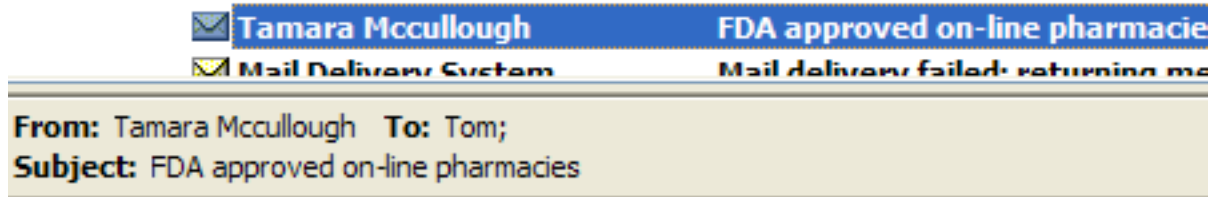


GIST descriptors



...

TEXT FEATURES

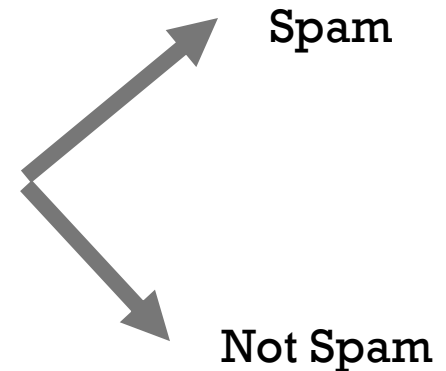


FDA approved on-line pharmacies.
Chose your product and site below:

Canadian pharmacy - Cialis Soft Tabs - \$5.78, **Viagra Professional** - \$1.38, Human Growth Hormone - \$43.37, Meridia - \$3.32, Tramadol

HerbalKing - Herbal pills for **Hair enlargement**. Techniques, products, dangerous pumps, exercises and surgeries.

Anatrim - Are you ready for Summer? Use **Anatrim**, the most powerful



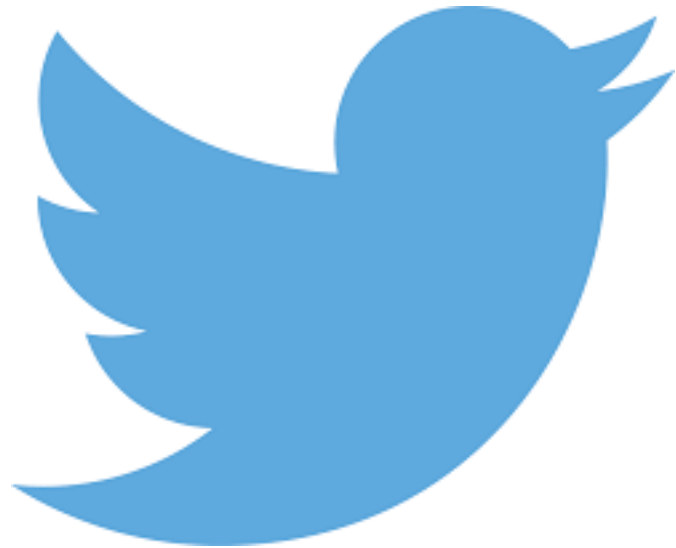
Bag of Words

$\begin{pmatrix} \textit{Viagra} \\ \textit{Soft} \\ \textit{Herbel} \\ \textit{Pills} \\ \textit{Are} \\ \dots \end{pmatrix}$

N-Grams

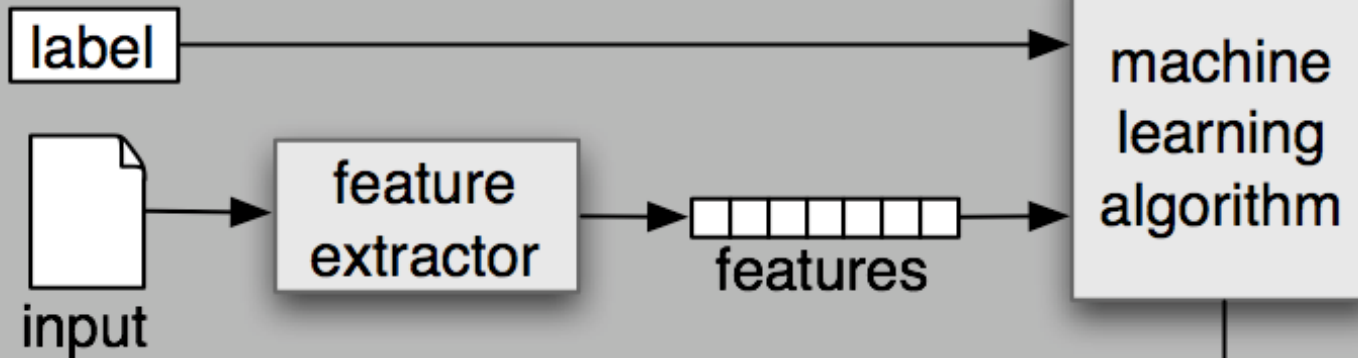
$\begin{pmatrix} \textit{herbel pills} \\ \textit{pills for} \\ \textit{for Hair} \\ \textit{Hair enlargement} \\ \textit{enlargement Techniques} \\ \dots \end{pmatrix}$

FEATURE TO PREDICT UNEMPLOYMENT

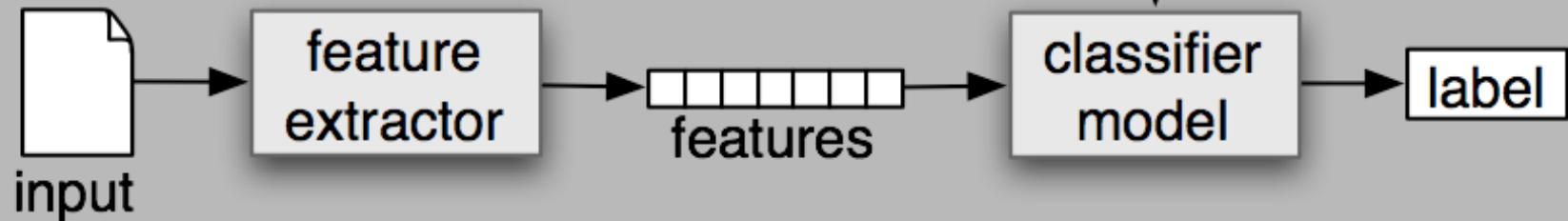


ML PIPELINE (SUPERVISED)

(a) Training



(b) Prediction



MANY CLASSIFIERS TO CHOOSE FROM

Decision Trees

K-nearest neighbor

Support Vector Machines

Logistic Regression

Naïve Bayes

Random Forrest

Bayesian network

Randomized Forests

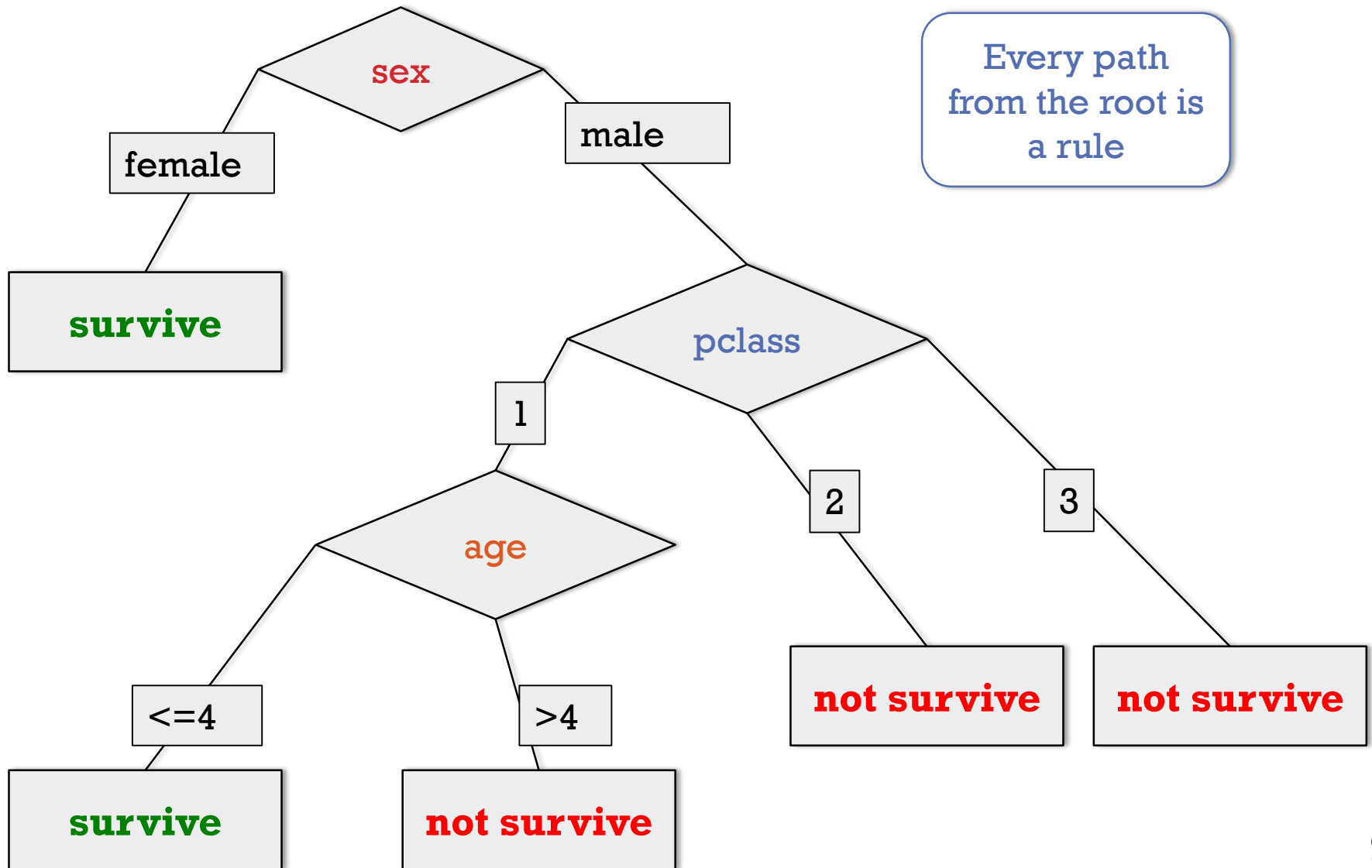
Boosted Decision Trees

RBMs

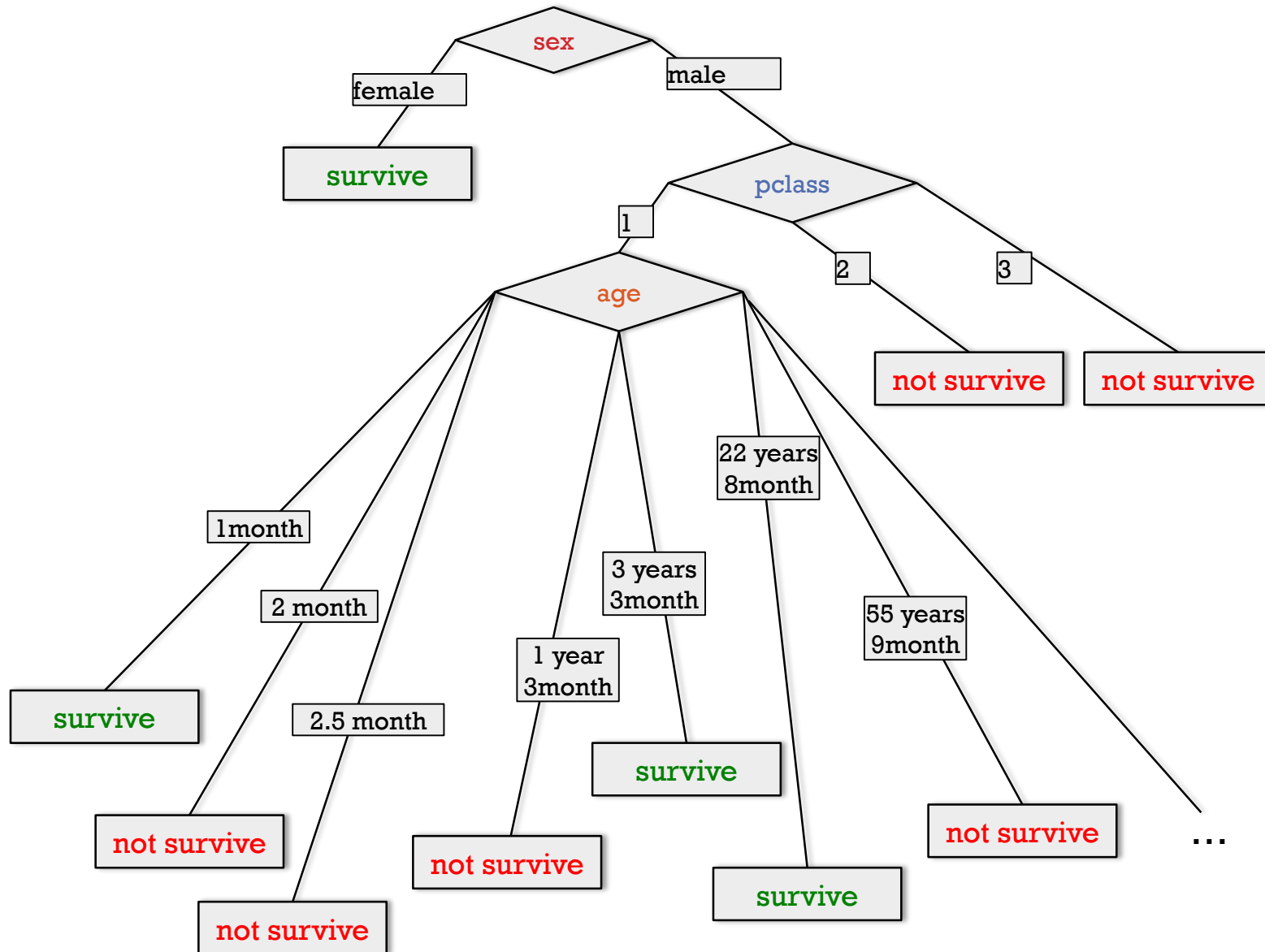
....

Which is the best one?

DECISION TREE



WHAT ABOUT THIS TREE?



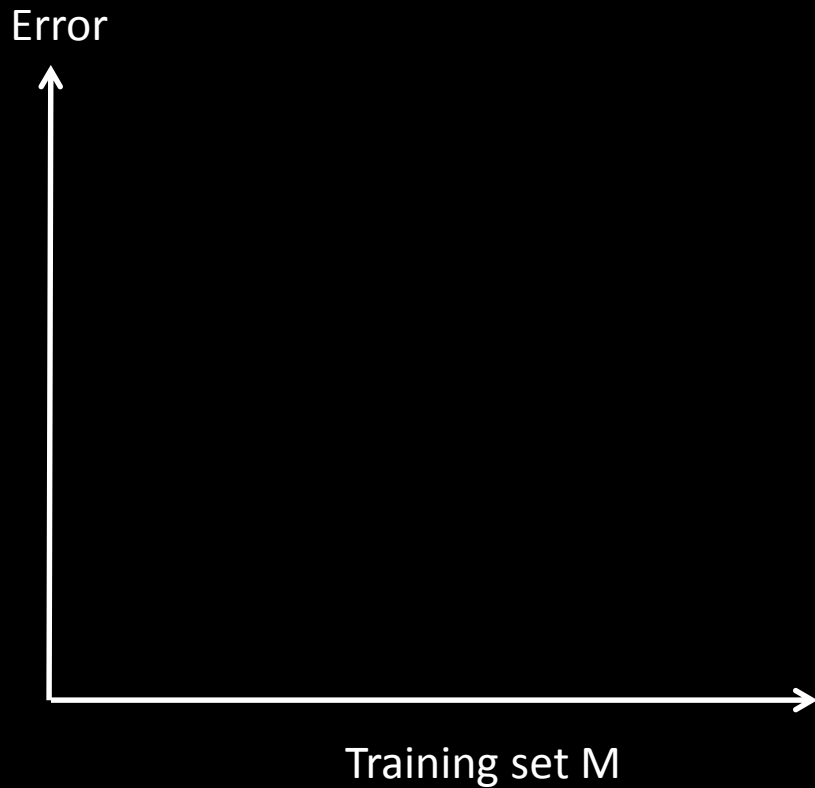
Machine Learning

Nightmare
series

What if your model has a high error?

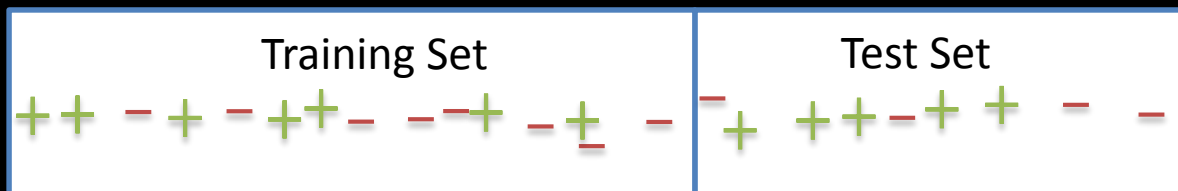
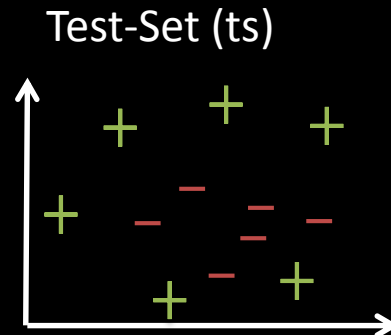
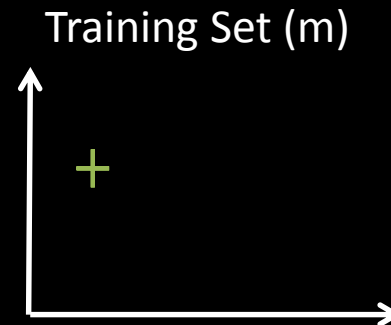
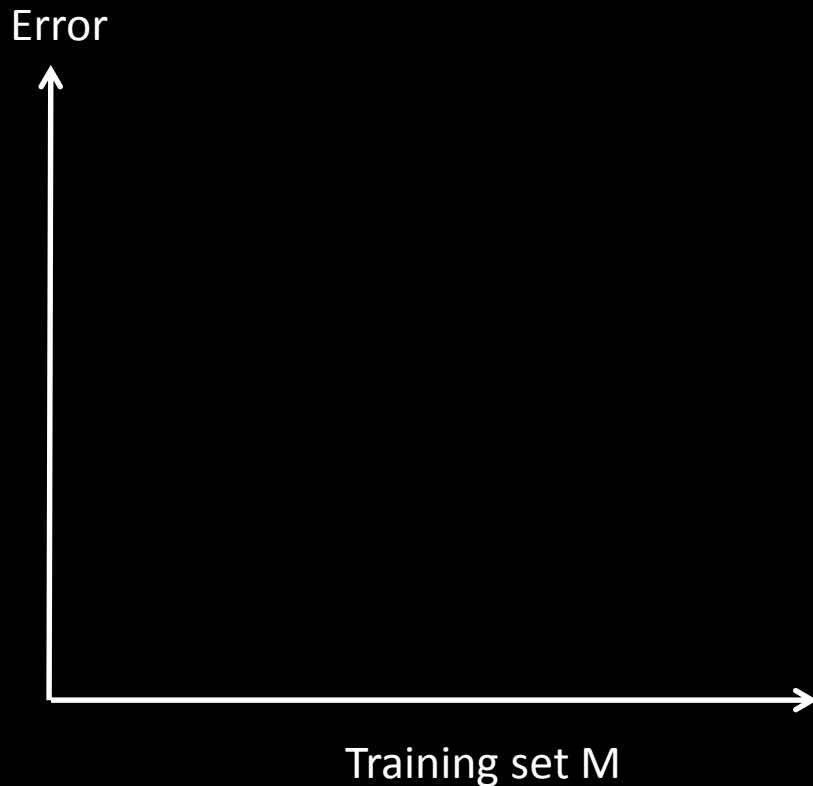
- Get more training examples
- Try smaller sets of features
- Try getting additional features
- Try adding polynomial features (kernels)
- Try decrease regularization
- Try increase regularization

Bias and Variance

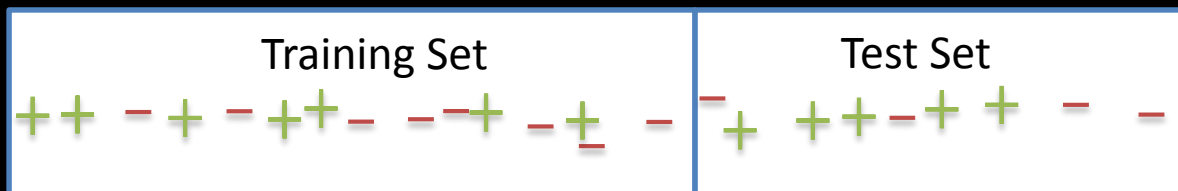
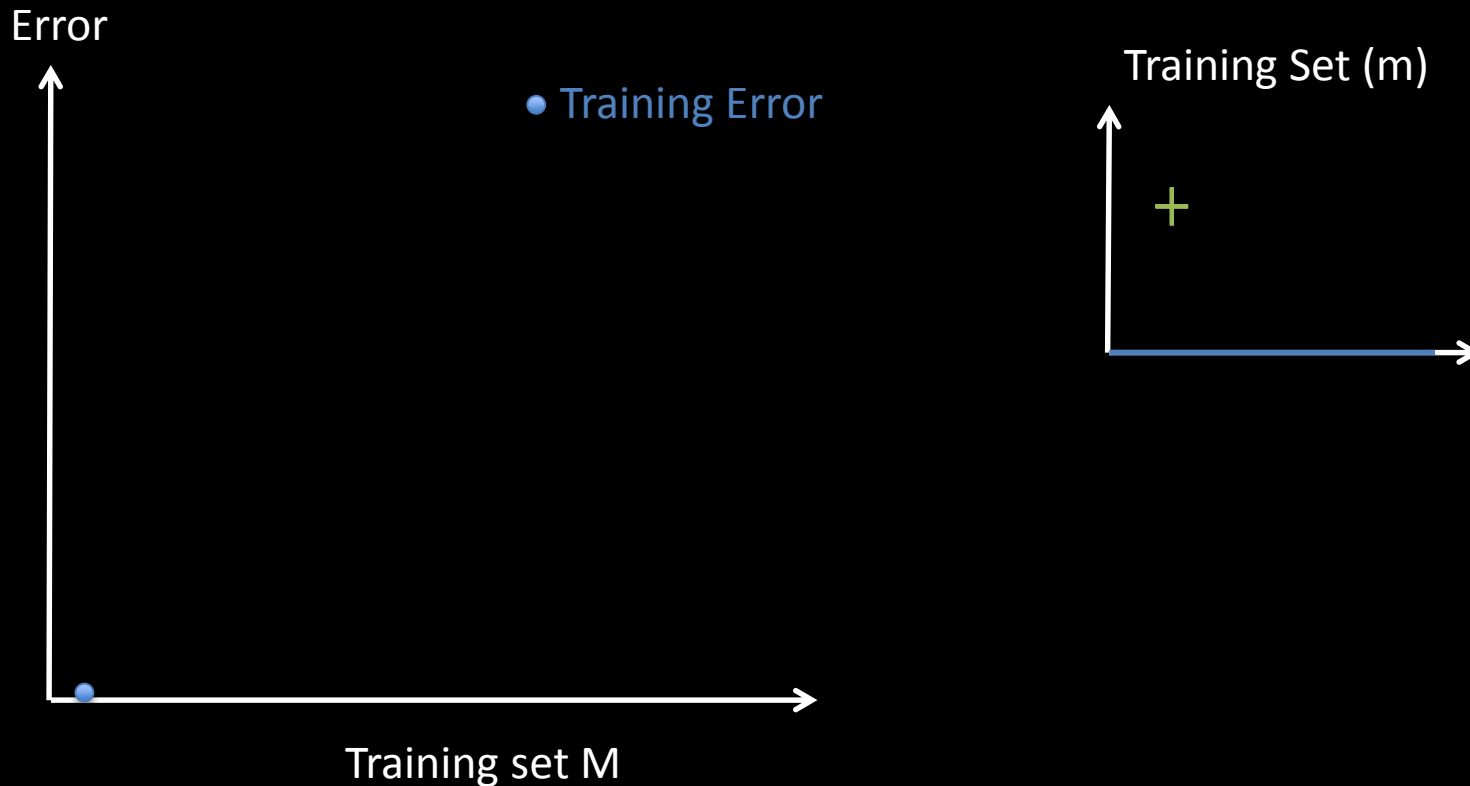


++ - + - ++ - - - + - + - - + ++ - + + - -

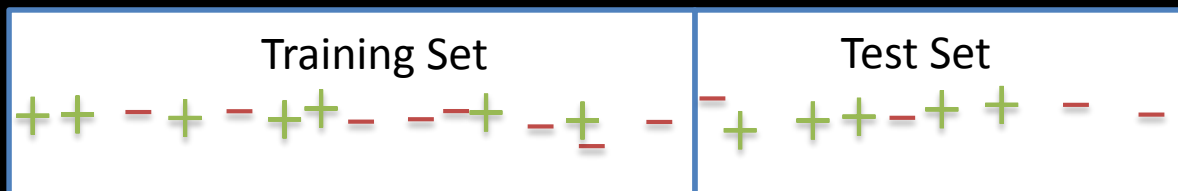
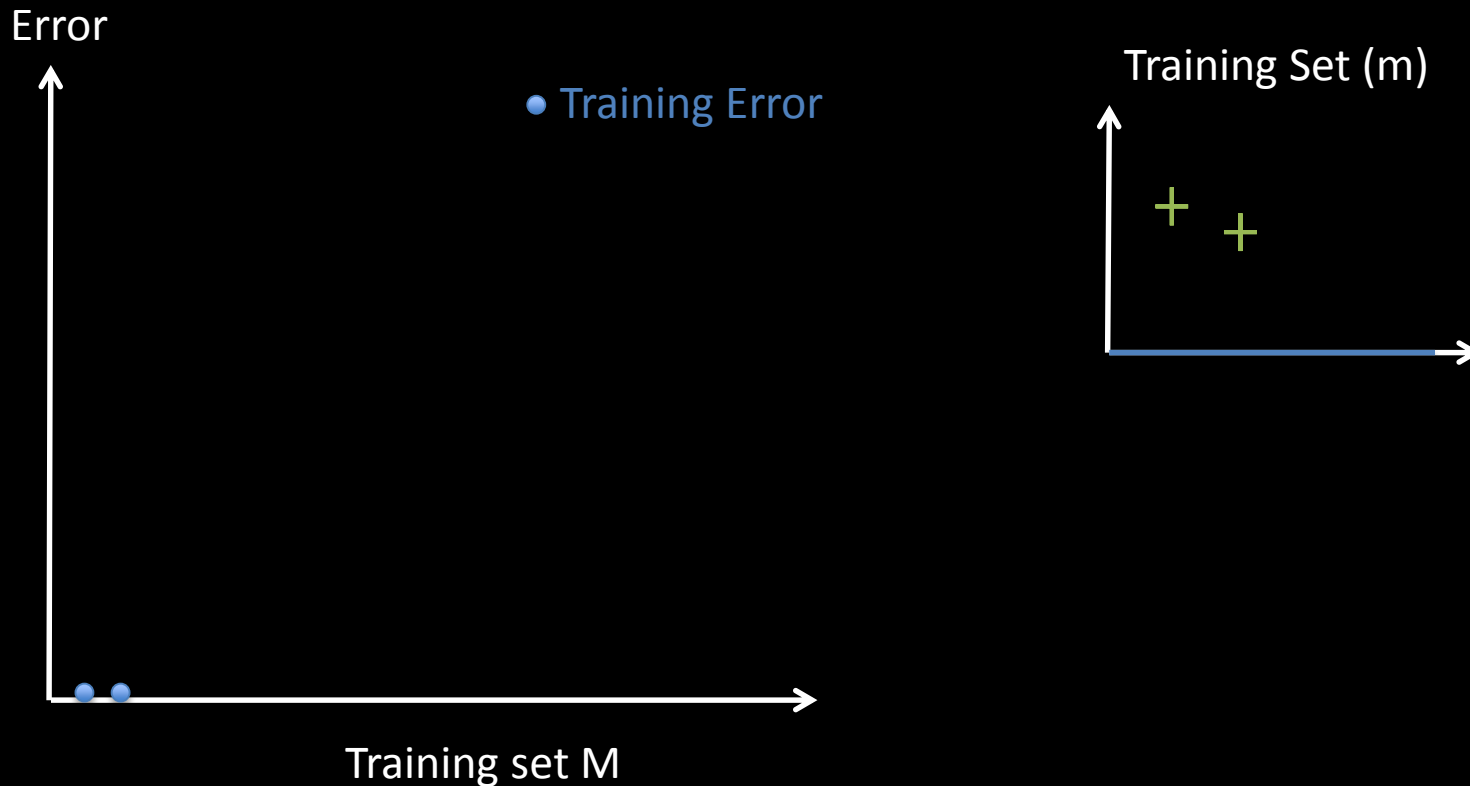
Bias and Variance



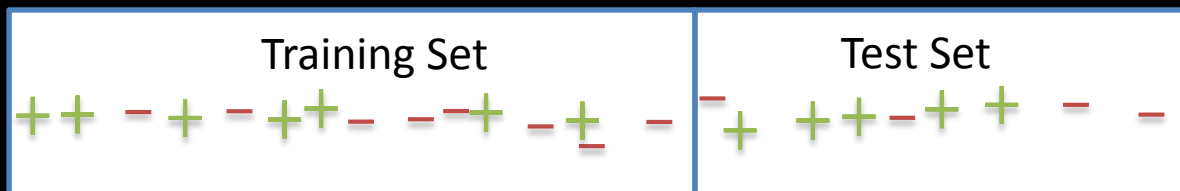
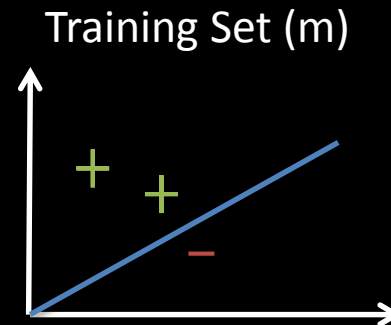
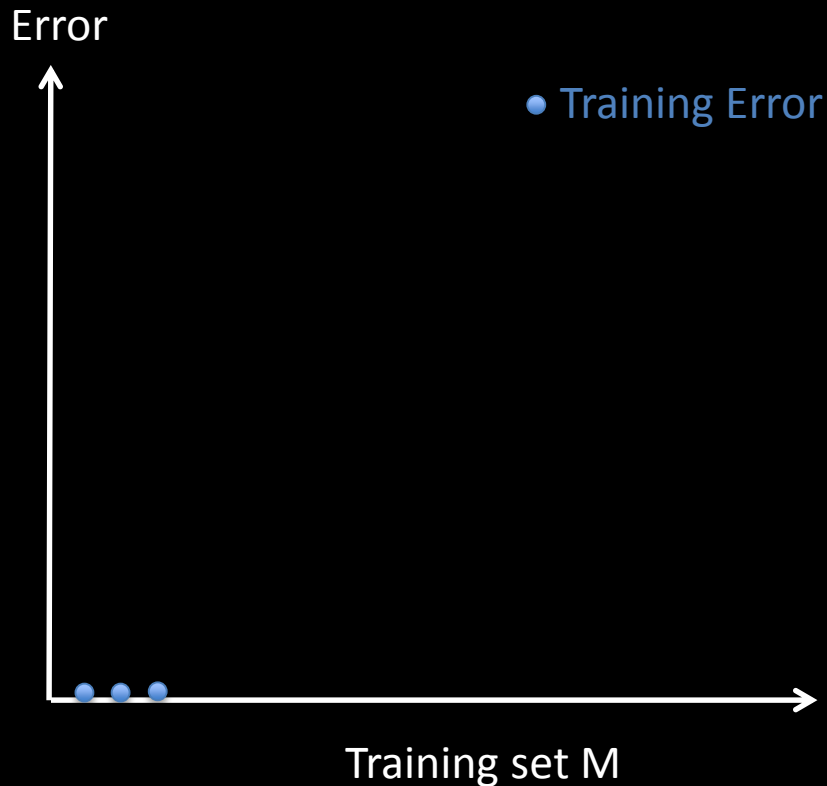
Bias and Variance



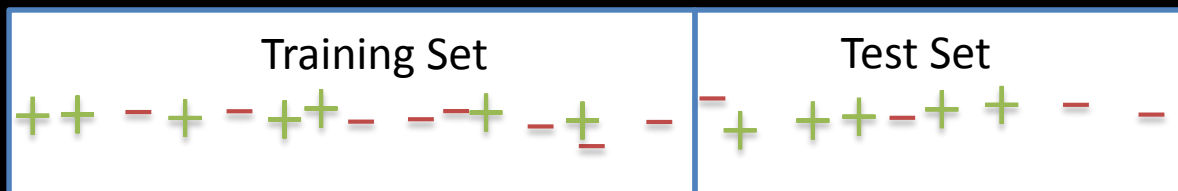
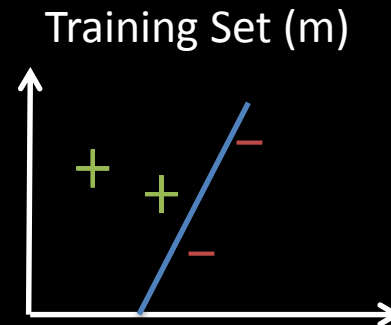
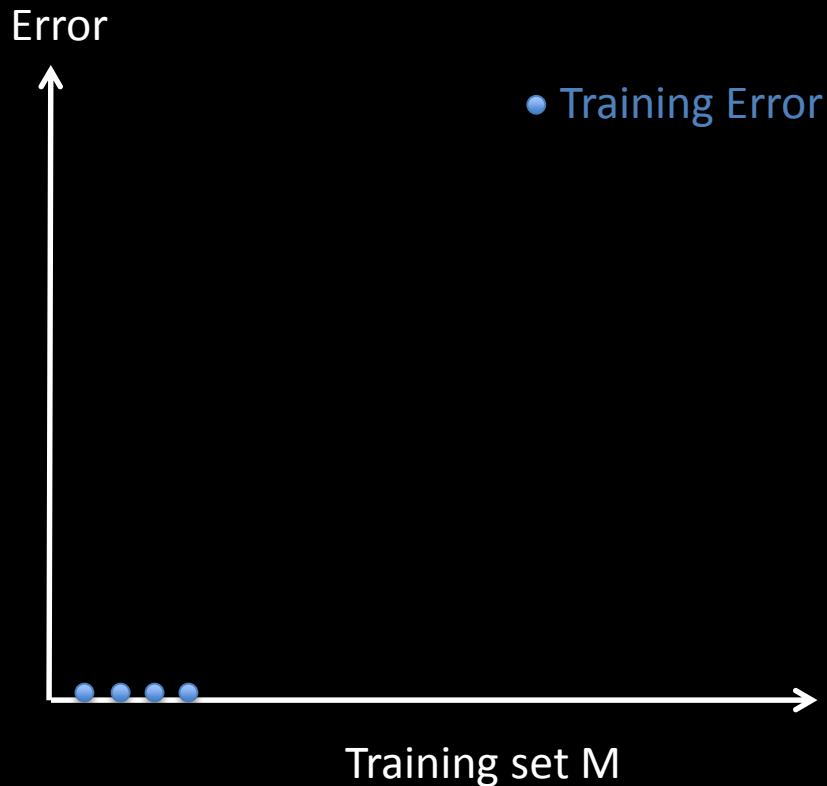
Bias and Variance



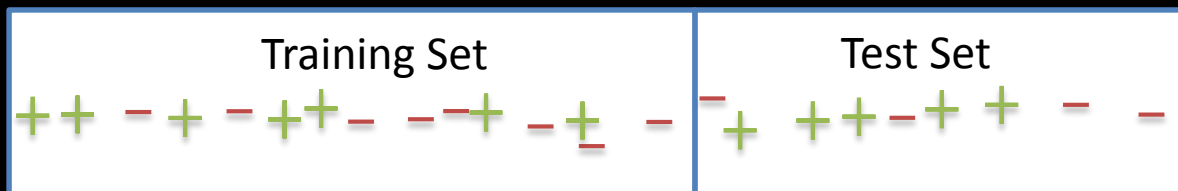
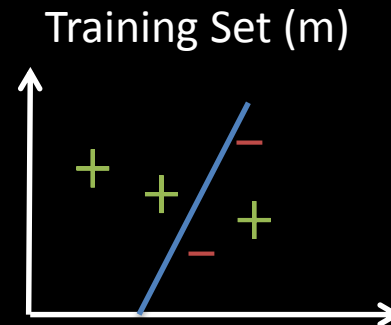
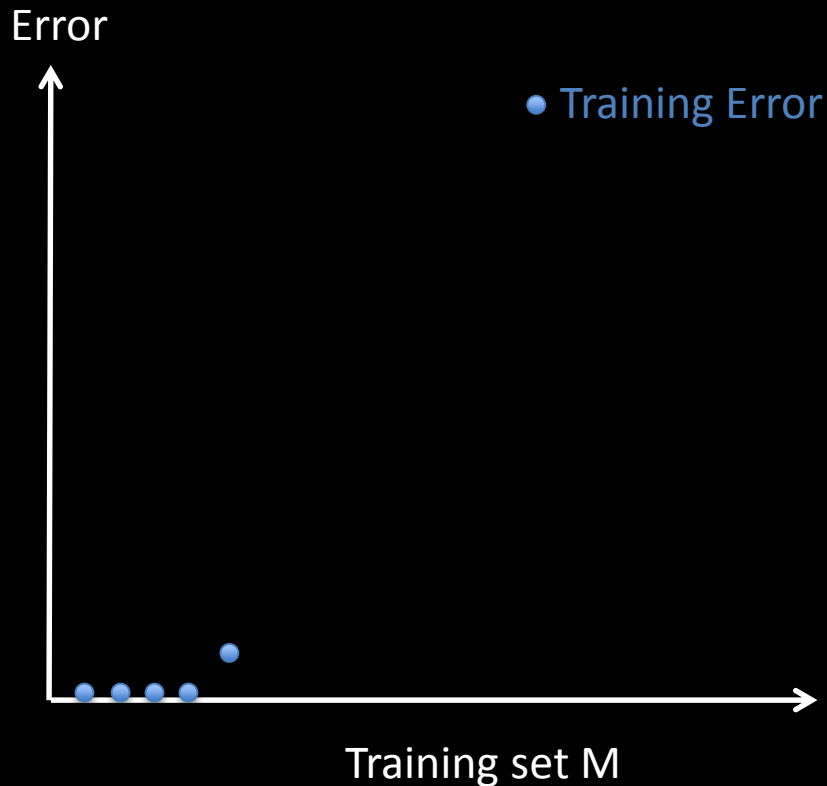
Bias and Variance



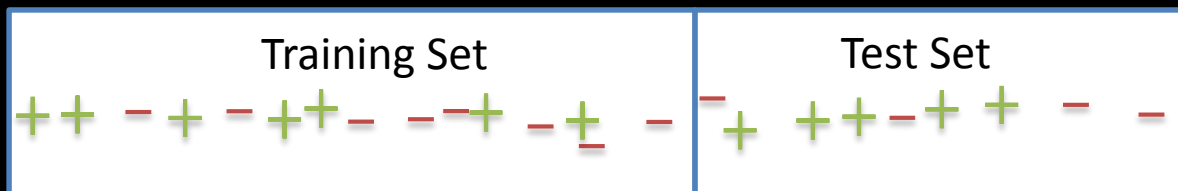
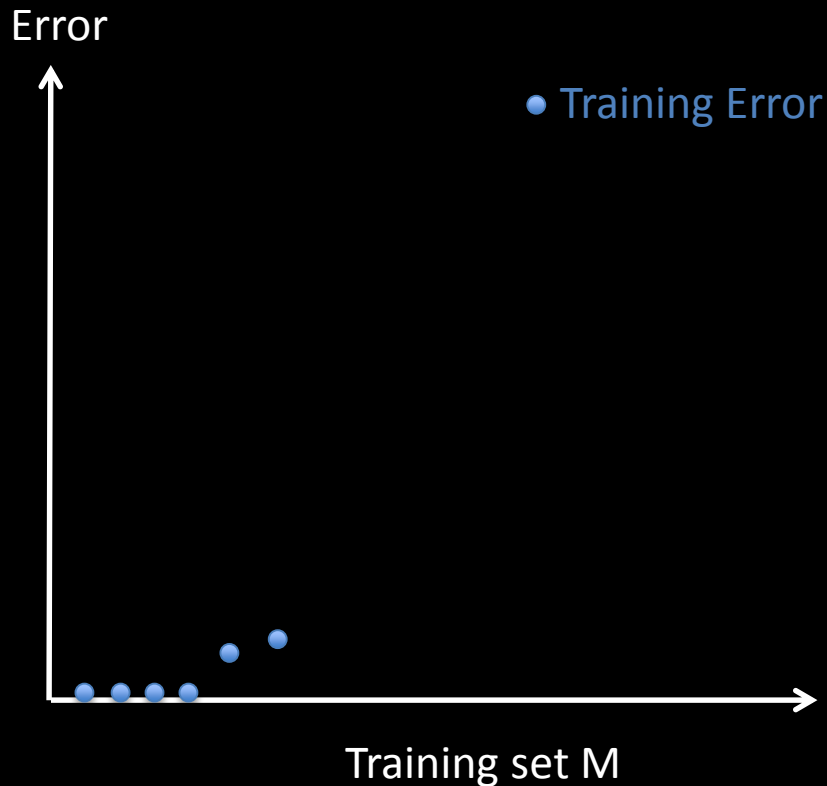
Bias and Variance



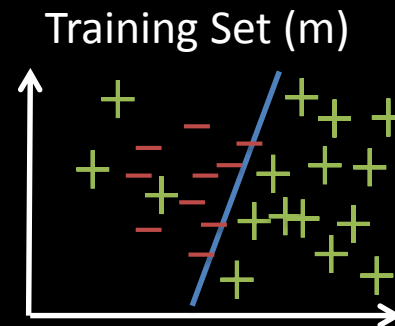
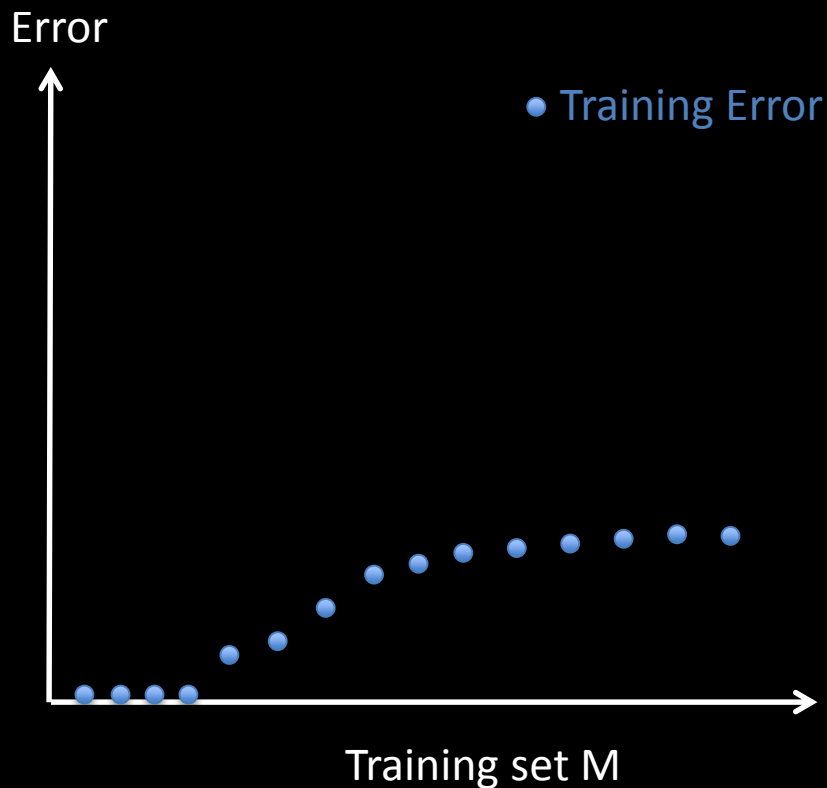
Bias and Variance



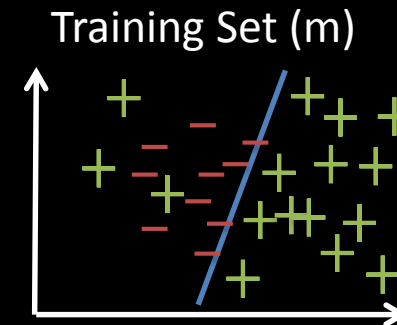
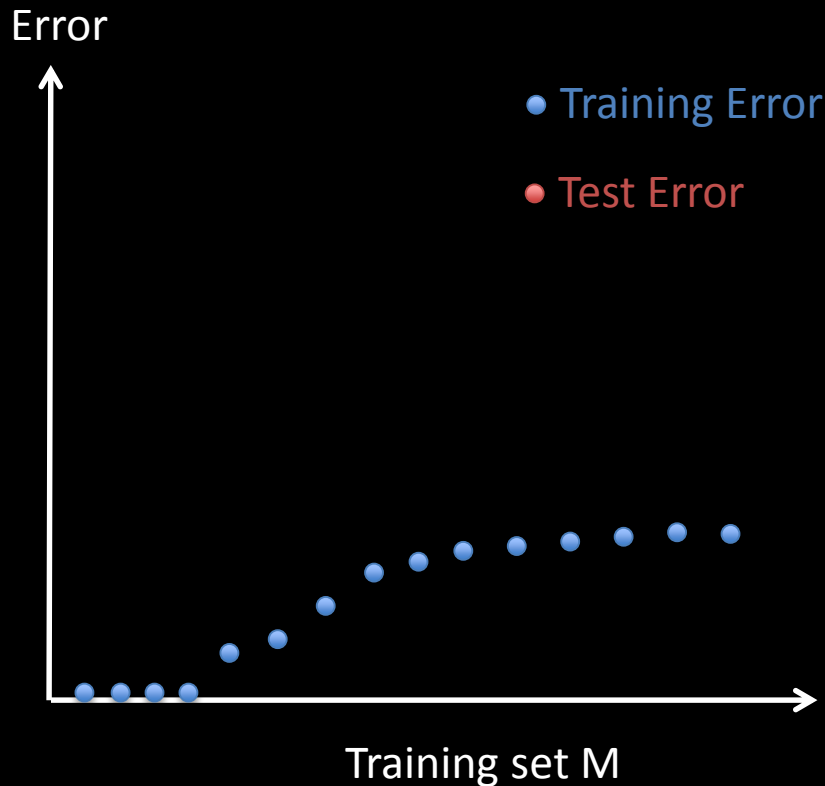
Bias and Variance



Bias and Variance



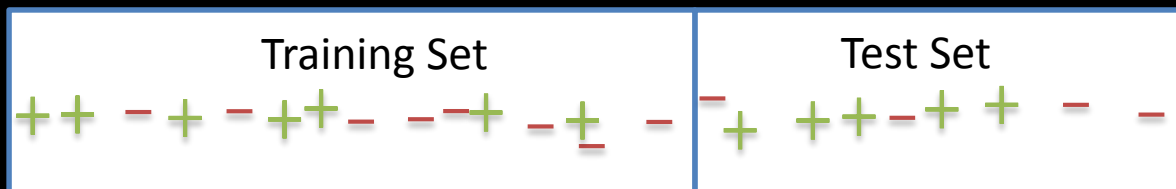
Bias and Variance



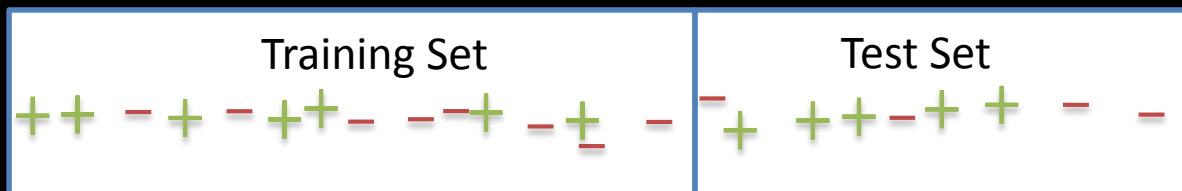
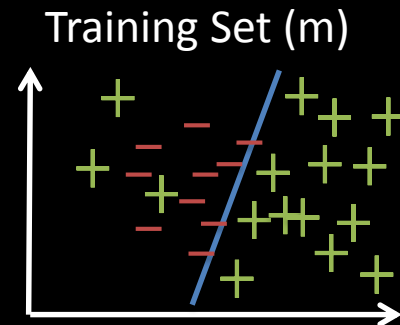
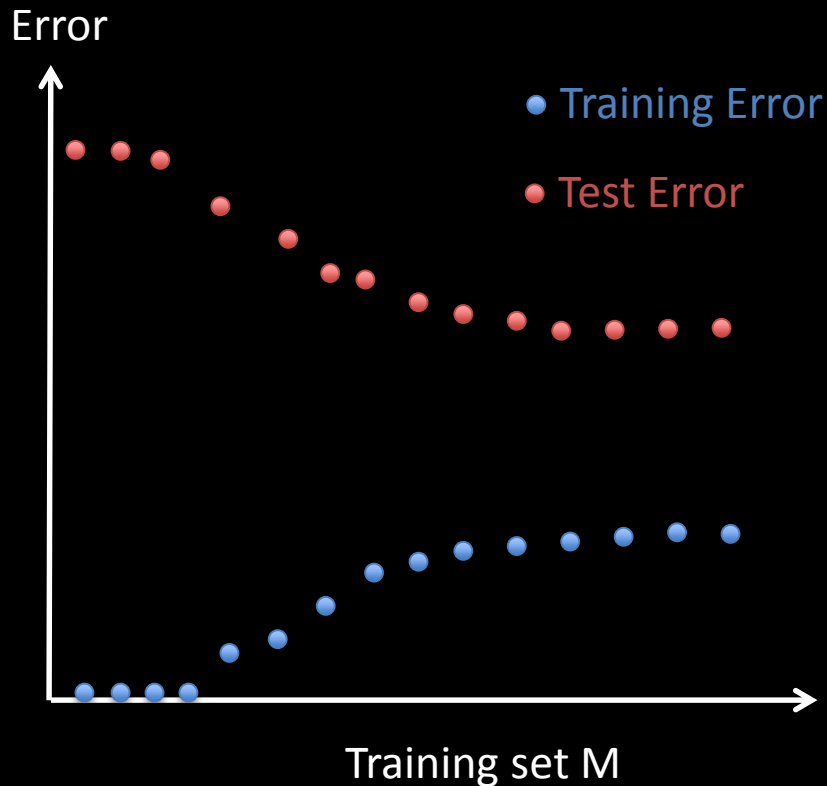
Clicker:

Training error

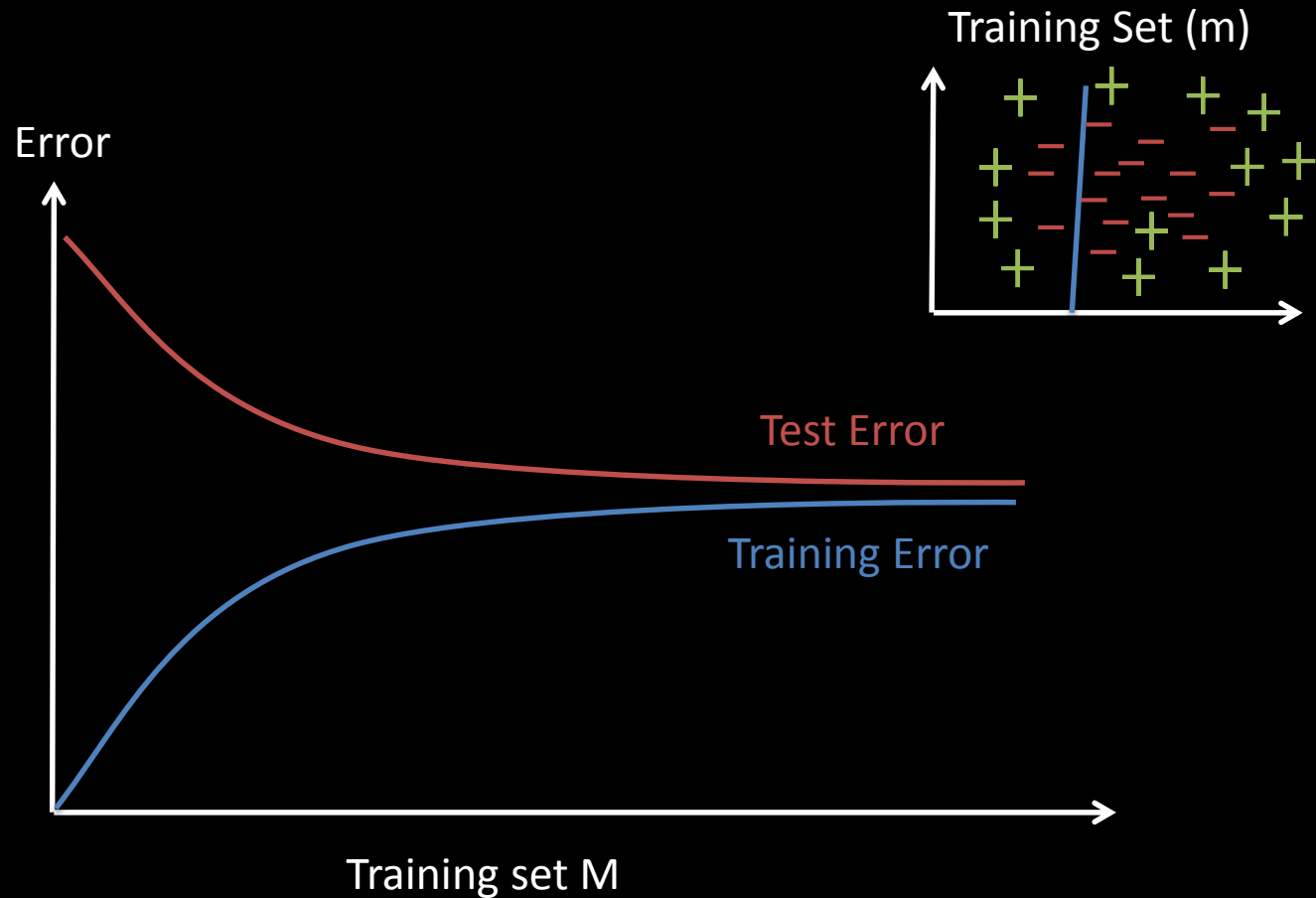
- a) decreases with M
- b) increases with M
- c) stays constant



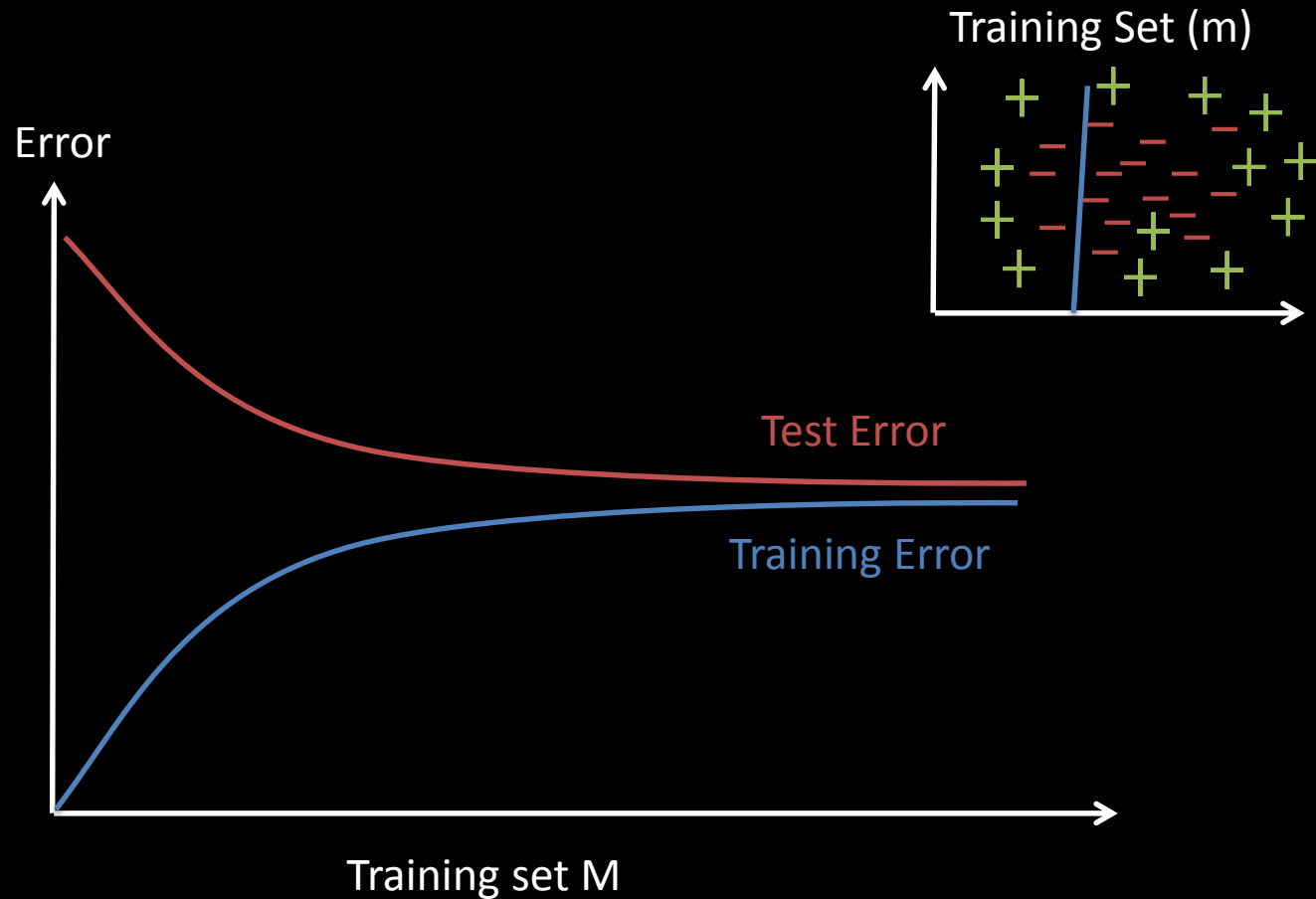
Bias and Variance



High Bias



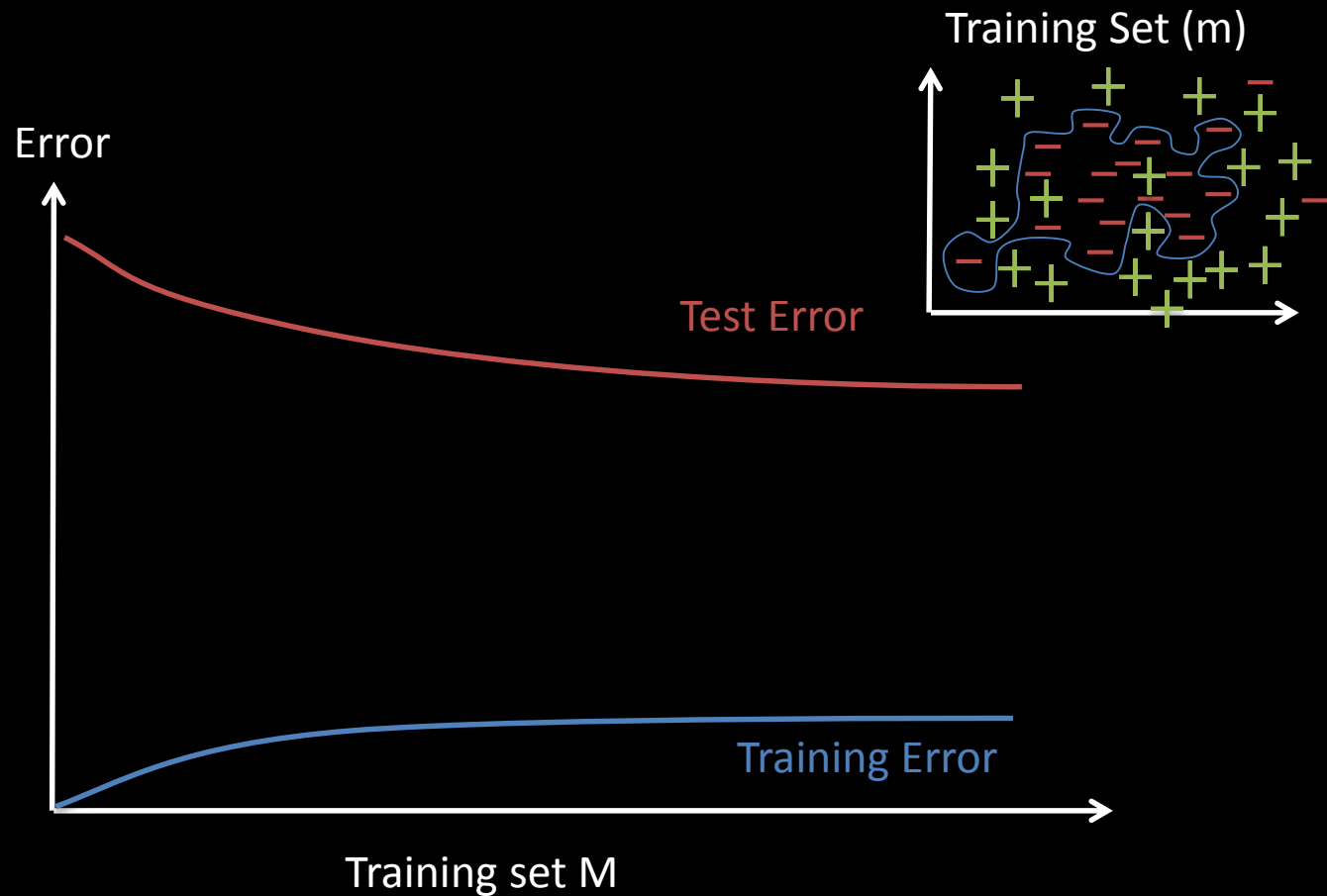
High Bias



Clicker: If you have high-bias, does more data help?

- a) No
- b) Yes

High Variance



Overfitting

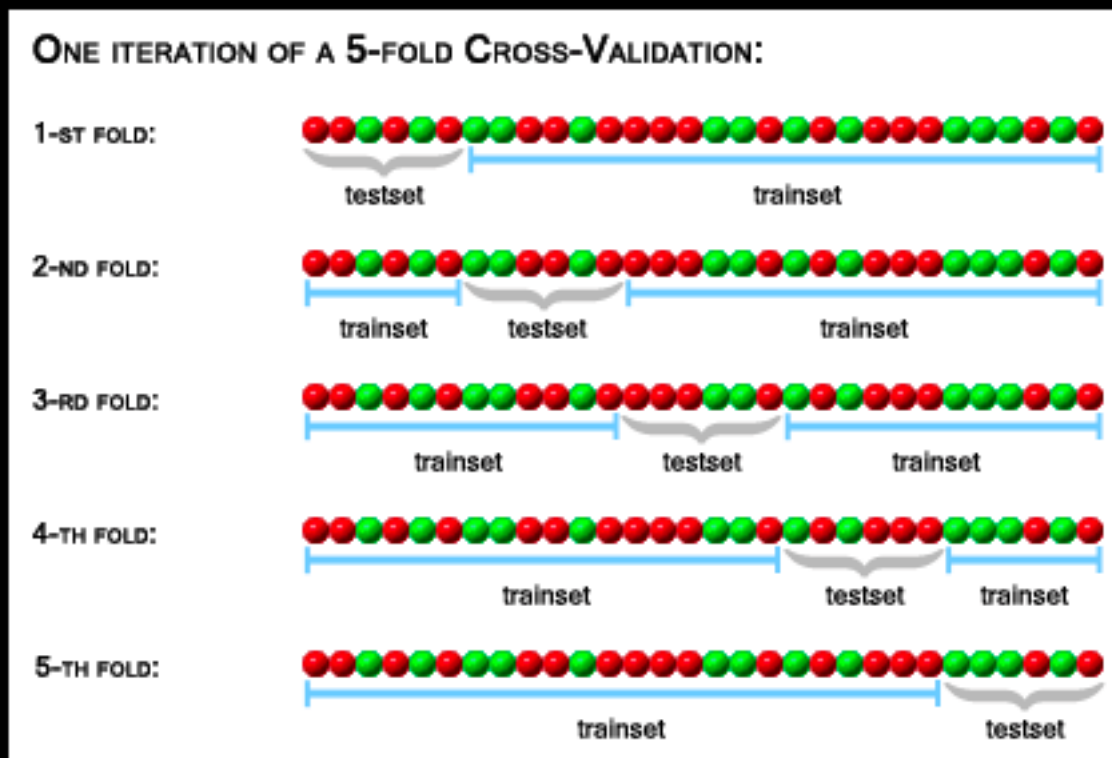
What if the knowledge and data we have are not sufficient to completely determine the correct classifier? Then we run the risk of just hallucinating a classifier (or parts of it) that is not grounded in reality, and is simply encoding random quirks in the data.

This problem is called *overfitting*, and is the bugbear of machine learning. When your learner outputs a classifier that is 100% accurate on the training data but only 50% accurate on test data, when in fact it could have output one that is 75% accurate on both, it has overfit.

Cross-validation

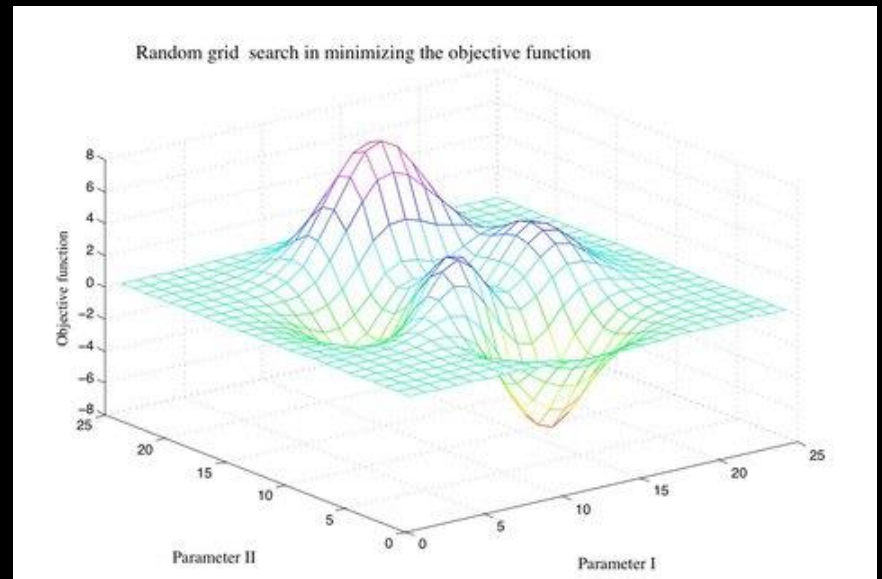
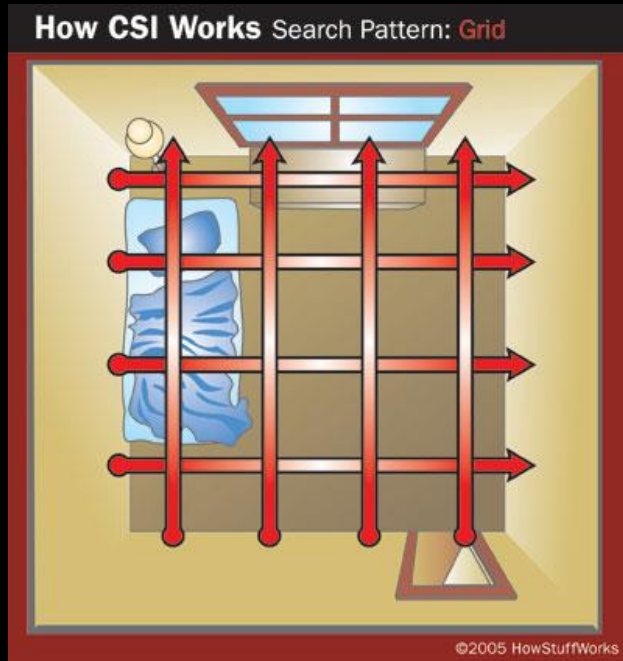
k-fold: split the data into k groups, train on every group except for one, which you test on.

Repeat for all groups



Parameter Tuning

Grid Search



Many classifiers to choose from

- Decision Trees
- K-nearest neighbor
- Support Vector Machines
- Logistic Regression
- Naïve Bayes
- Random Forrest
- Bayesian network
- Randomized Forests
- Boosted Decision Trees
- RBMs
-