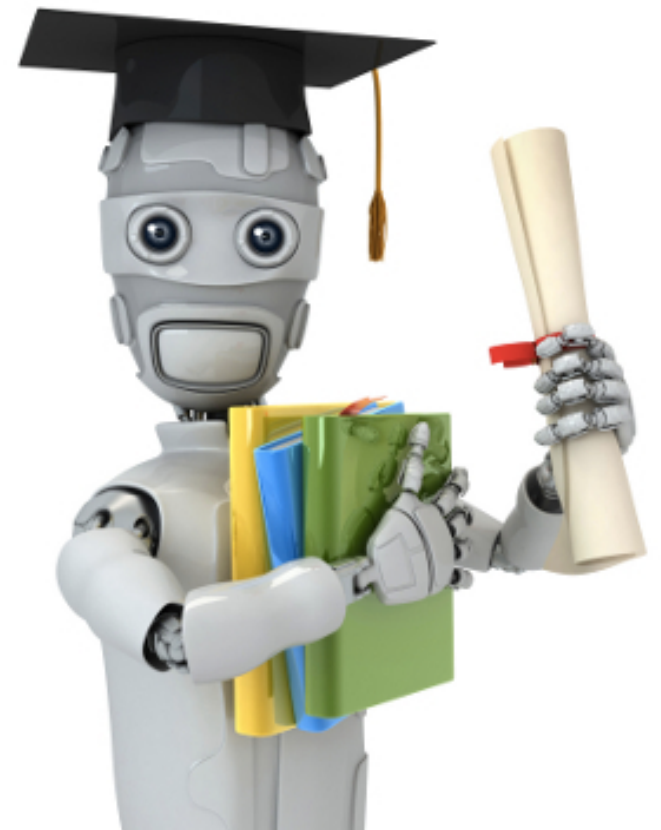


MACHINE LEARNING OVERVIEW





PETER FLACH

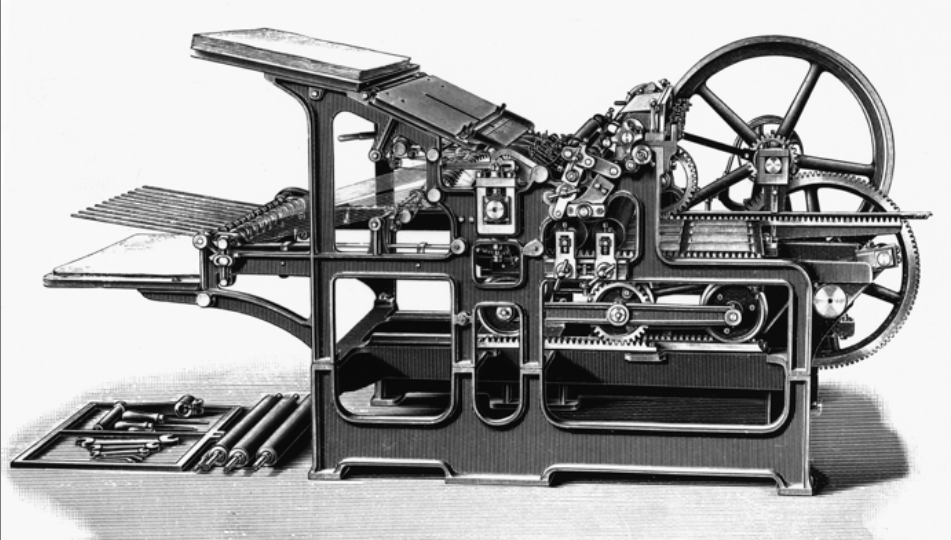
Machine Learning

The Art and Science of Algorithms
that Make Sense of Data

CAMBRIDGE



VS



MACHINE LEARNING PROBLEMS

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

MACHINE LEARNING PROBLEMS

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

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

$$y = f(x)$$

output prediction
 function

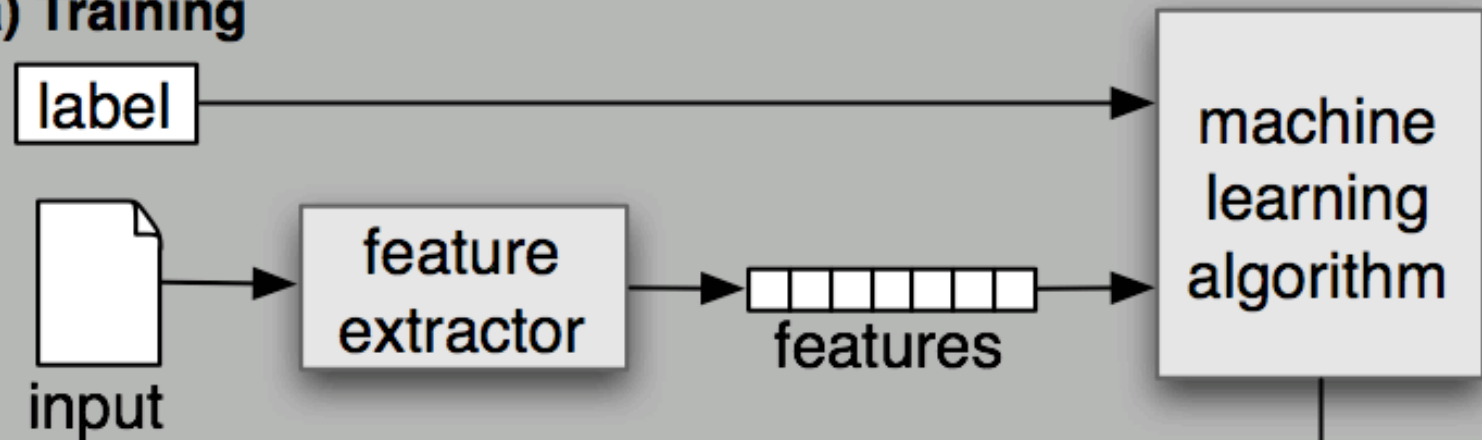
features

Training: given a *training set* of labeled examples $\{(x_1, y_1), \dots, (x_N, y_N)\}$, estimate the prediction function f by minimizing the prediction error on the training set

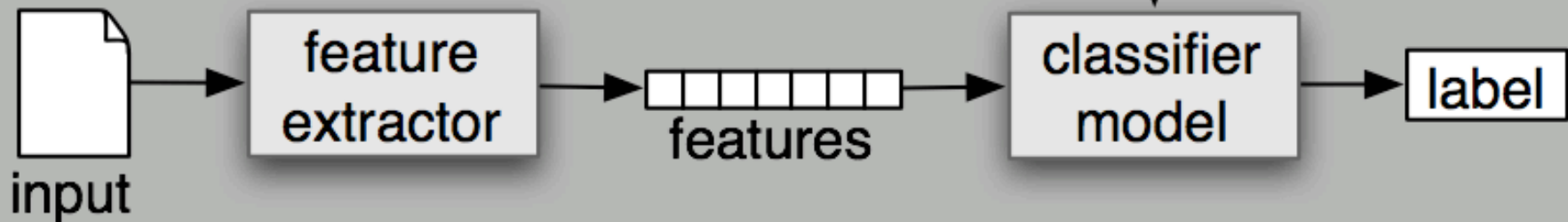
Testing: apply f to a never before seen *test example* x and output the predicted value $y = f(x)$

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
- ...

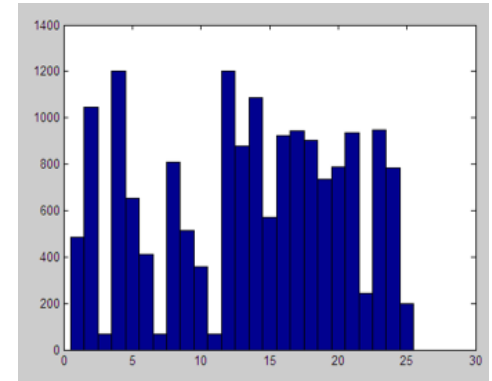
Customer State	Product Type	Product Weight	Volume (L*H*D)	Month	Delivery Time

IMAGE FEATURES

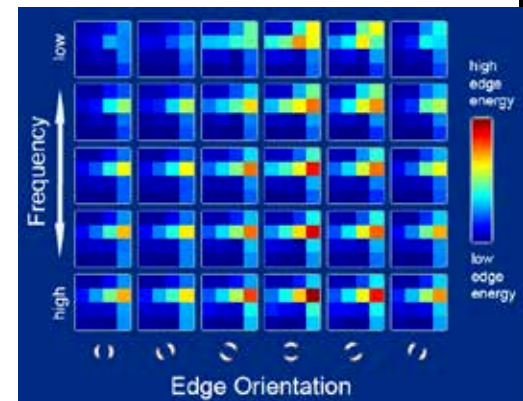
Raw pixels



Histograms



GIST descriptors



...

TEXT FEATURES

✉ Tamara Mccullough FDA approved on-line pharmacie
✉ Mail Delivery System Mail delivery failed: returning me

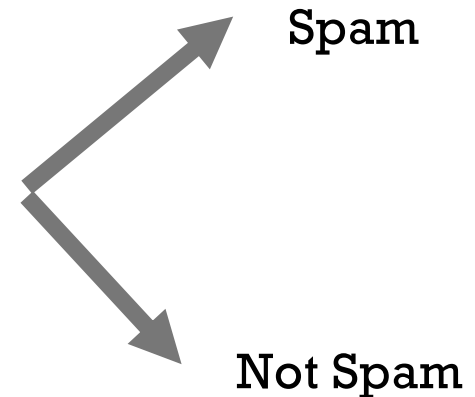
From: Tamara Mccullough **To:** Tom;
Subject: FDA approved on-line pharmacies

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

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

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

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



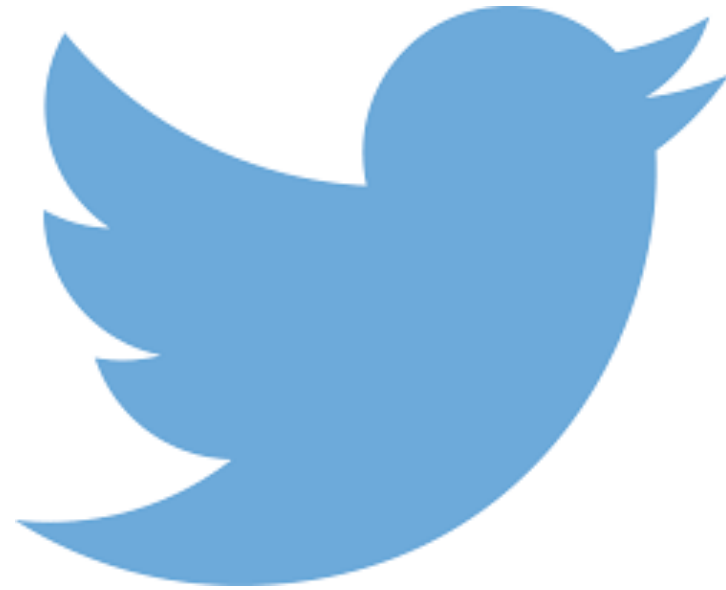
Bag of Words

Viagra
Soft
Herbel
Pills
Are
...

N-Grams

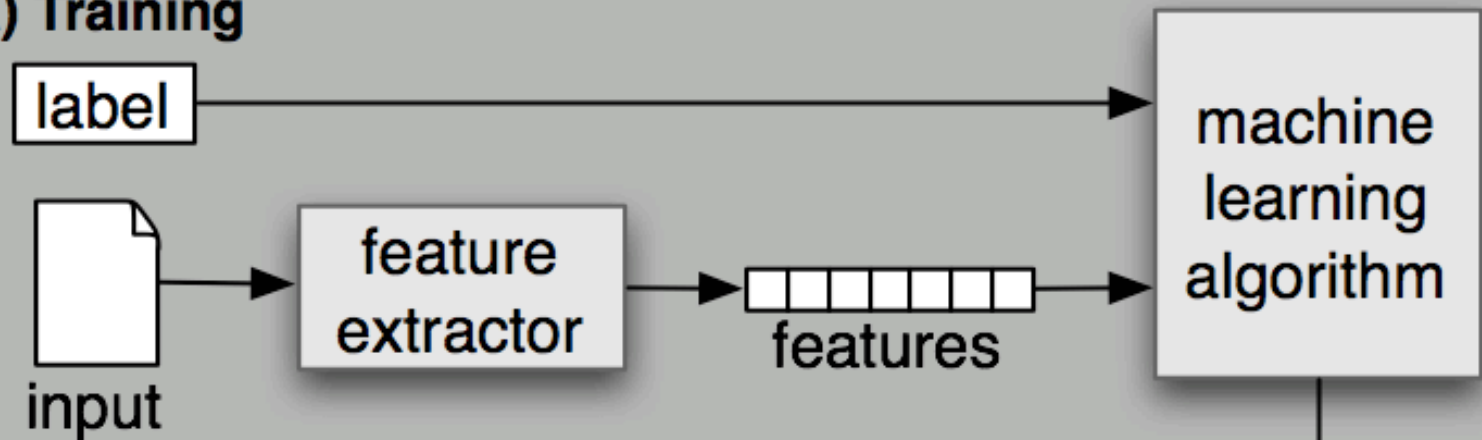
herbel pills
pills for
for Hair
Hair enlargement
enlargement Techniques
...

FEATURE TO PREDICT UNEMPLOYMENT

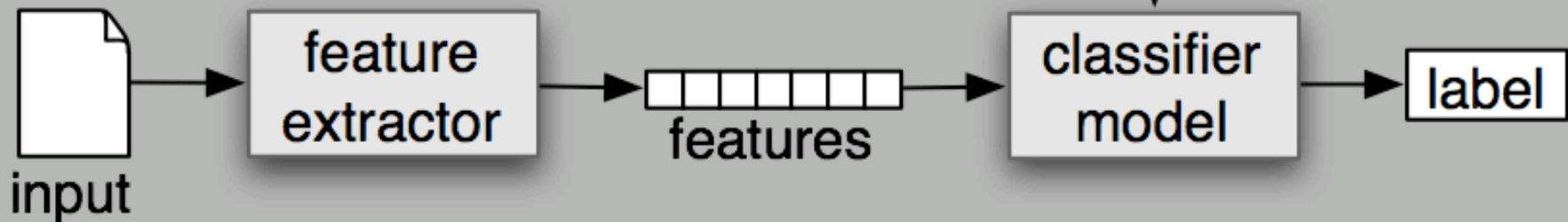


ML PIPELINE (SUPERVISED)

(a) Training



(b) Prediction



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		

MANY CLASSIFIERS TO CHOOSE FROM

K-nearest neighbor

Support Vector Machines

Decision Trees

Random Forrest

(Gradient) Boosted Decision Trees

Logistic Regression

Naïve Bayes

Bayesian network

RBM

....

Which is the best one?

MANY CLASSIFIERS TO CHOOSE FROM

K-nearest neighbor

Support Vector Machines

Decision Trees

Random Forrest

(Gradient) Boosted Decision Trees

Logistic Regression

Naïve Bayes

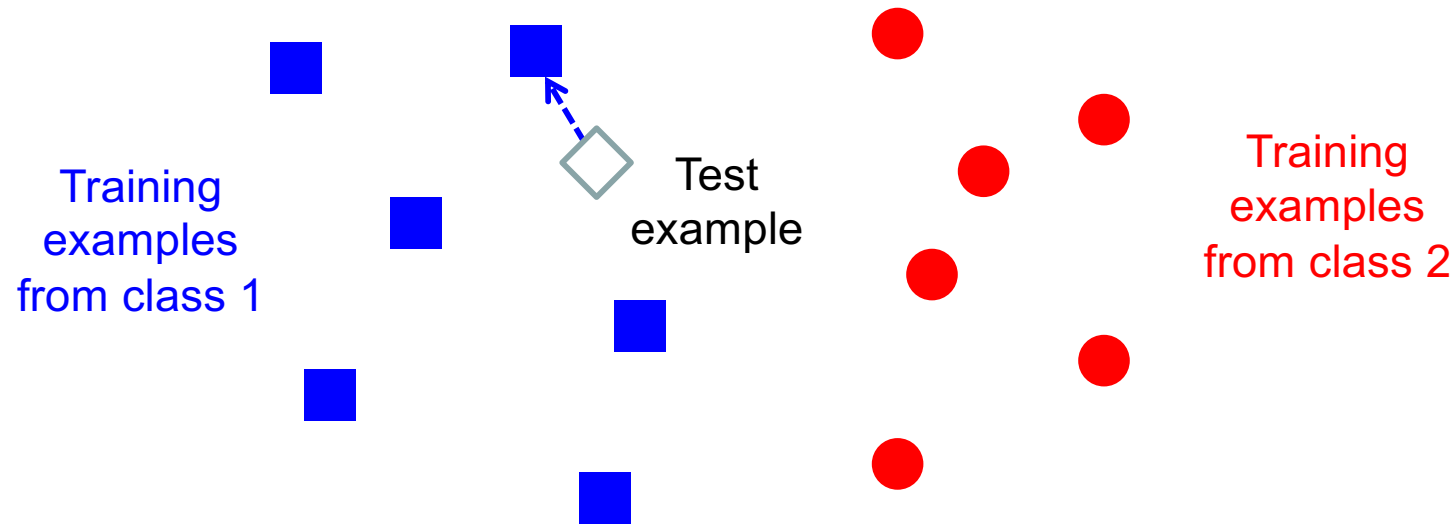
Bayesian network

RBM

....

Which is the best one?

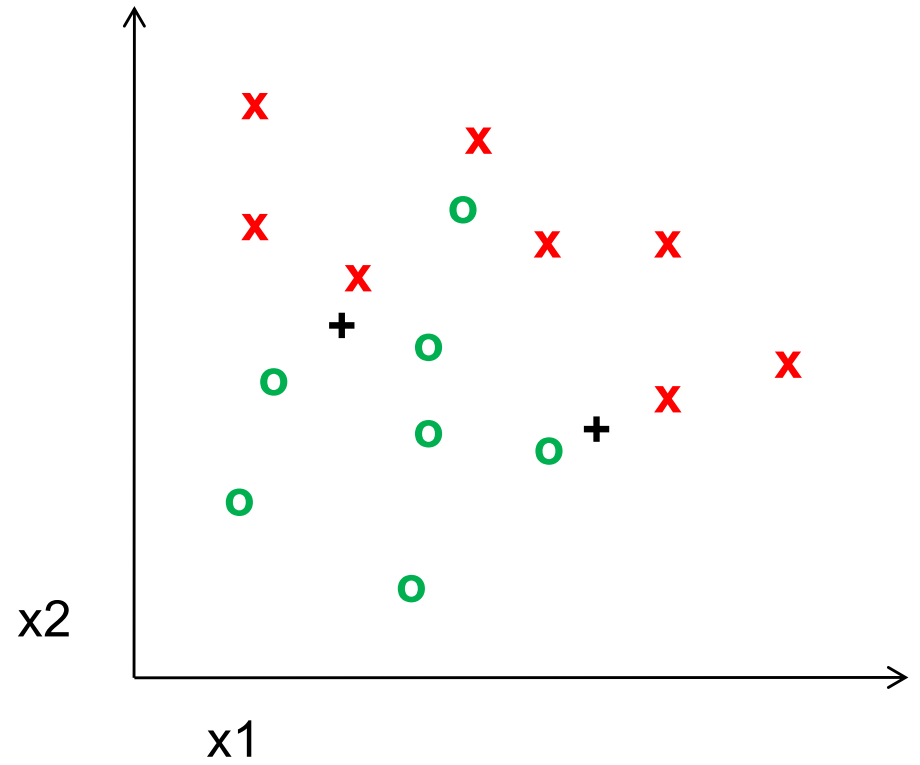
CLASSIFIERS: NEAREST NEIGHBOR



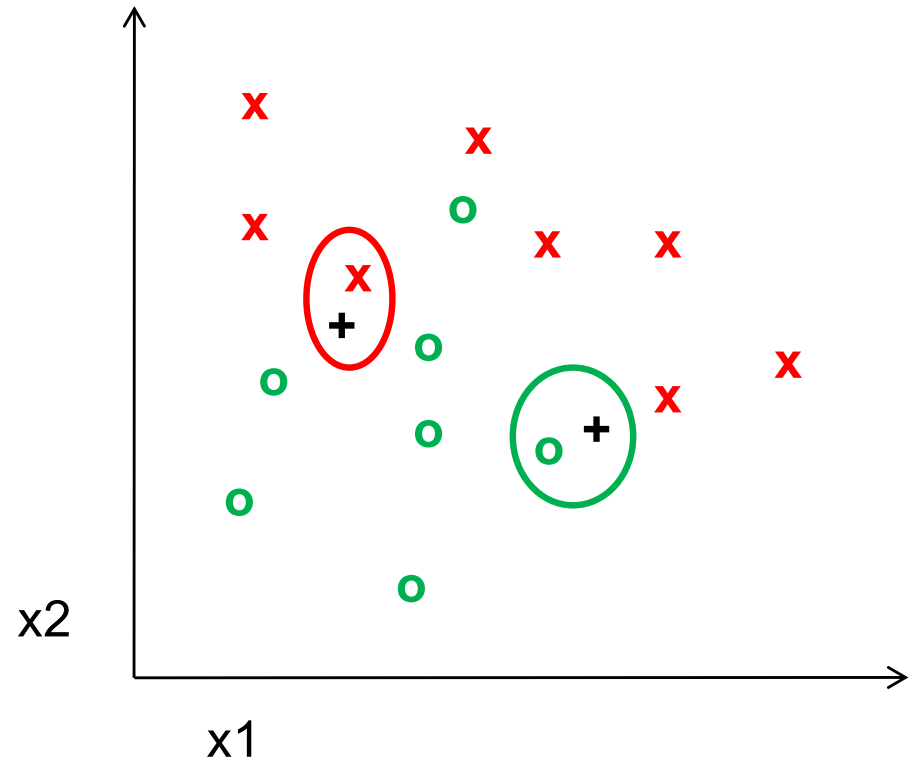
$f(\mathbf{x}) = \text{label of the training example nearest to } \mathbf{x}$

- All we need is a distance function for our inputs
- No training required!

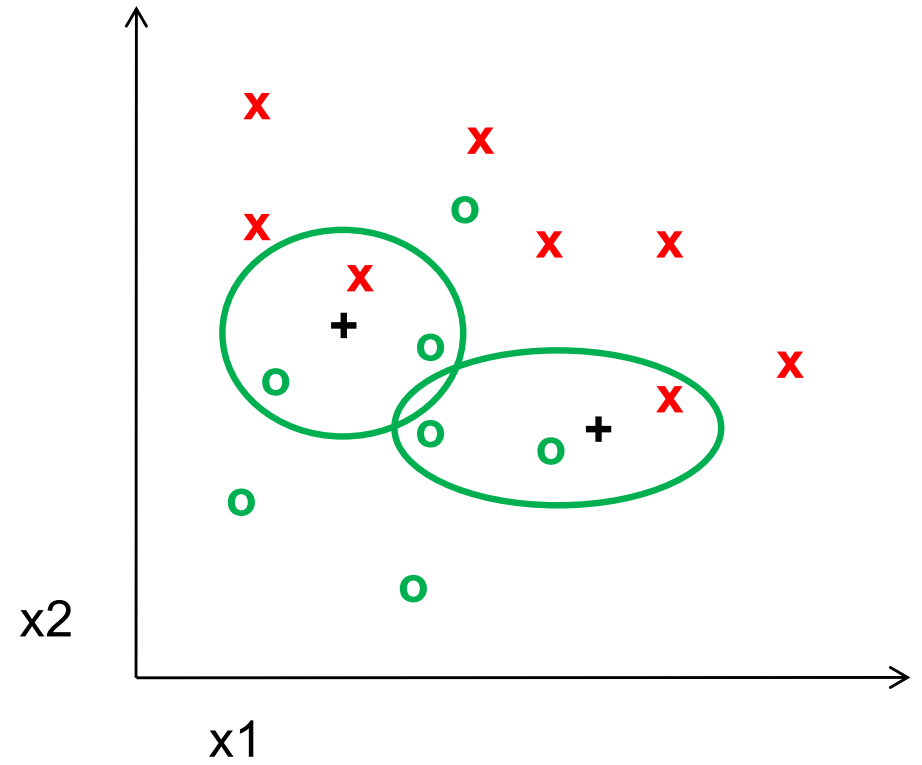
K-NEAREST NEIGHBOR



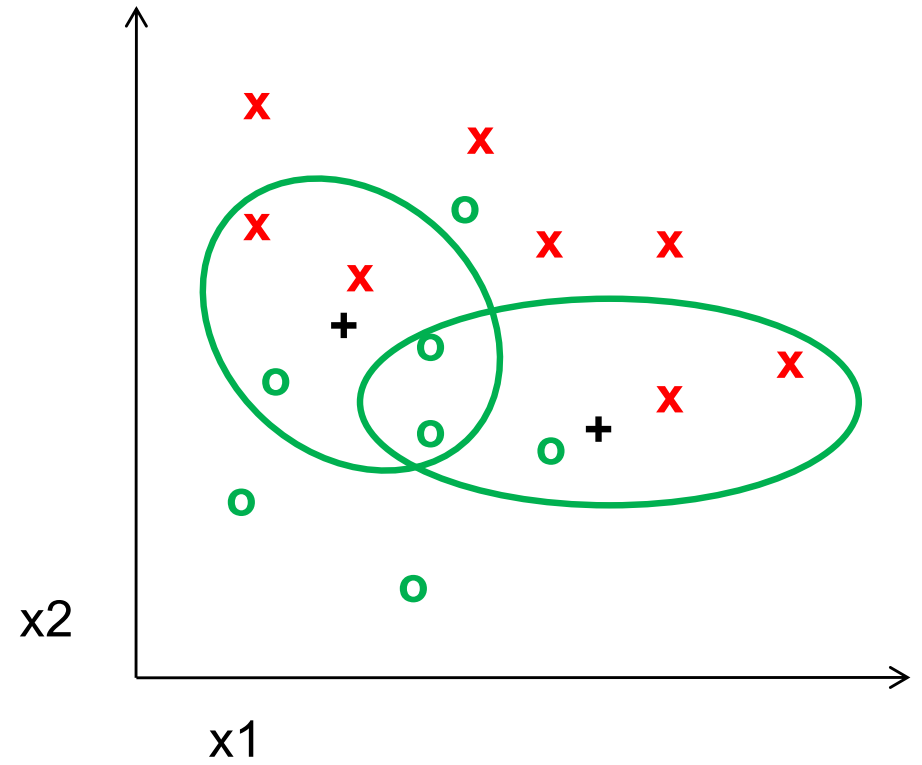
1-NEAREST NEIGHBOR



3-NEAREST NEIGHBOR

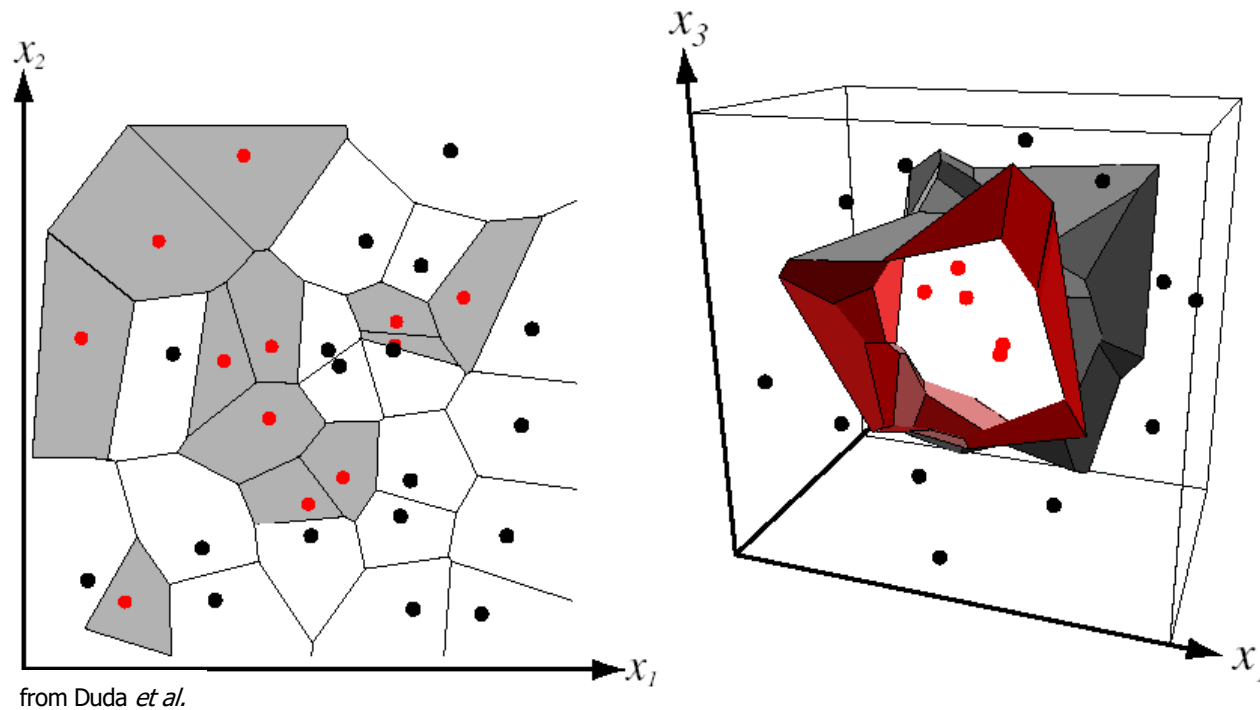


5-NEAREST NEIGHBOR



DECISION BOUNDARIES KNN

Assign label of nearest training data point to each test data point



Voronoi partitioning of feature space
for two-category 2D and 3D data

MANY CLASSIFIERS TO CHOOSE FROM

K-nearest neighbor

Support Vector Machines

Which is the best one?

Decision Trees

Random Forrest

(Gradient) Boosted Decision Trees

Logistic Regression

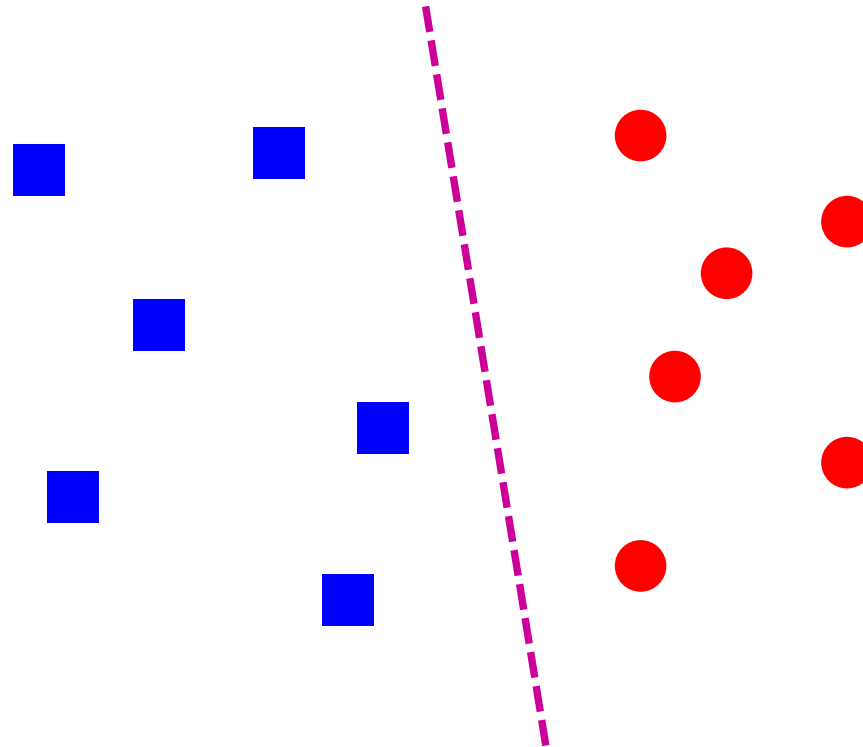
Naïve Bayes

Bayesian network

RBM

....

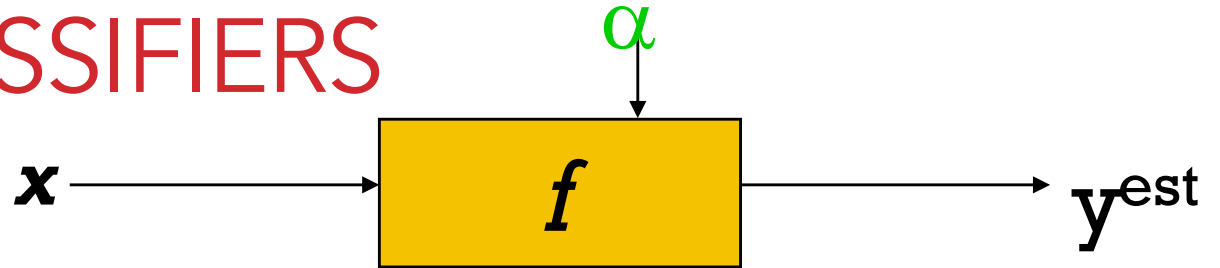
CLASSIFIERS: LINEAR



Find a *linear function* to separate the classes:

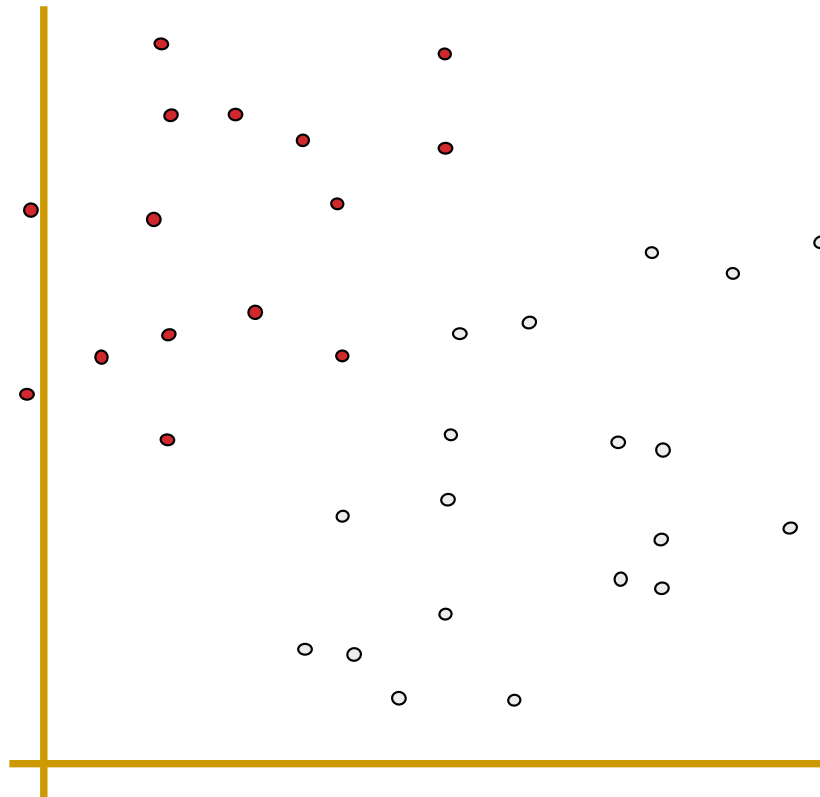
$$f(x) = \text{sgn}(w \cdot x + b)$$

LINEAR CLASSIFIERS



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1



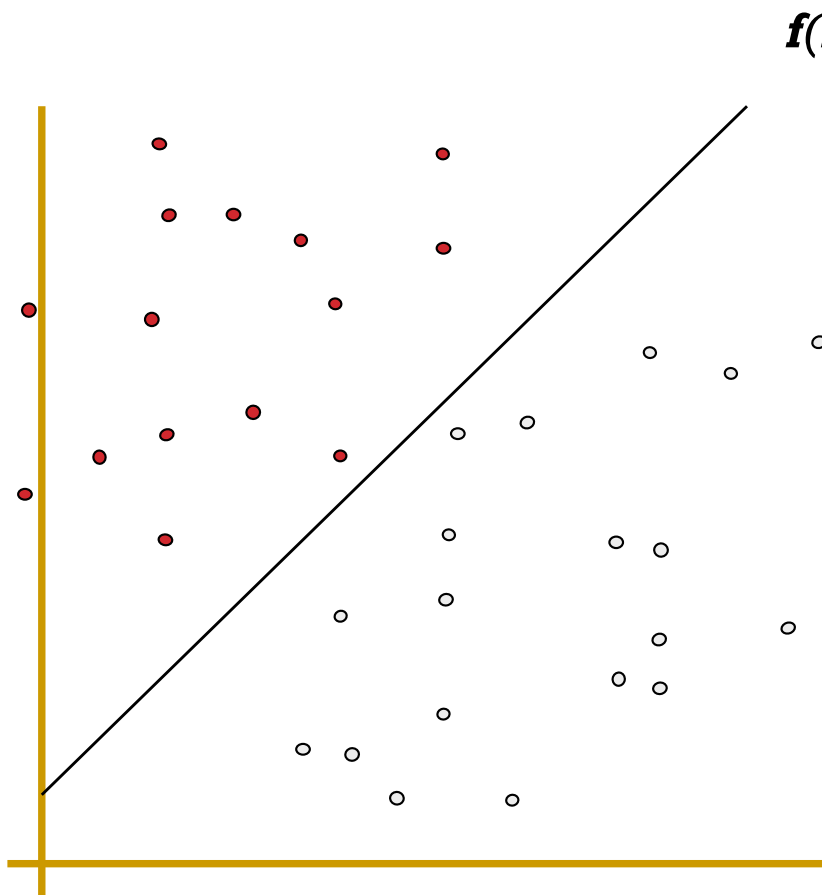
How would you classify this data?

LINEAR CLASSIFIERS

α



- denotes +1
- denotes -1



How would you classify this data?

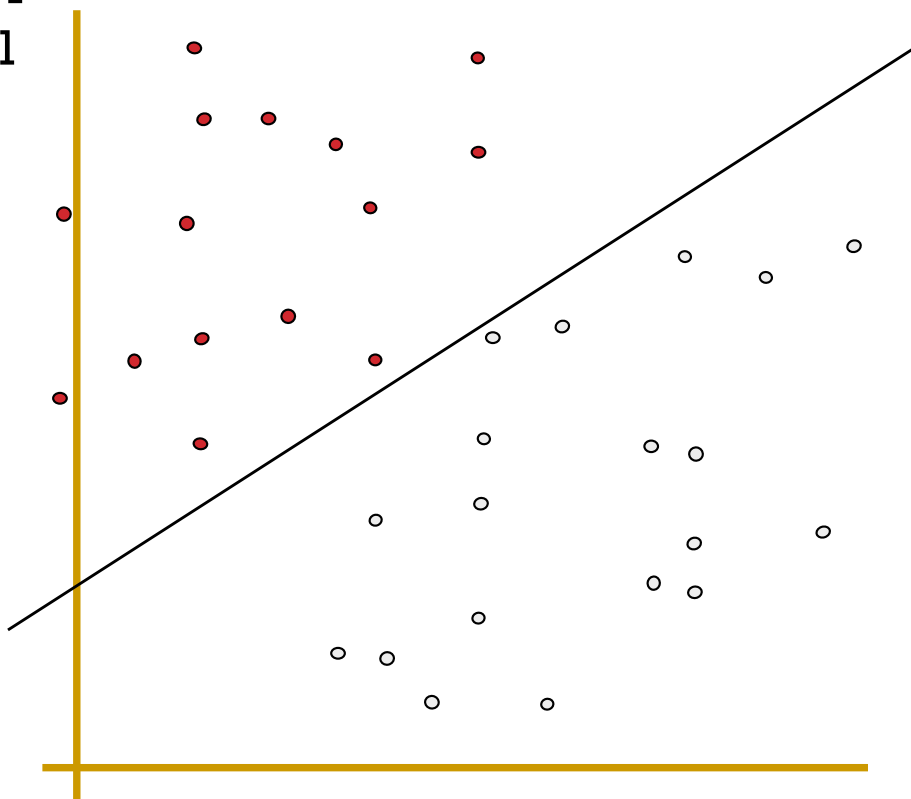
LINEAR CLASSIFIERS

α



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1



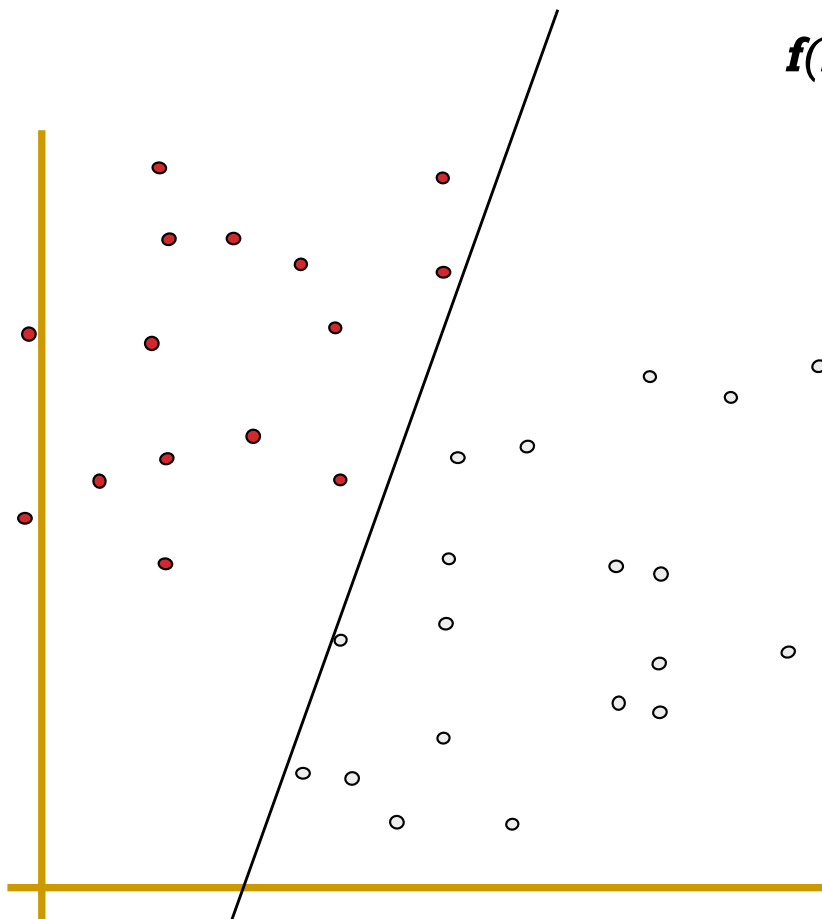
How would you classify this data?

LINEAR CLASSIFIERS

α



- denotes +1
- denotes -1



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

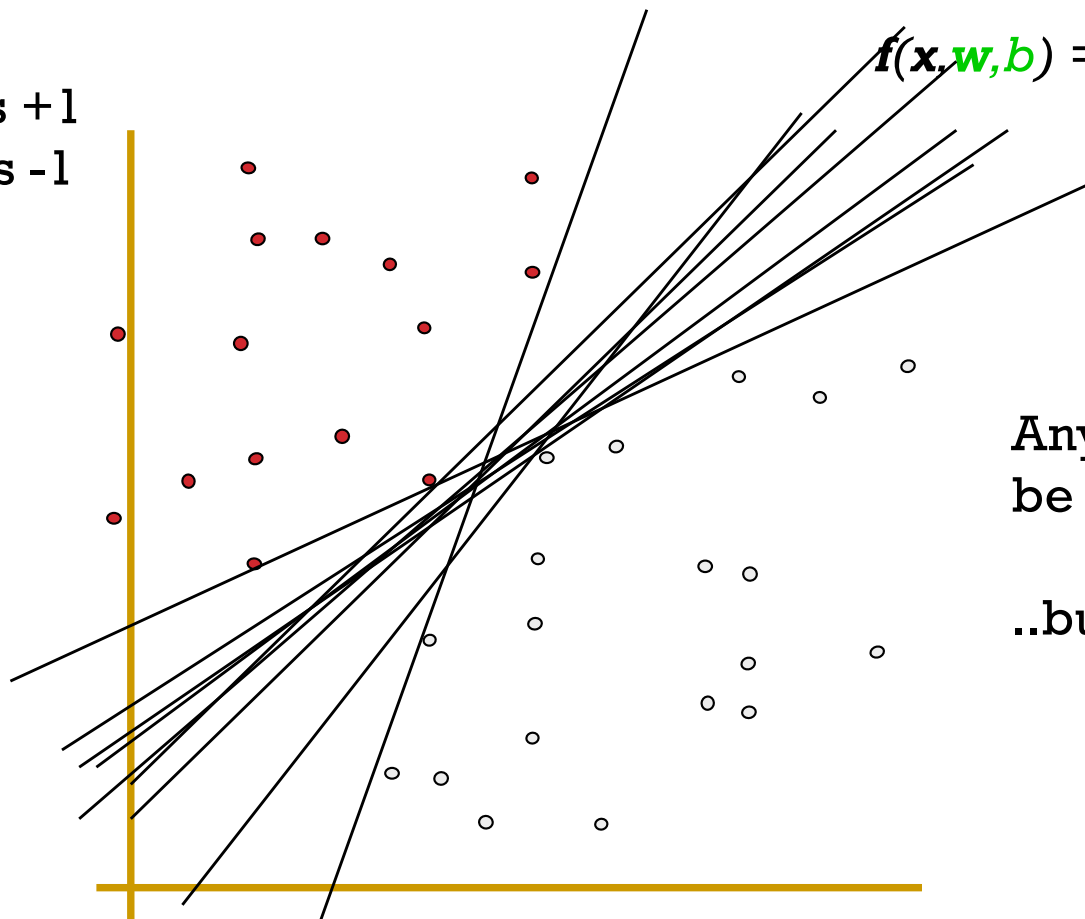
How would you classify this data?

LINEAR CLASSIFIERS

α



- denotes +1
- denotes -1



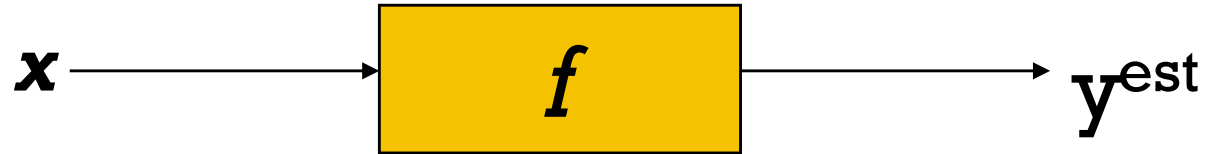
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

Any of these would be fine..

..but which is best?

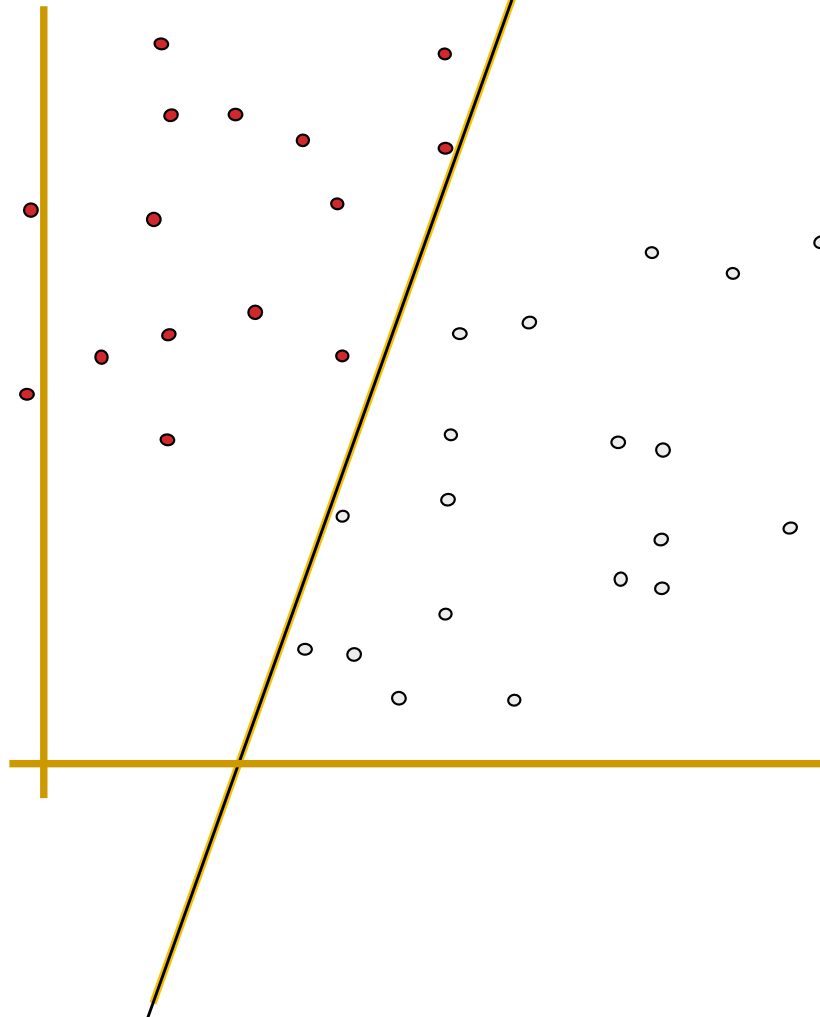
CLASSIFIER MARGIN

α



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1



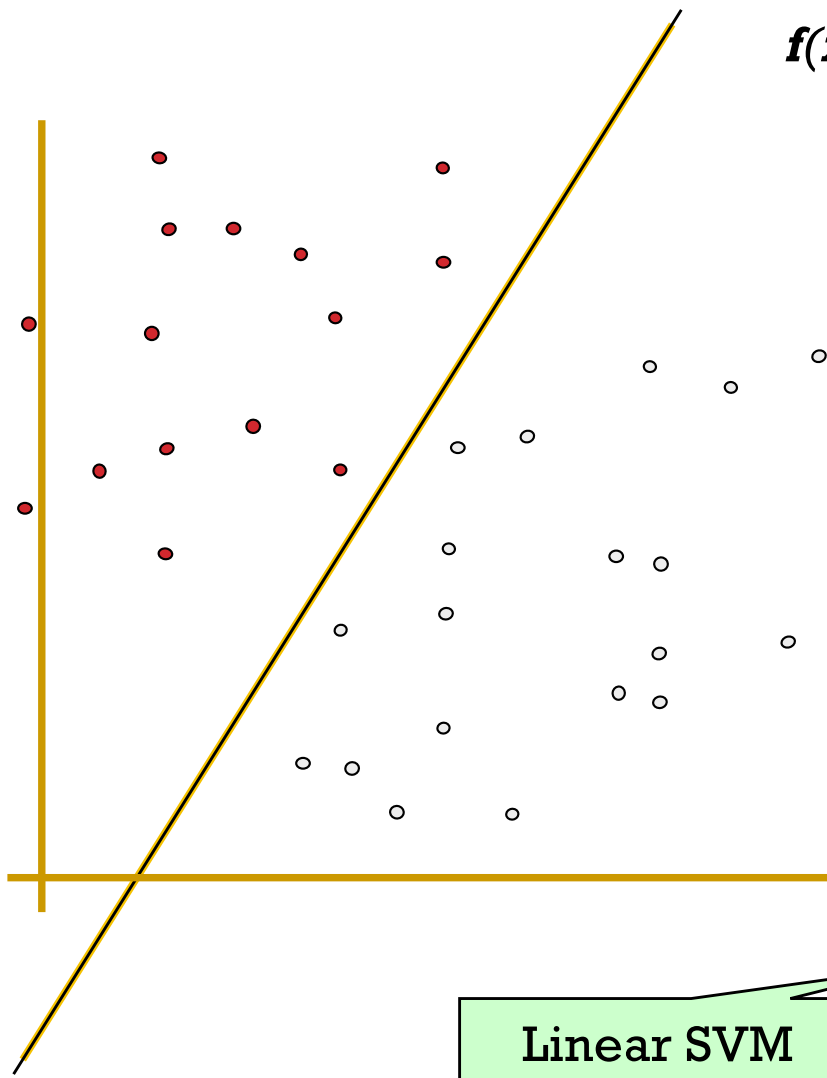
Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

MAXIMUM MARGIN

α



- denotes +1
- denotes -1

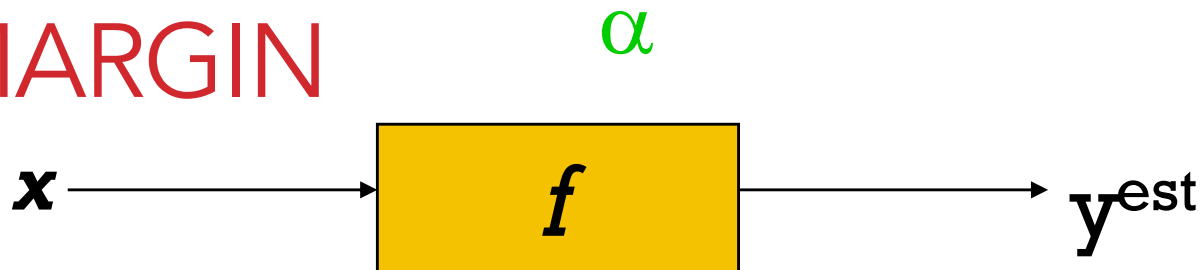


$$f(x, w, b) = \text{sign}(w \cdot x - b)$$

The maximum margin linear classifier is the linear classifier with the, um, maximum margin. This is the simplest kind of SVM (Called an LSVM)

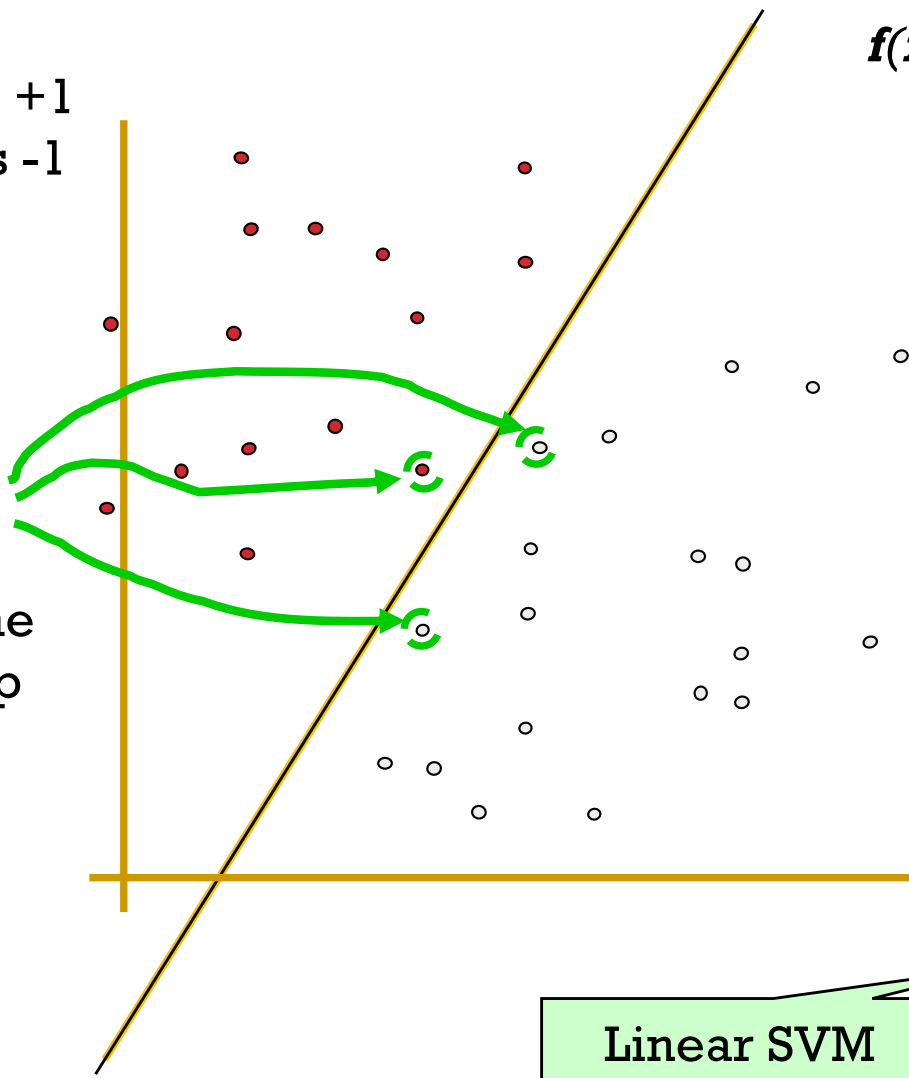
Linear SVM

MAXIMUM MARGIN



- denotes +1
- denotes -1

Support Vectors are those datapoints that the margin pushes up against

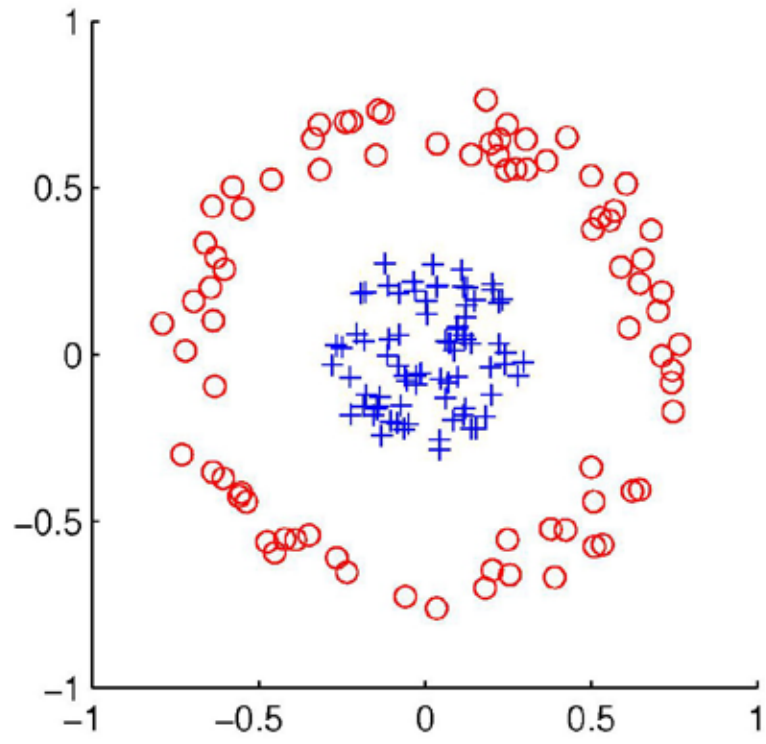


$$f(x, w, b) = \text{sign}(w \cdot x - b)$$

The **maximum margin linear classifier** is the linear classifier with the, um, maximum margin. This is the simplest kind of SVM (Called an **LSVM**)

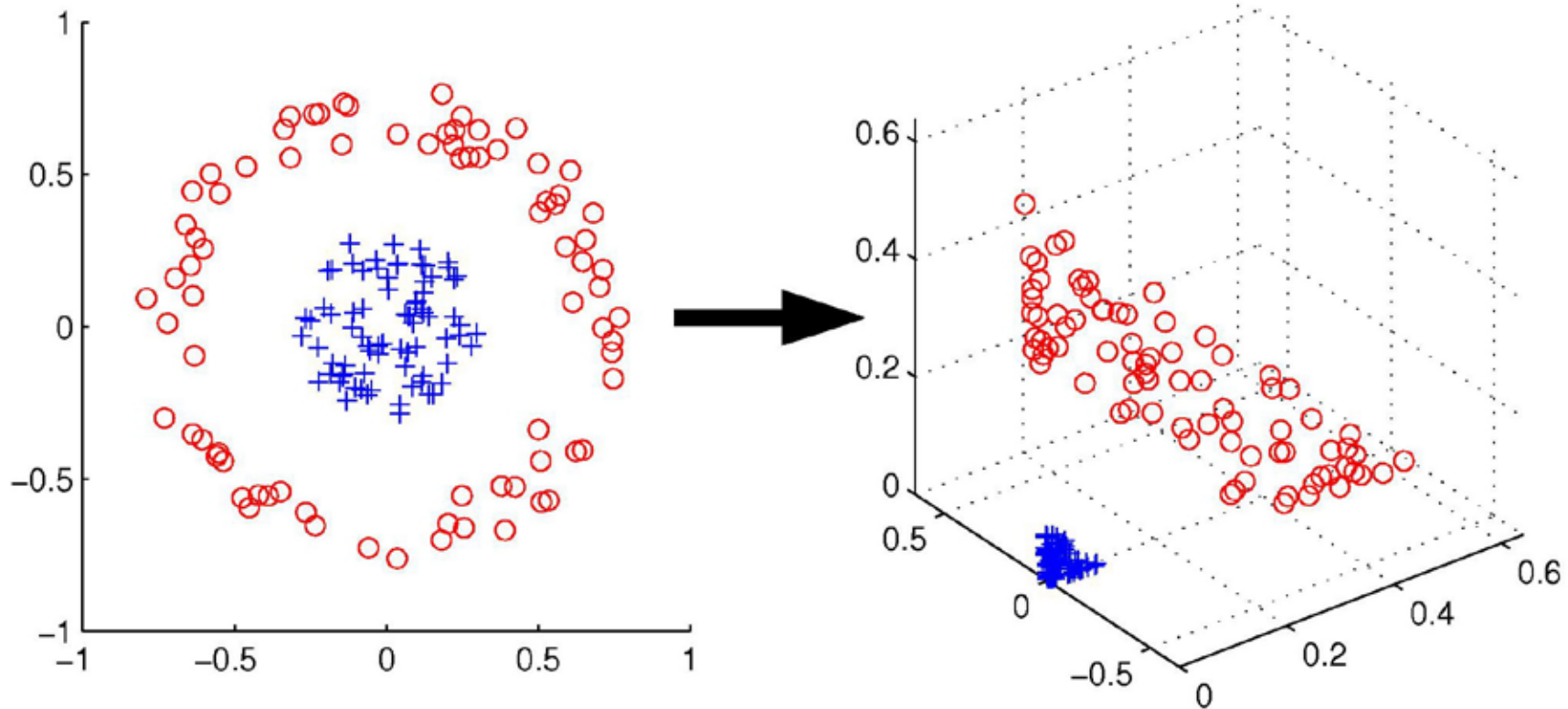
Linear SVM

THE KERNEL TRICK



~

THE KERNEL TRICK



$$\begin{aligned} \phi : \quad \mathcal{R}^2 &\longrightarrow \mathcal{R}^3 \\ (x_1, x_2) &\longmapsto (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2) \end{aligned}$$

[<http://www.cs.berkeley.edu/~jordan/courses/281B-spring04/lectures/lec3.pdf>]

SVM with a polynomial Kernel visualization

Created by:
Udi Aharoni

<https://www.youtube.com/watch?v=3liCbRZPrZA>

MANY CLASSIFIERS TO CHOOSE FROM

K-nearest neighbor

Support Vector Machines

Decision Trees

Random Forrest

(Gradient) Boosted Decision Trees

Logistic Regression

Naïve Bayes

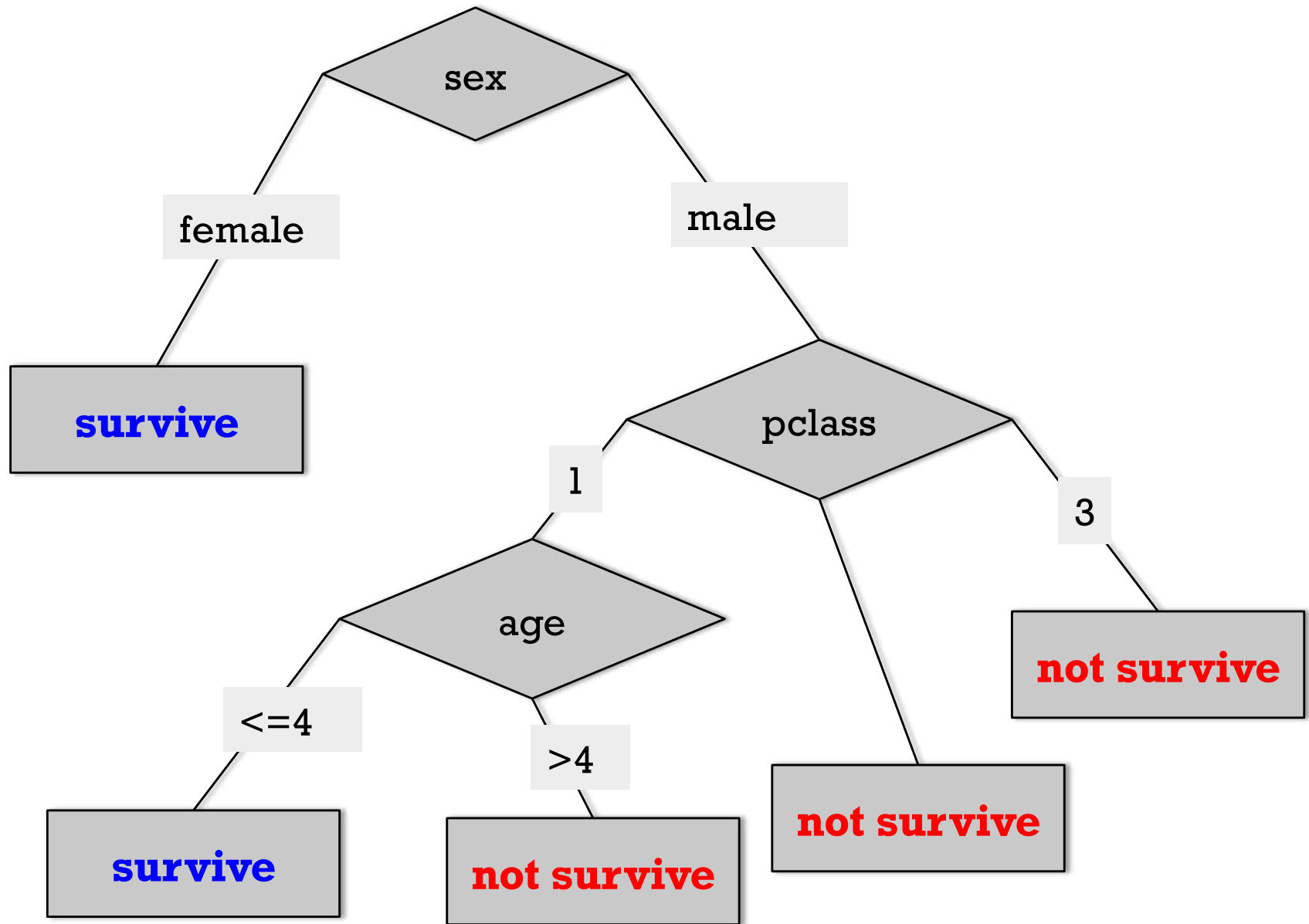
Bayesian network

RBM

....

Which is the best one?

DECISION TREES



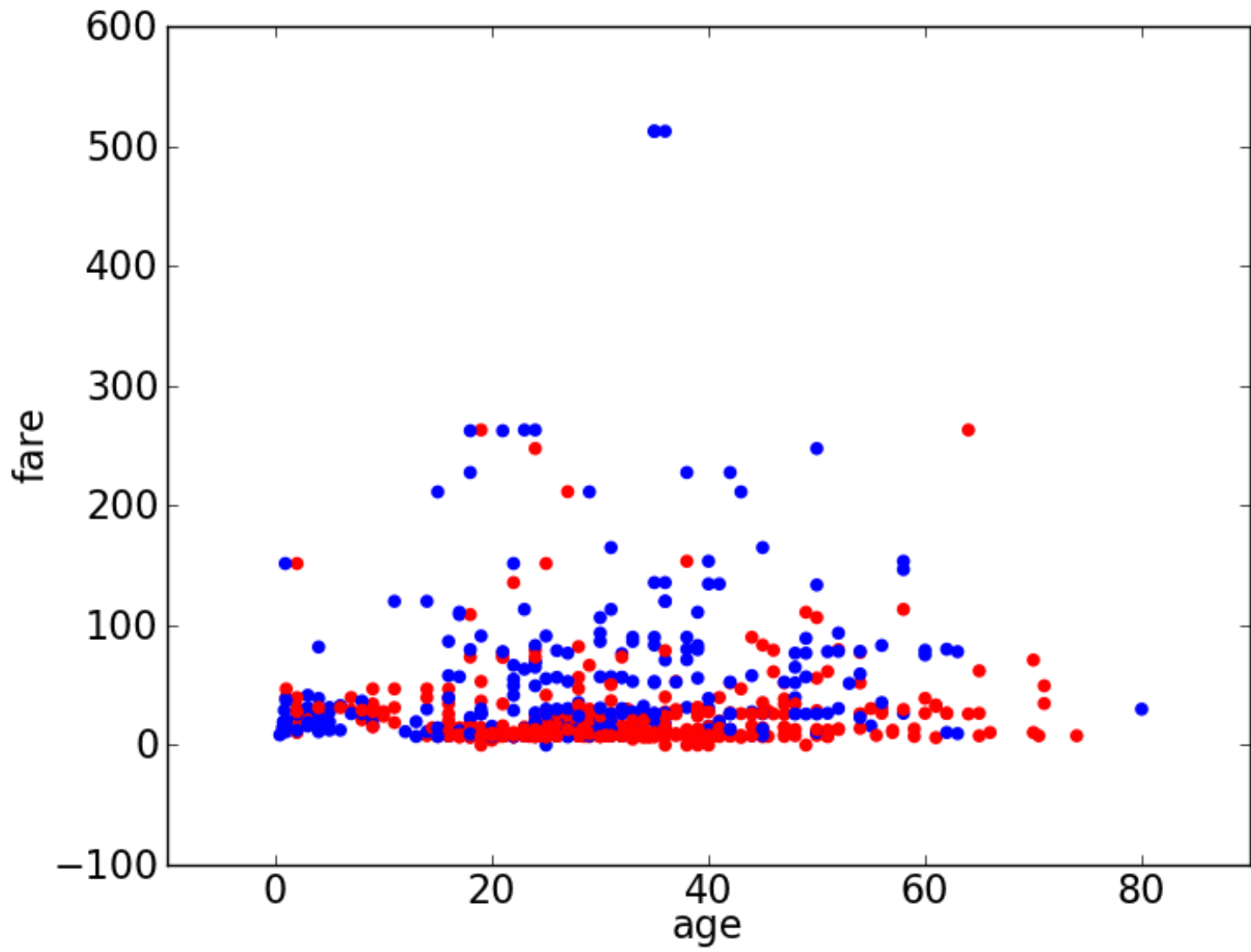
BUILDING A DECISION TREE (ID3 ALGORITHM)

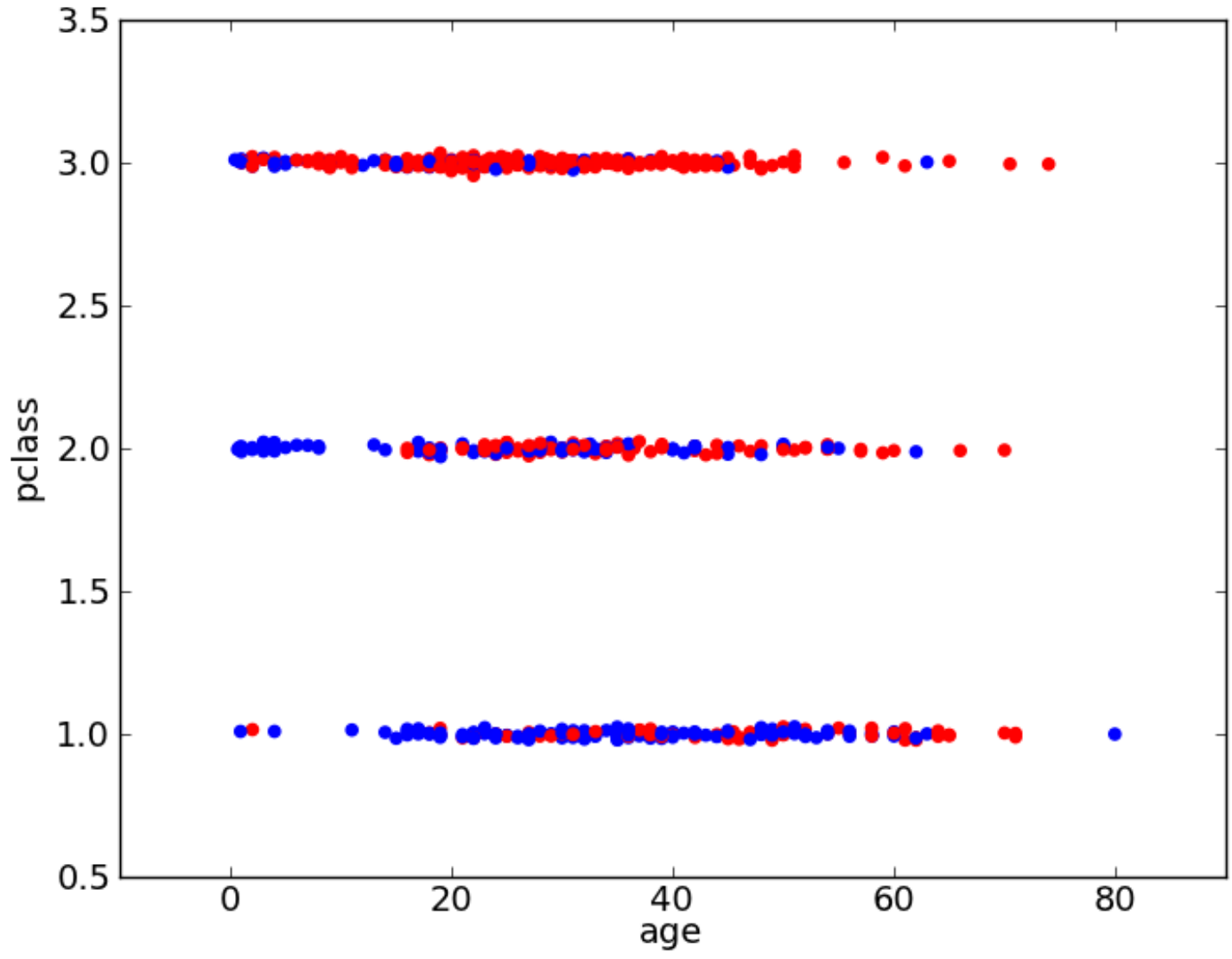
1. Assume attributes are discrete
 - Discretize continuous attributes
2. Choose the attribute with the highest Information Gain
3. Create branches for each value of attribute
4. Examples partitioned based on selected attributes
5. Repeat with remaining attributes
6. Stopping conditions
 - All examples assigned the same label
 - No examples left

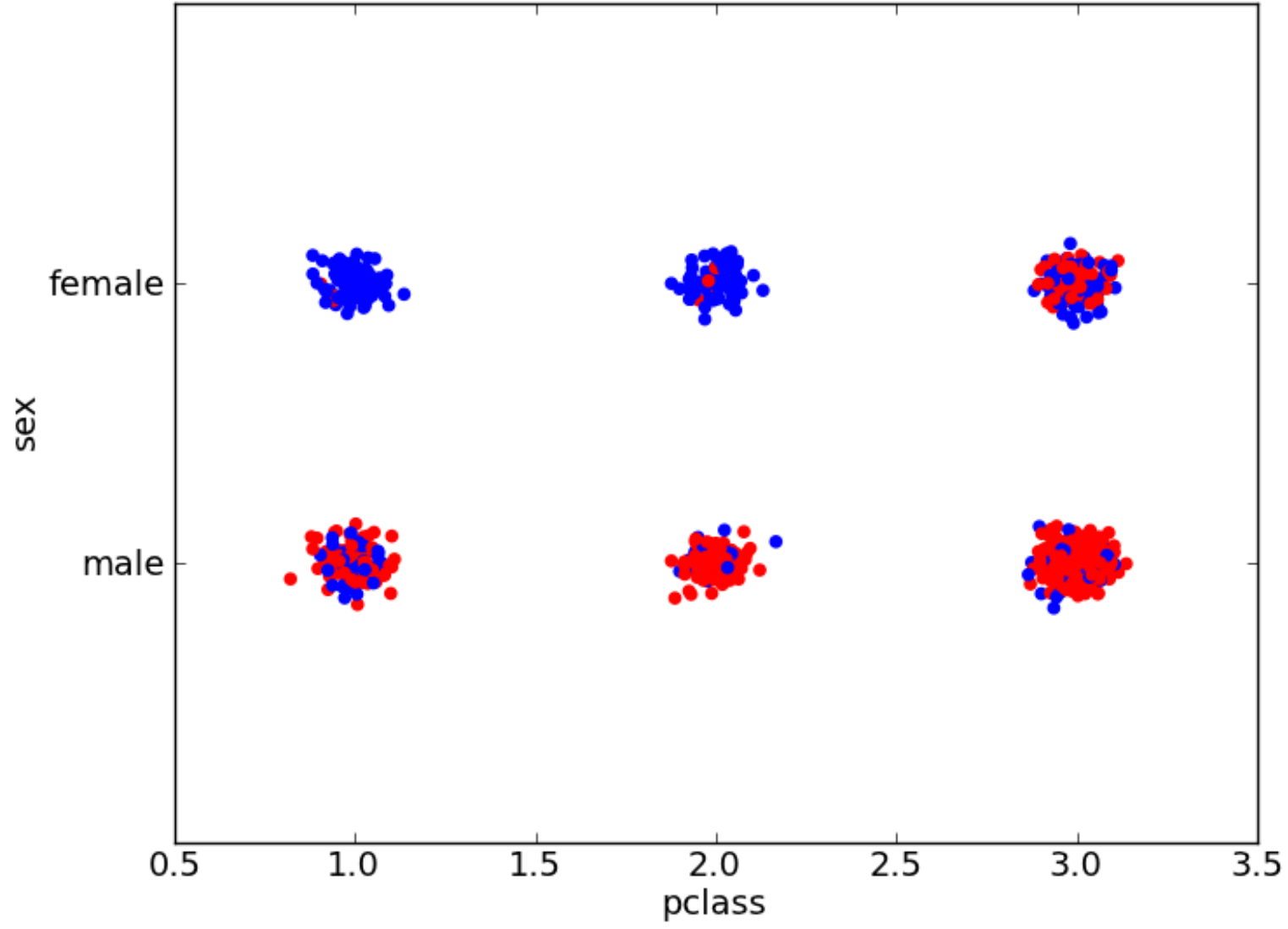
HIGHEST INFORMATION GAIN

Titanic Data

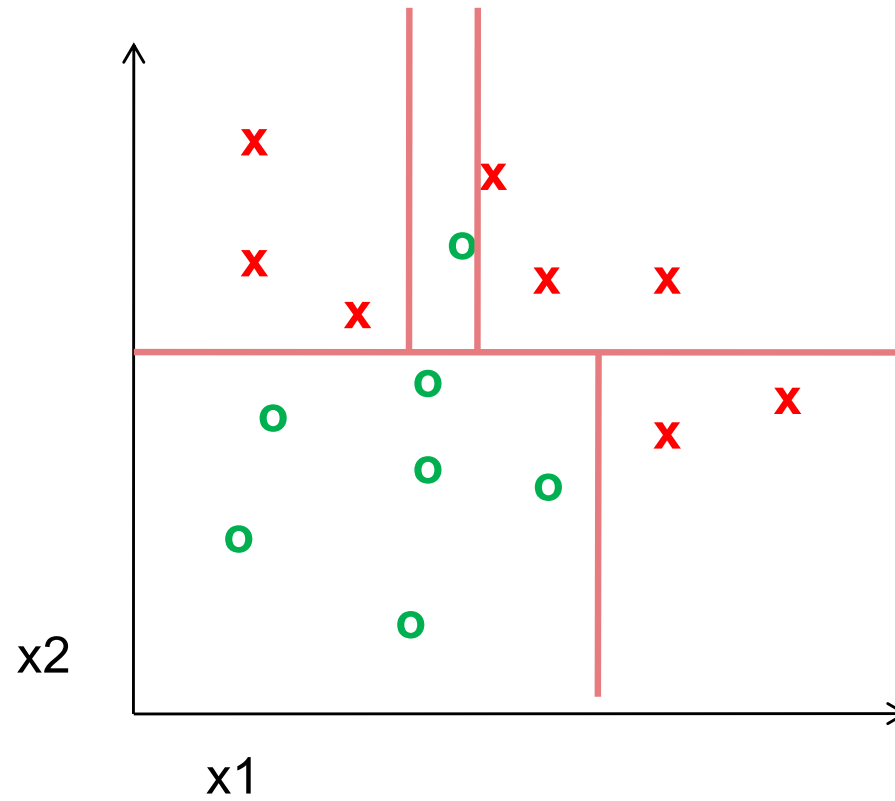
survived	pclass	sex	age	sibsp	parch	fare	cabin	embarked
0	3	male	22	1	0	7.25		S
1	1	female	38	1	0	71.2833	C85	C
1	3	female	26	0	0	7.925		S
1	1	female	35	1	0	53.1	C123	S
0	3	male	35	0	0	8.05		S
0	3	male		0	0	8.4583		Q
0	1	male	54	0	0	51.8625	E46	S
0	3	male	2	3	1	21.075		S
1	3	female	27	0	2	11.1333		S
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 BOUNDARIES: DECISION TREES



MANY CLASSIFIERS TO CHOOSE FROM

K-nearest neighbor

Support Vector Machines

Decision Trees

Random Forrest

(Gradient) Boosted Decision Trees

Logistic Regression

Naïve Bayes

Bayesian network

RBM

....

Which is the best one?

Ensemble Methods

ENSEMBLE METHODS

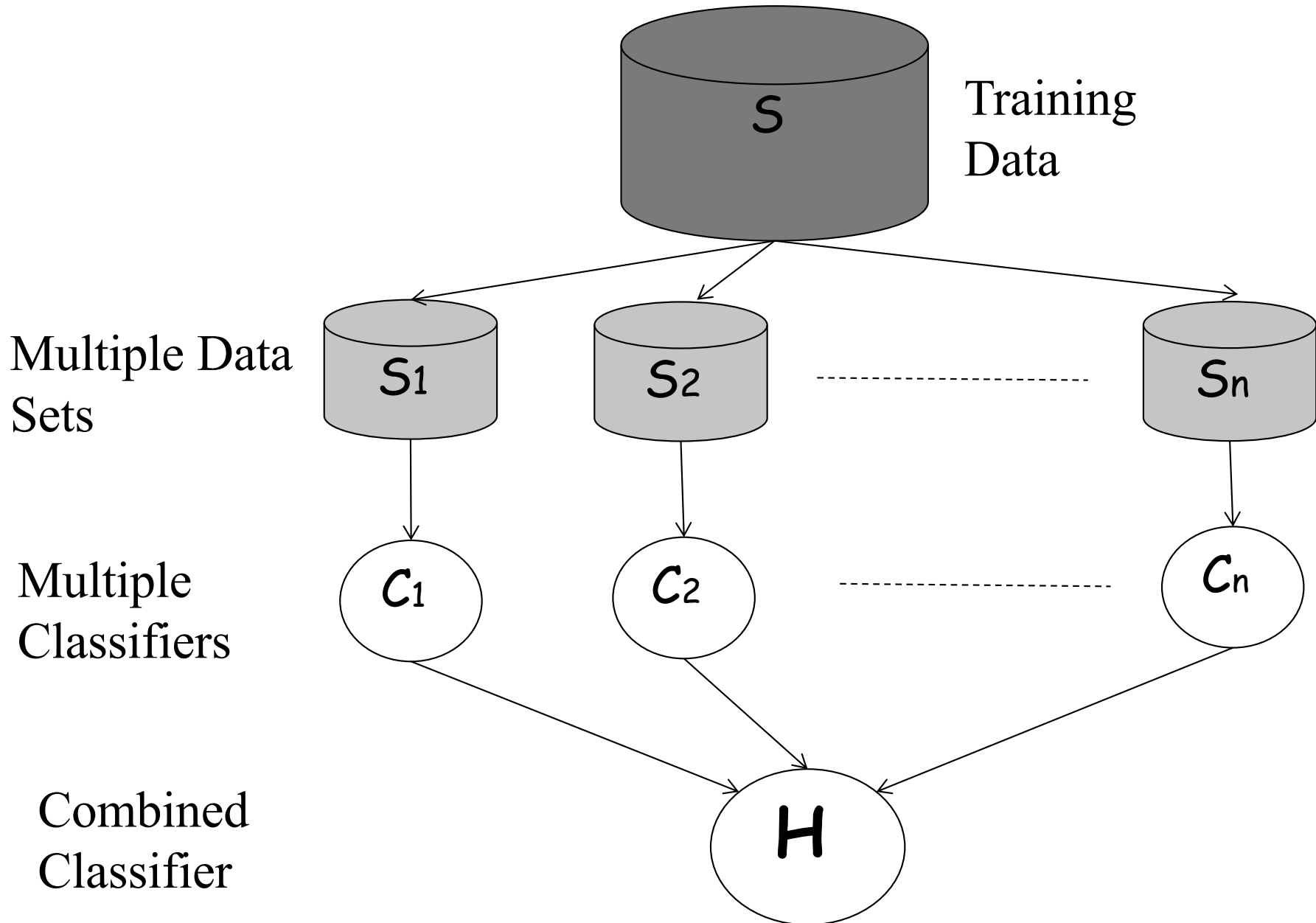
Bagging (Breiman 1994,...)

Random forests (Breiman 2001,...)

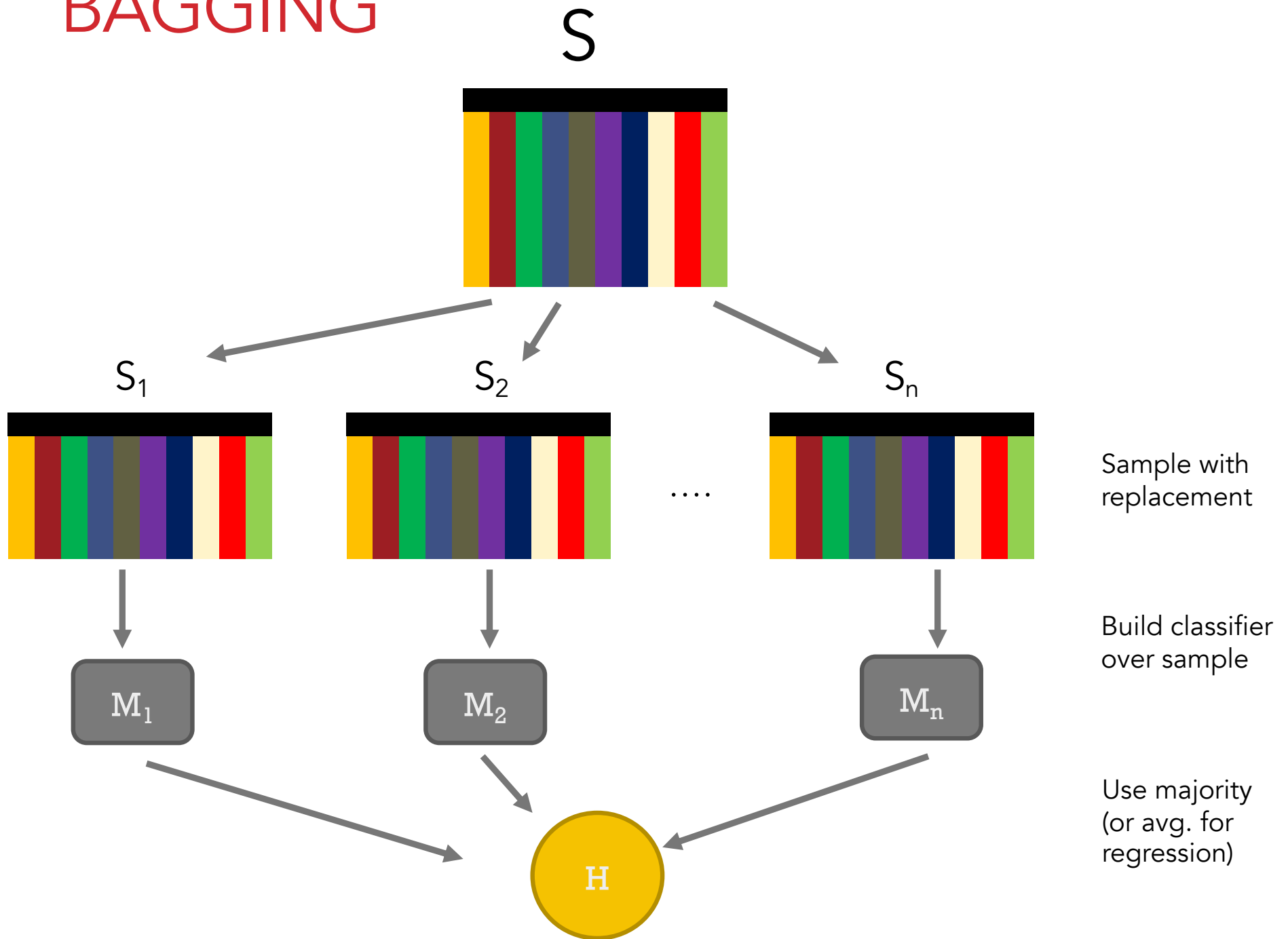
Boosting (Freund and Schapire 1995, Friedman et al. 1998,...)

Predict class label for unseen data by aggregating a set of predictions (classifiers learned from the training data).

GENERAL IDEA



BAGGING



BAGGING

- Can help a lot if data is noisy.
- Bagging works because it reduces variance by voting/averaging
 - In some pathological hypothetical situations the overall error might increase
 - Usually, the more classifiers the better
- **If the learning algorithm is unstable, then Bagging almost always improves performance**
 - Learning algorithm is unstable: if small changes to the training set cause large changes in the learned classifier.
 - Some candidates: Decision tree, decision stump, regression tree, linear regression, SVMs

ENSEMBLE METHODS

Bagging (Breiman 1994,...)

Random forests (Breiman 2001,...)

Boosting (Freund and Schapire 1995, Friedman et al. 1998,...)

THE RANDOM FORESTS ALGORITHM

Given a training set S

For $i = 1$ to k do:

Build subset S_i by sampling with replacement from S

Learn tree T_i from S_i

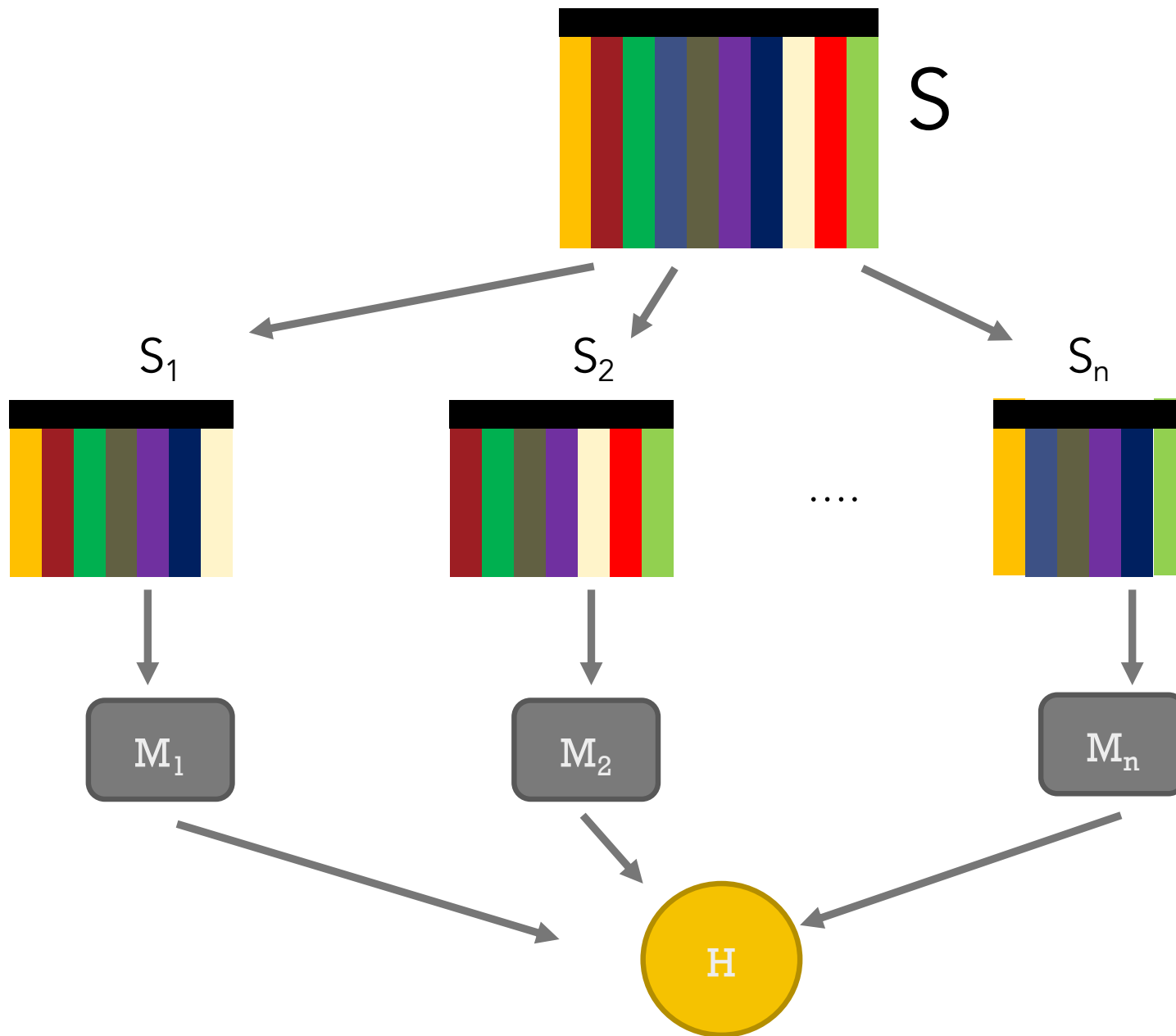
At each node:

Choose best split from random subset of F features

Each tree grows to the largest extent, and no pruning

Make predictions according to majority vote of the set of k trees.

RANDOM FORREST



Sample with Replacement & select random subset of features*

Build classifier over sample

Use majority Vote for classification (or avg. for regression)

* Normally done for each node of the decision tree – not once

ENSEMBLE METHODS

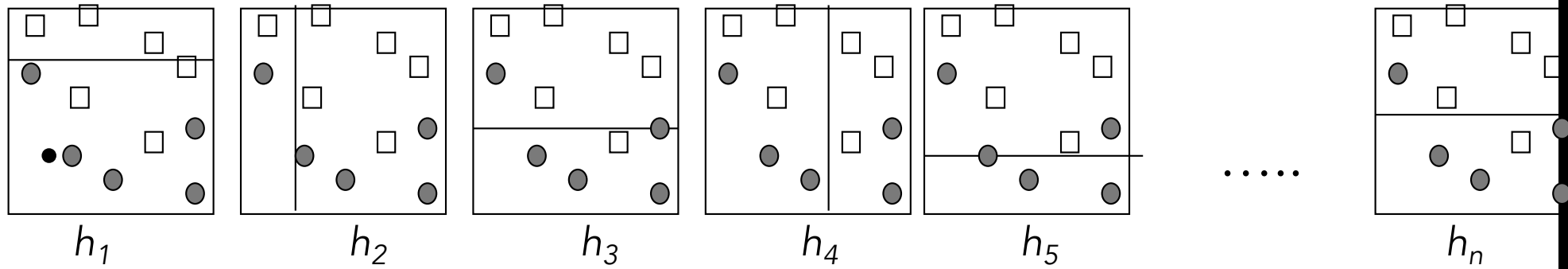
Bagging (Breiman 1994,...)

Random forests (Breiman 2001,...)

Boosting (Freund and Schapire 1995, Friedman et al. 1998,...)

ADABOOST - CORE IDEA

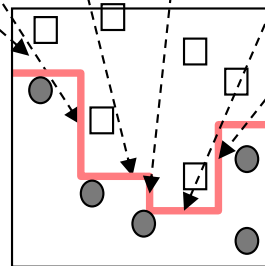
Take a set of weak classifiers (normally they should do better than guessing)



Classification
Result

θ_1 θ_2 θ_3 θ_4 θ_5 θ_n

Weight the result of each classify
with θ

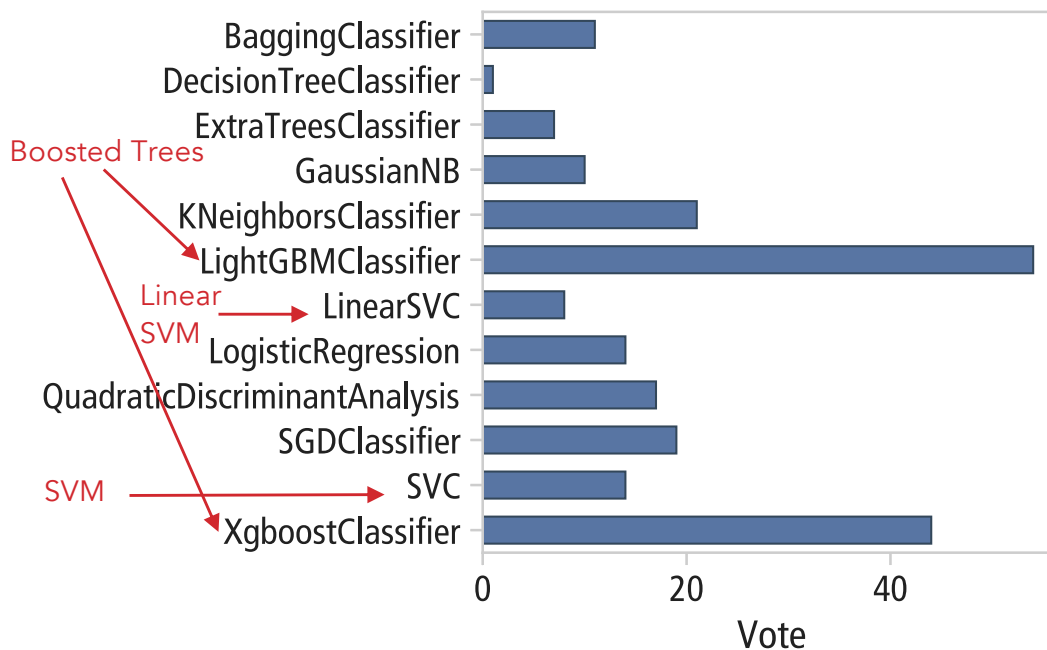


Combine to form the
Final strong classifier

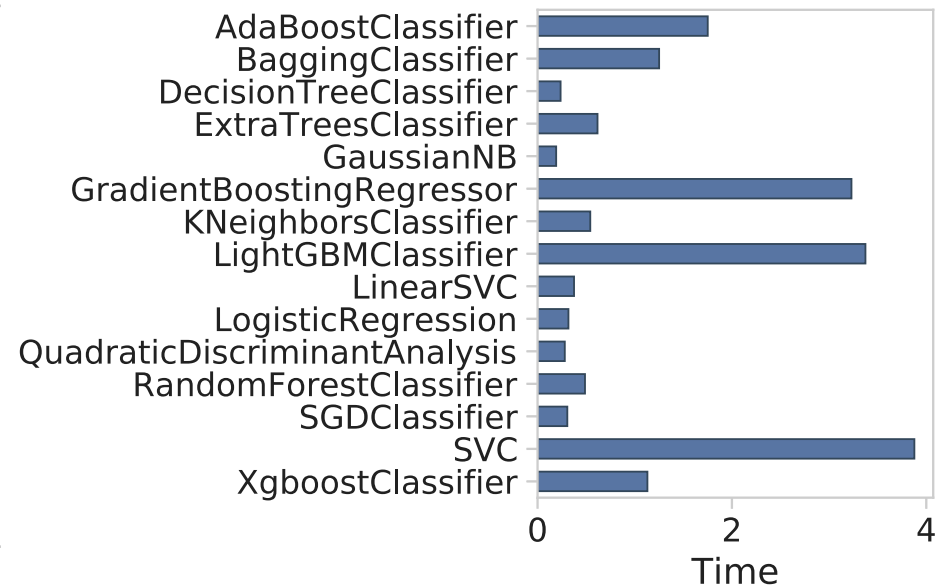
$$H(x) = \text{sign} \left(\sum_{i=1}^n \theta_i h_i(x) \right)$$

PERFORMANCE OF DIFFERENT ML MODEL FAMILIES

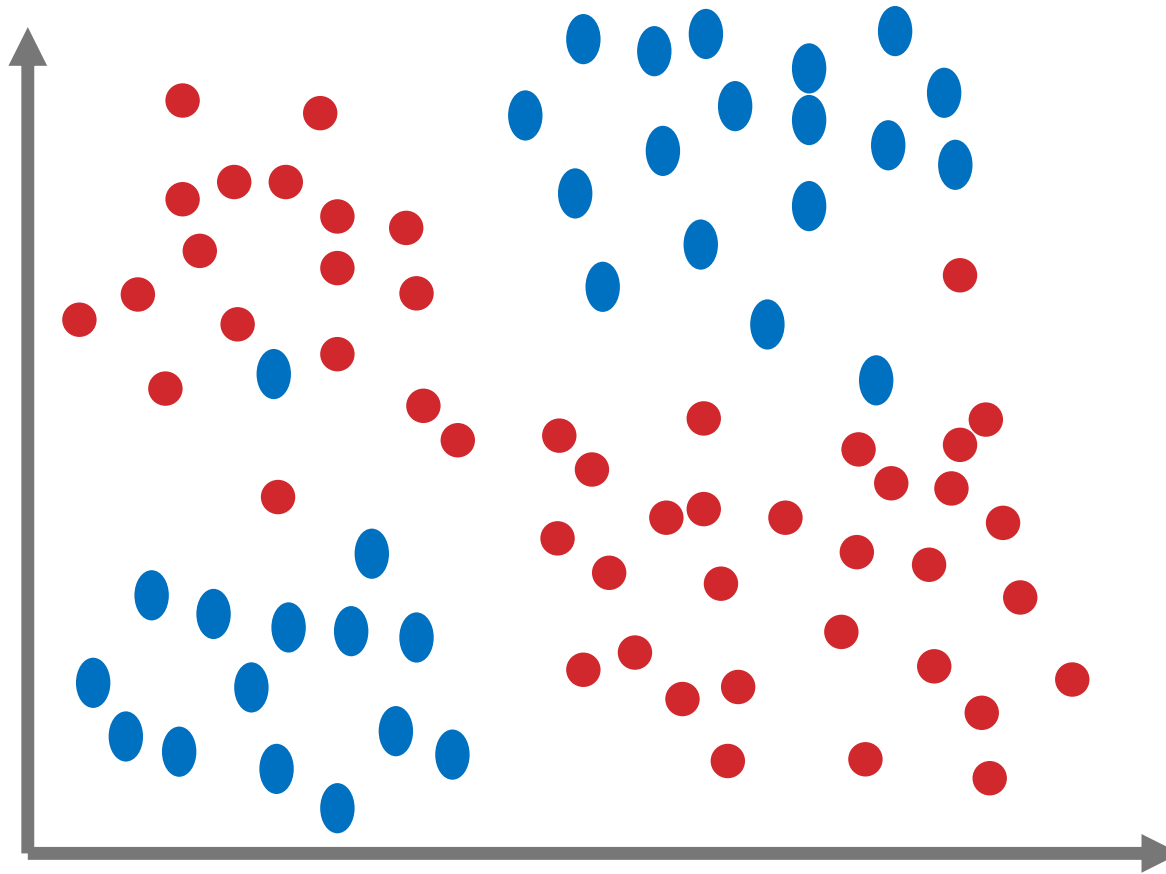
How often ranked 1st



Relative Training Time



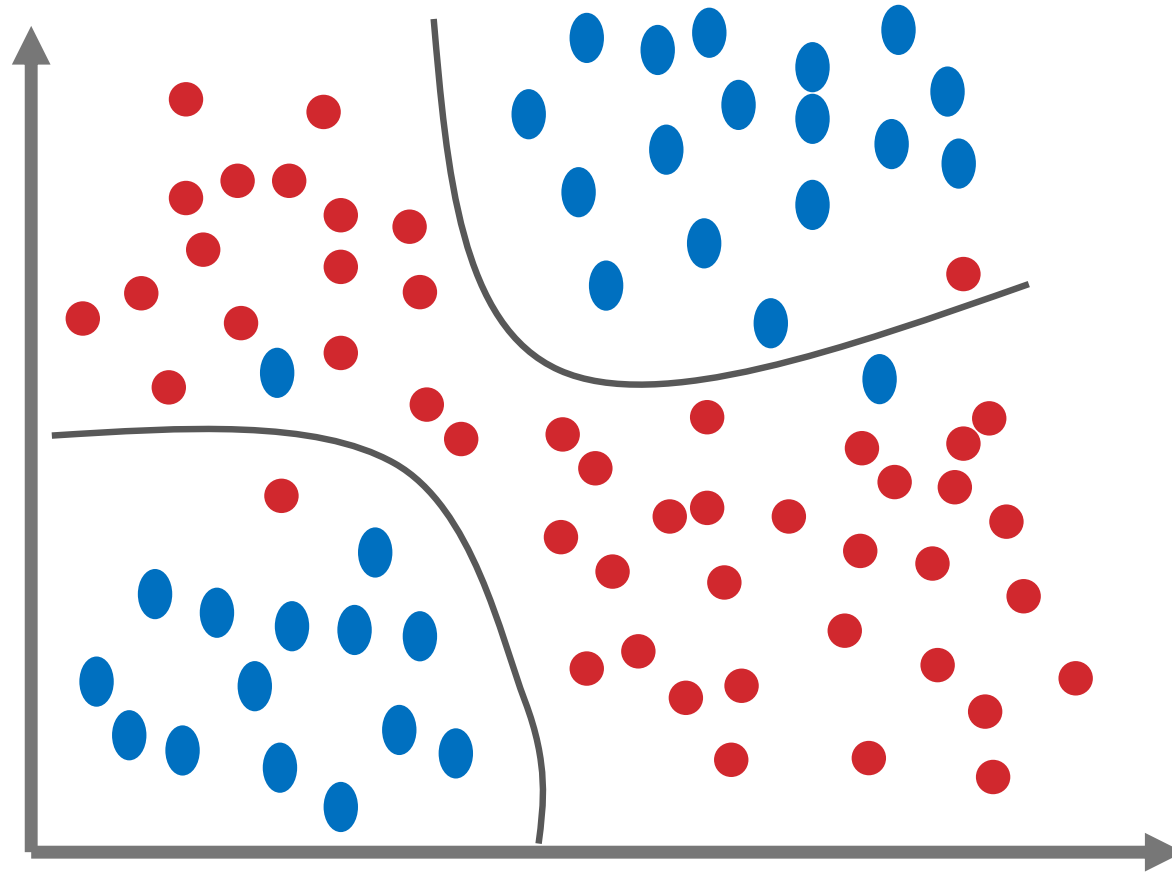
IN-CLASS TASK



How would you draw the expected decision boundary for

- Random Forreest
- SVM w/ kernel and regularization
- 1-KNN

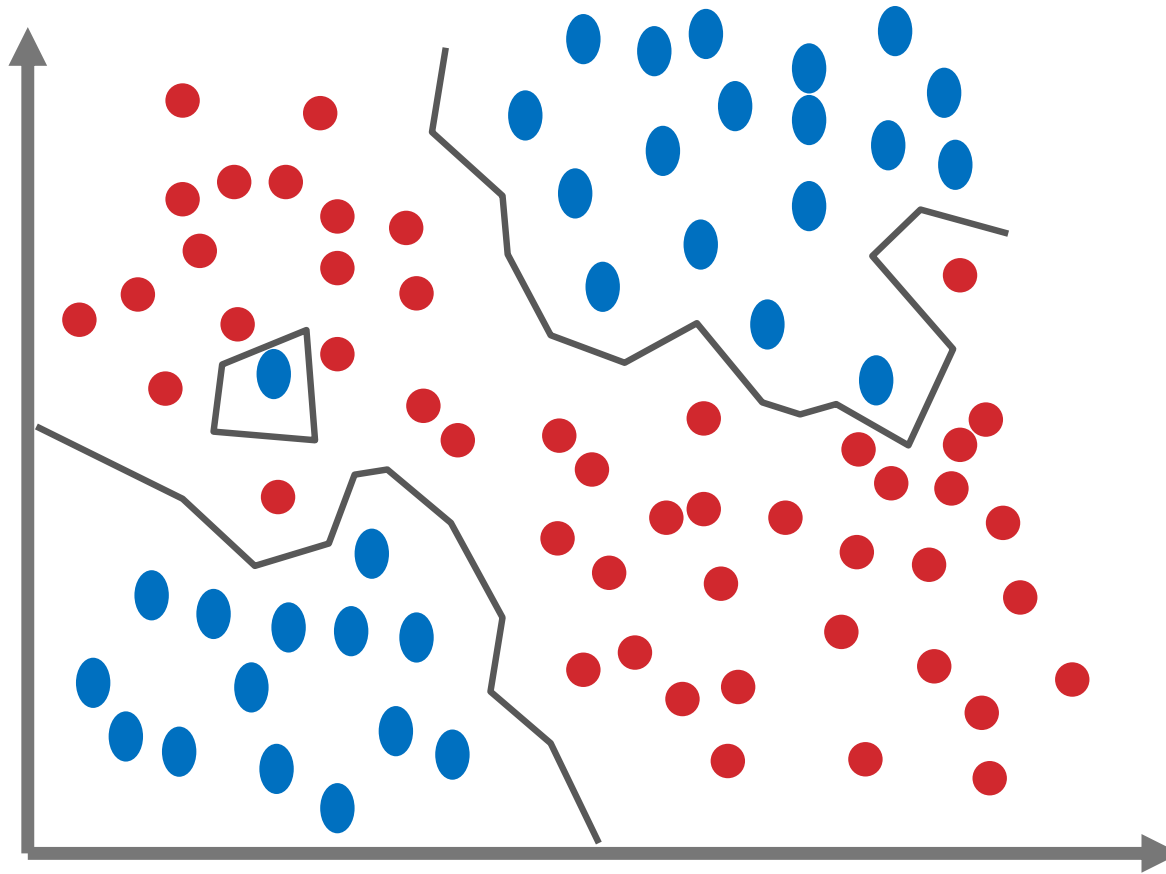
CLICKER



The decision boundary looks like the one of:

- a) Random Forrest
- b) SVM w/ kernel and regularization
- c) 1-KNN

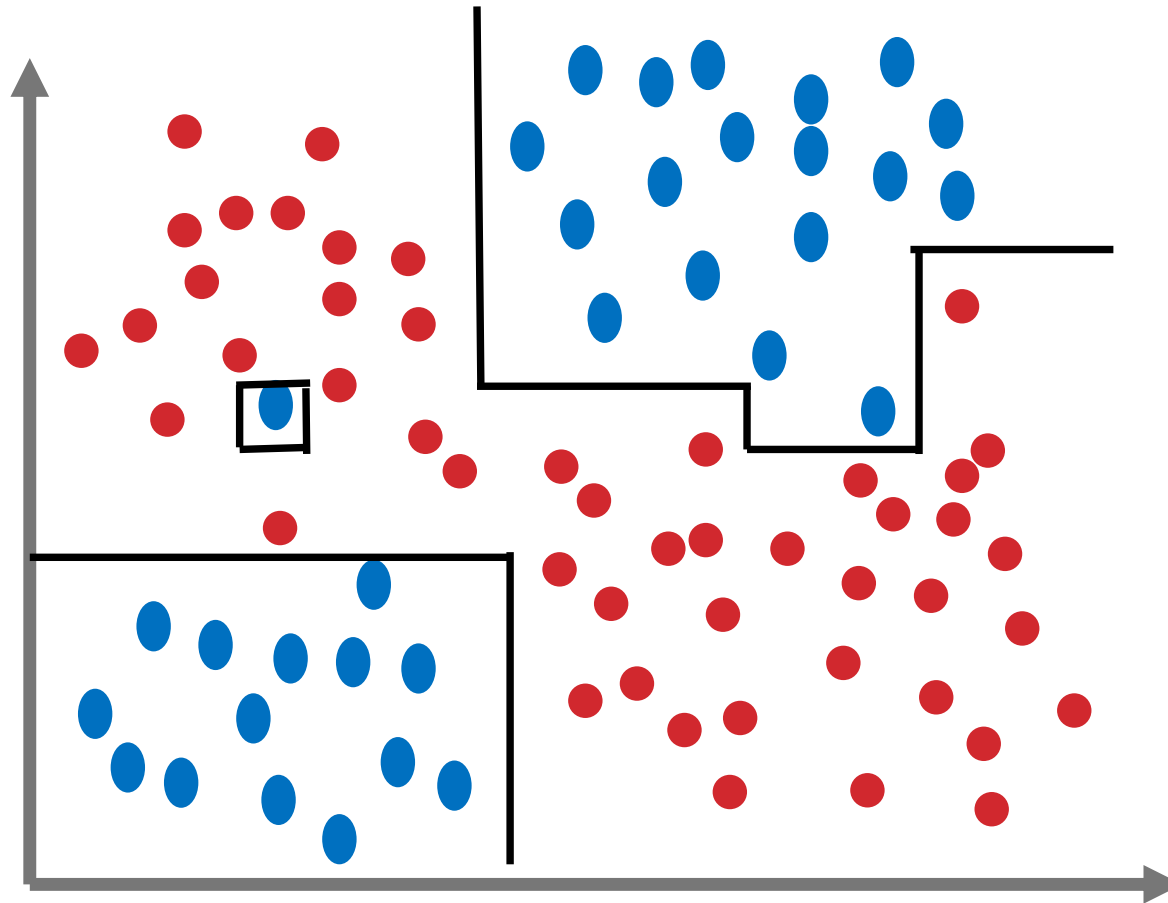
CLICKER



The decision boundary looks like the one of:

- a) Random Forrest
- b) SVM w/ kernel and regularization
- c) 1-KNN

RANDOM FORREST



The decision boundary looks like the one of:

- a) Random Forrest
- b) SVM w/ kernel and regularization
- c) 1-KNN

Machine Learning

*Nightmare
SERIES*

What if your model has a high error?

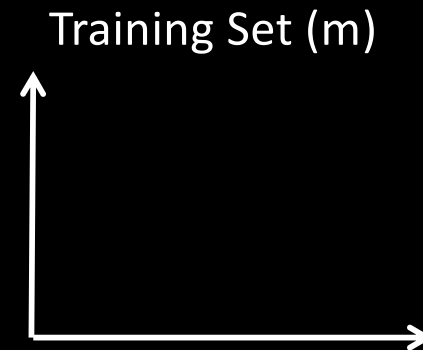
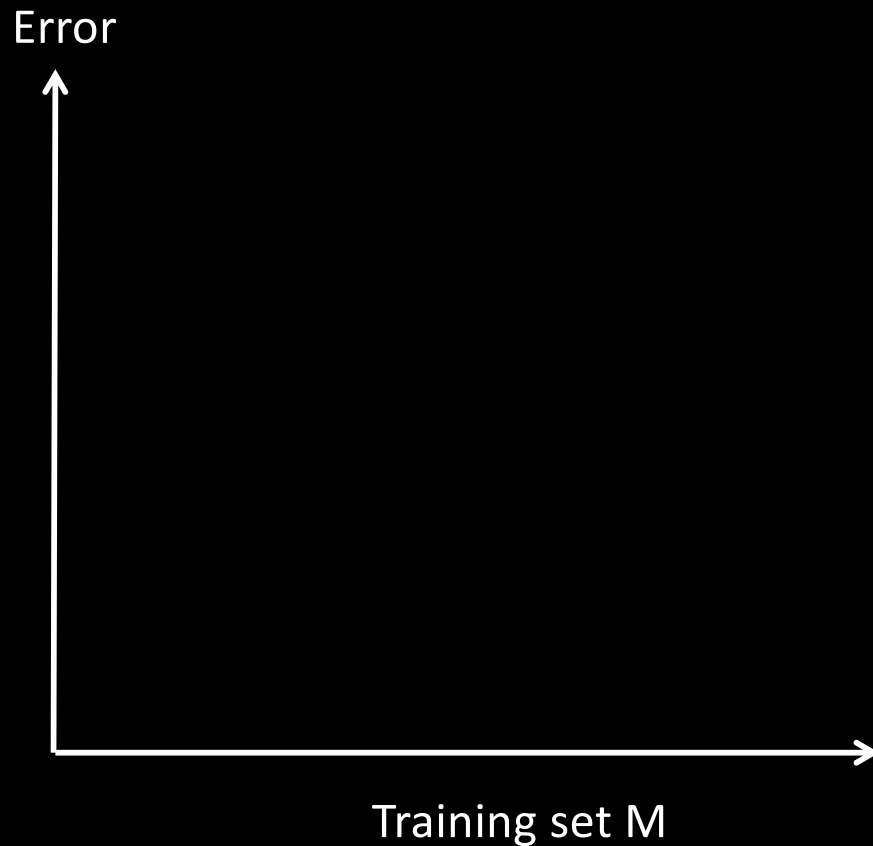
- Try getting more training examples
- Try smaller sets of features
- Try getting additional features
- Try creating features from existing features (kernels)
- Try decrease regularization
- Try increase regularization

Bias and Variance

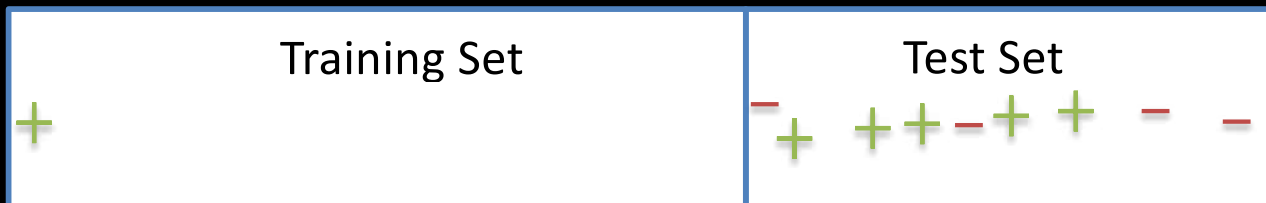
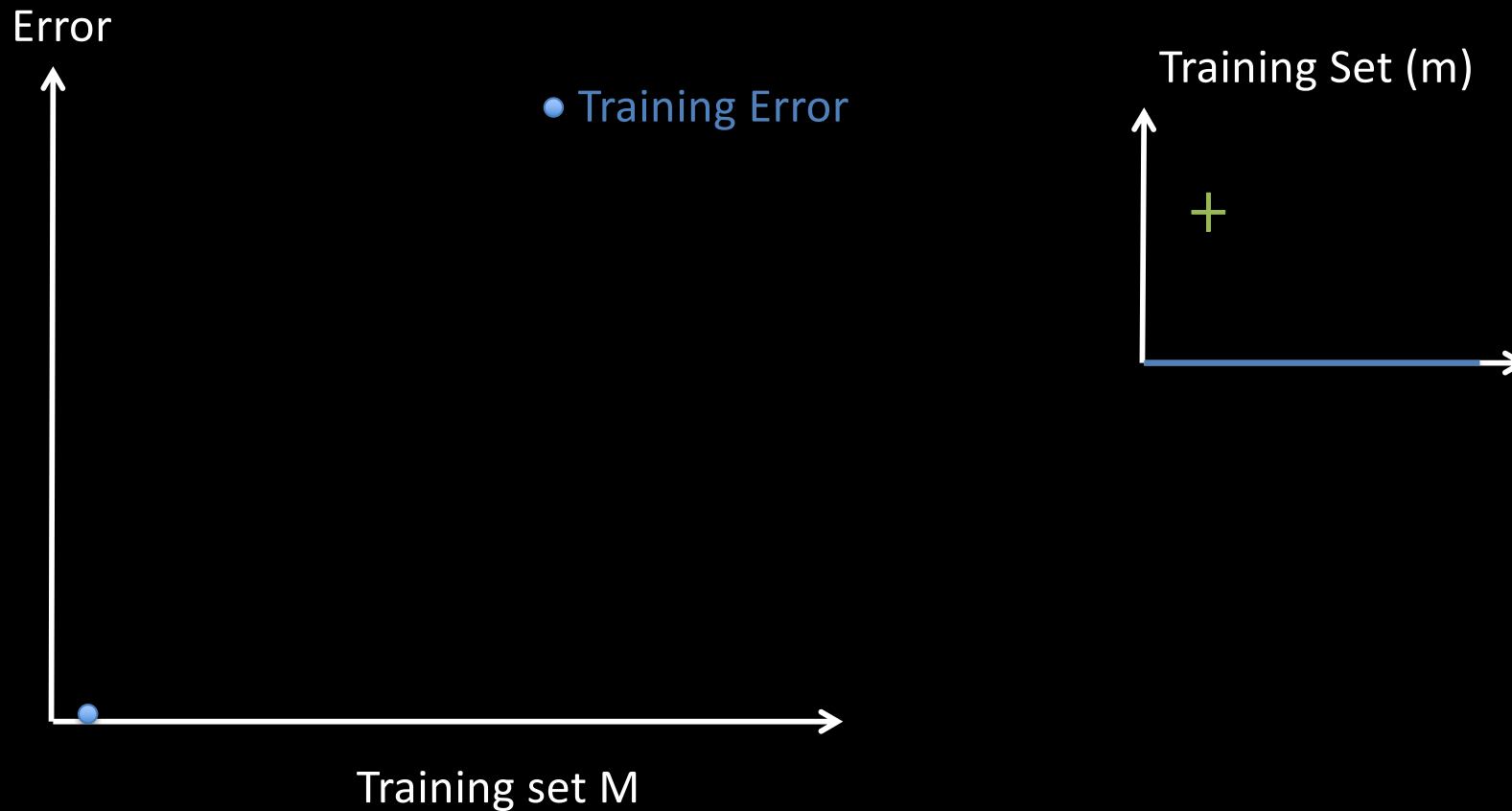


++ - + - ++ - - + - + - - + ++ - + + - -

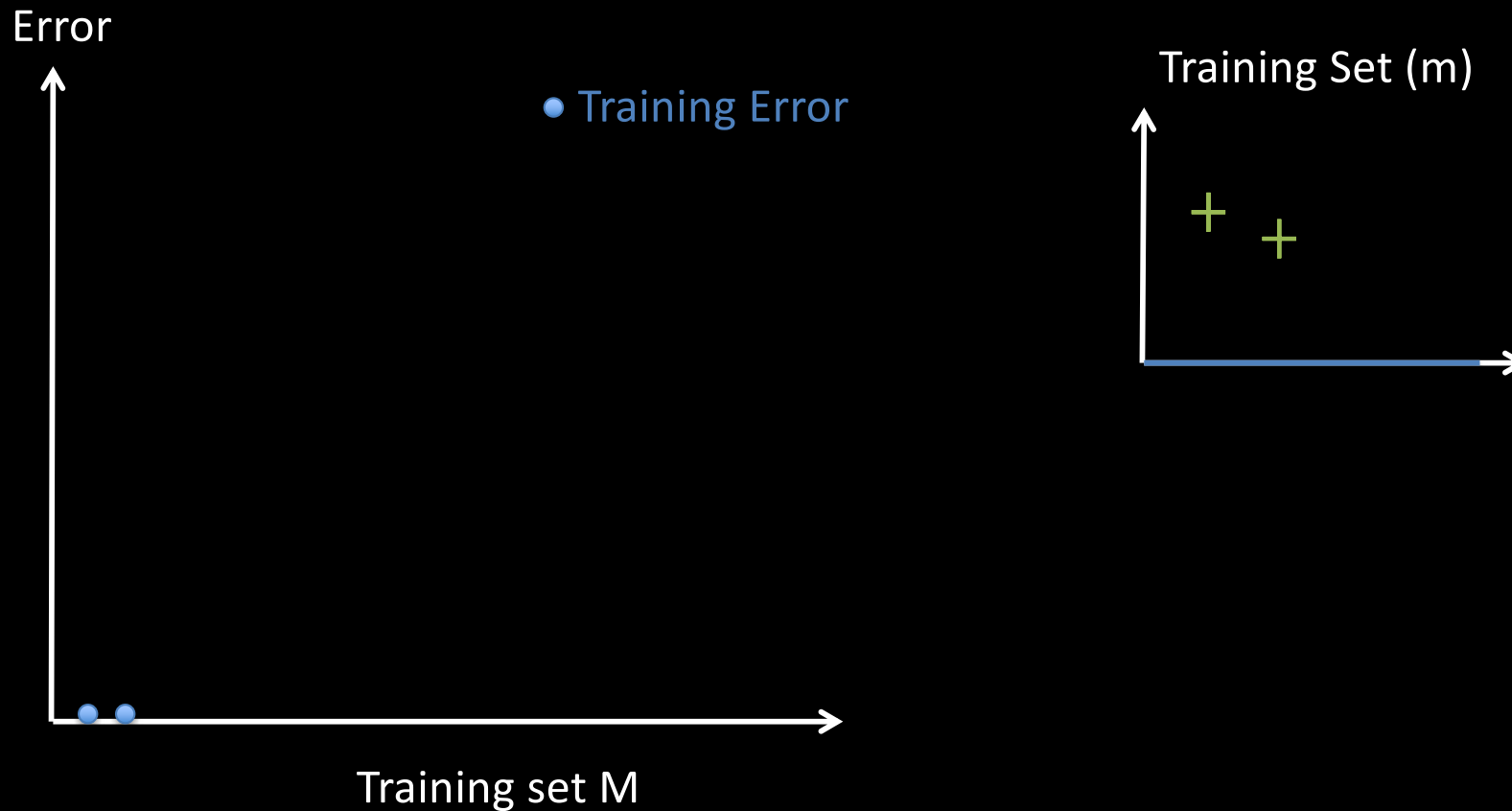
Bias and Variance



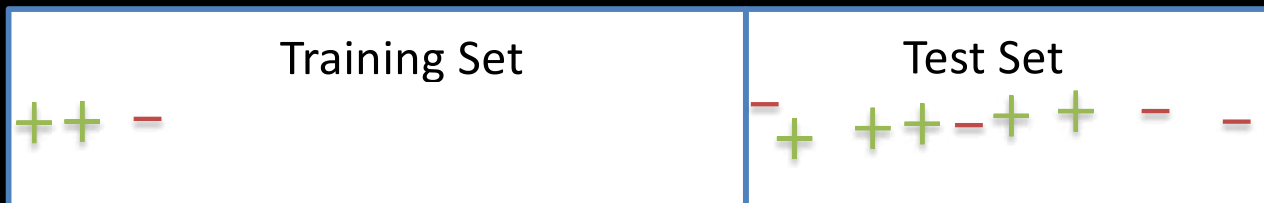
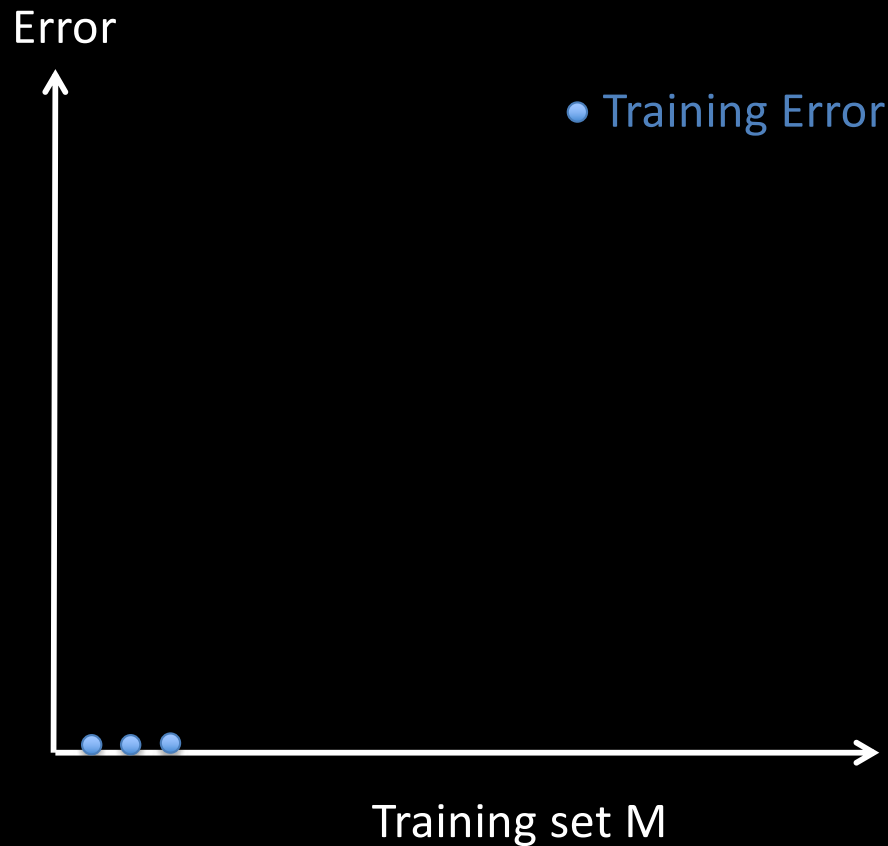
Bias and Variance



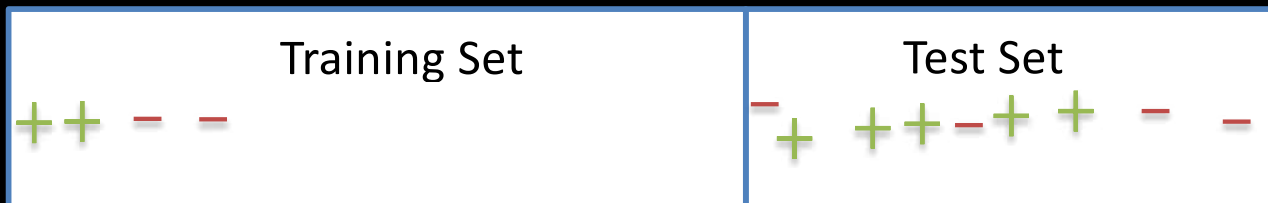
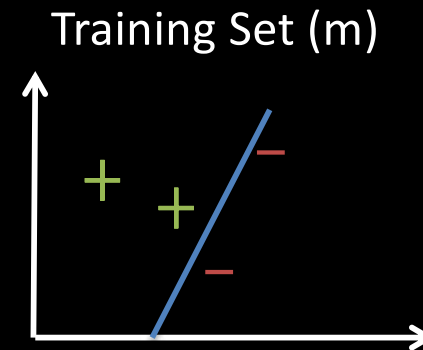
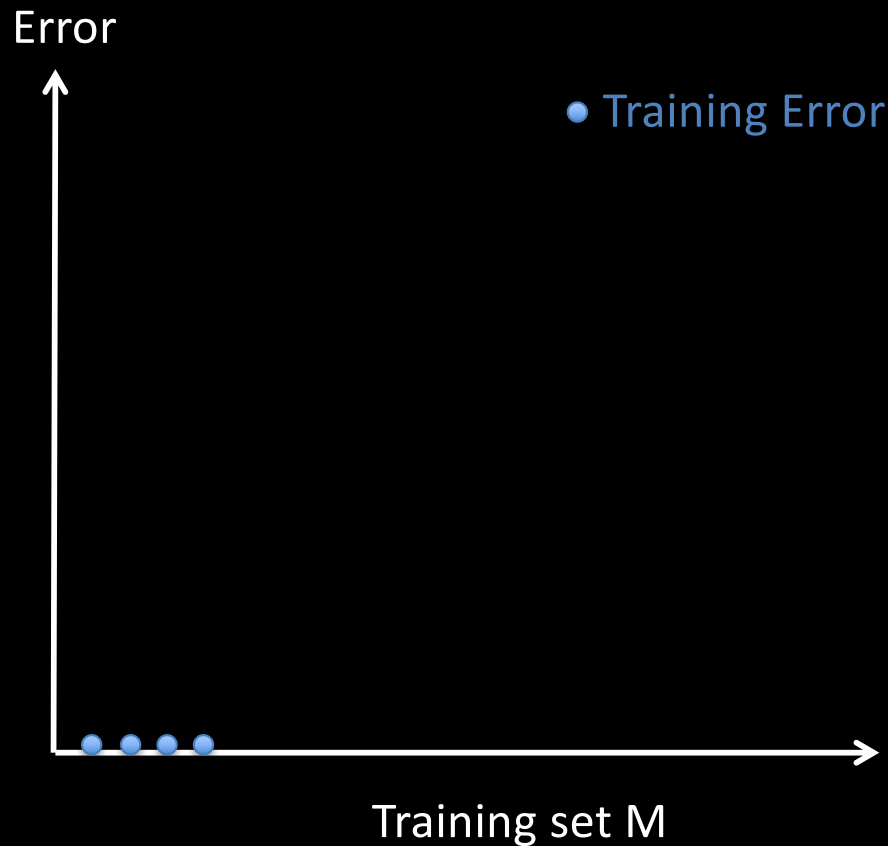
Bias and Variance



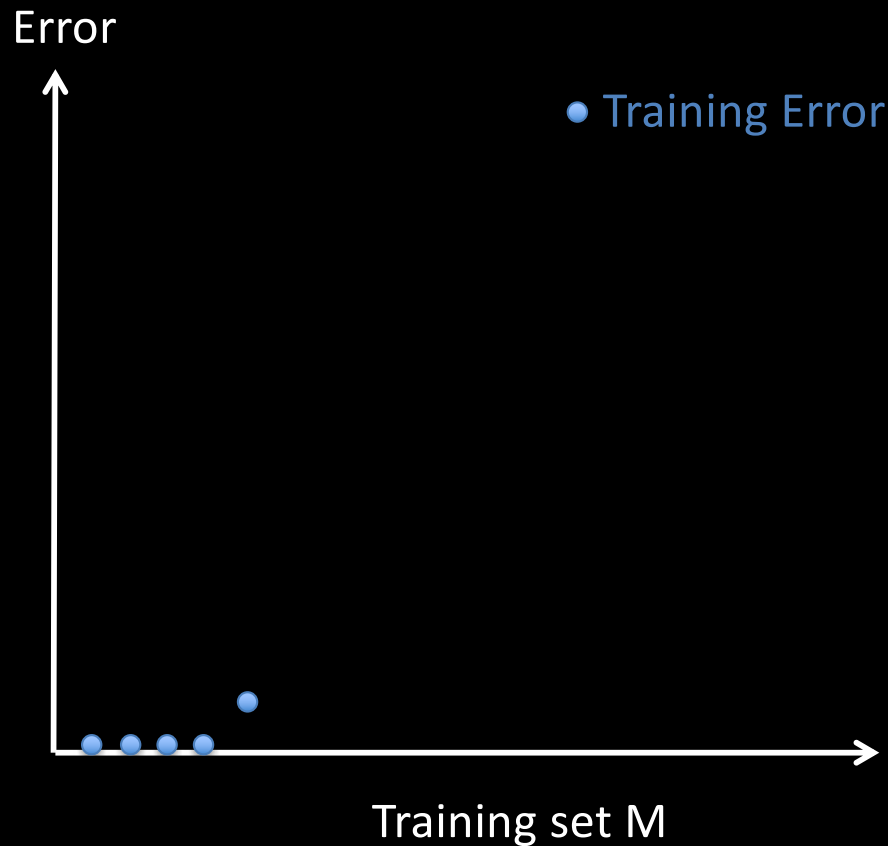
Bias and Variance



Bias and Variance



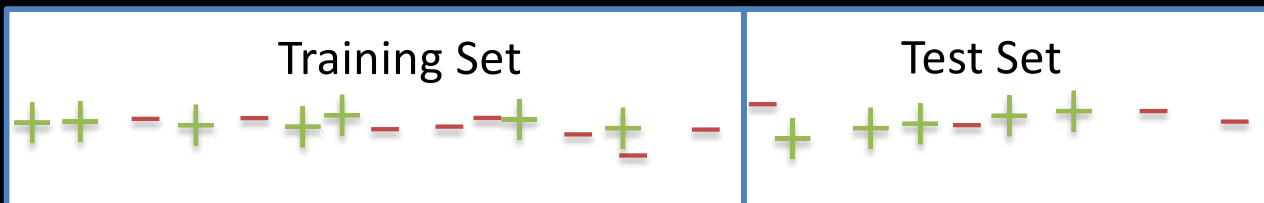
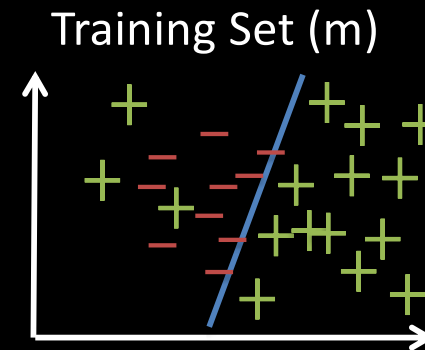
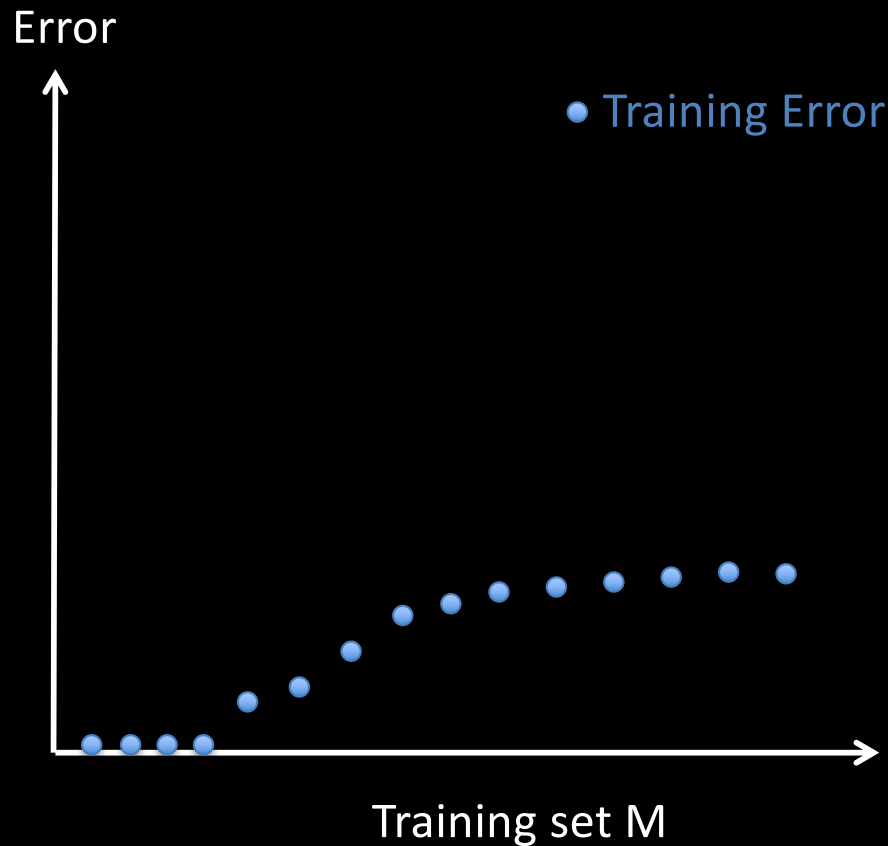
Bias and Variance



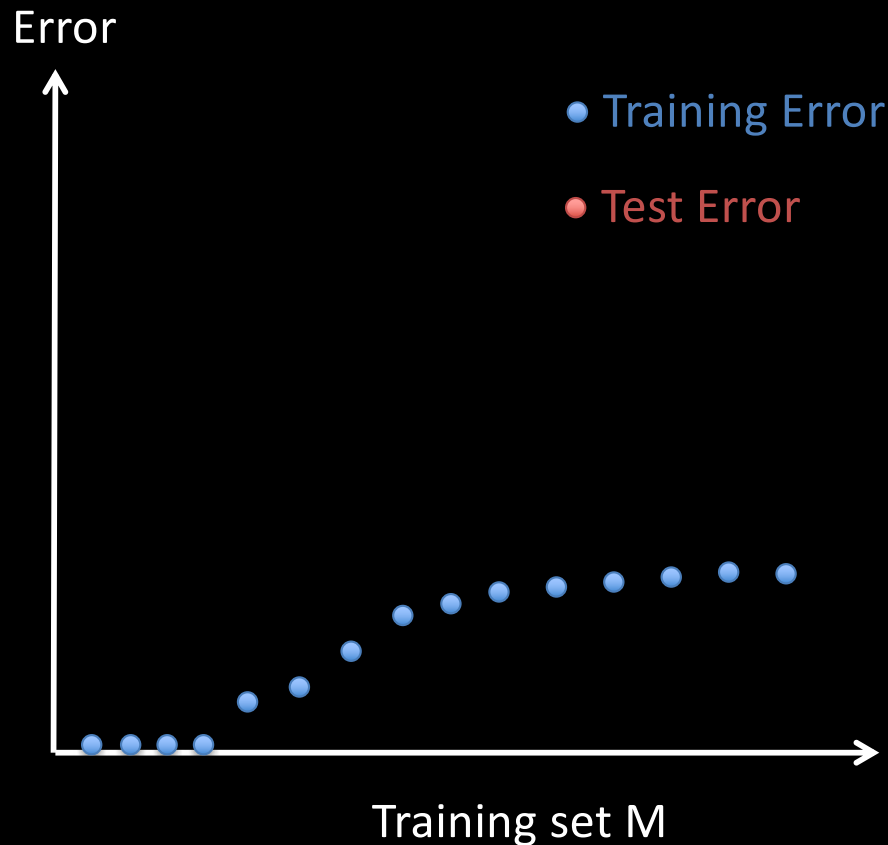
Bias and Variance



Bias and Variance



Bias and Variance



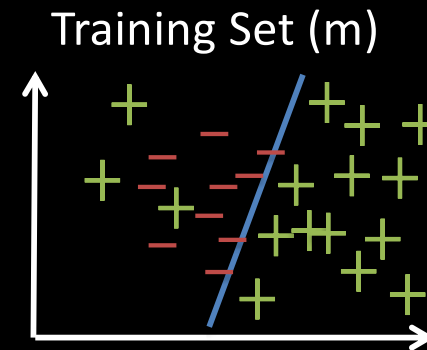
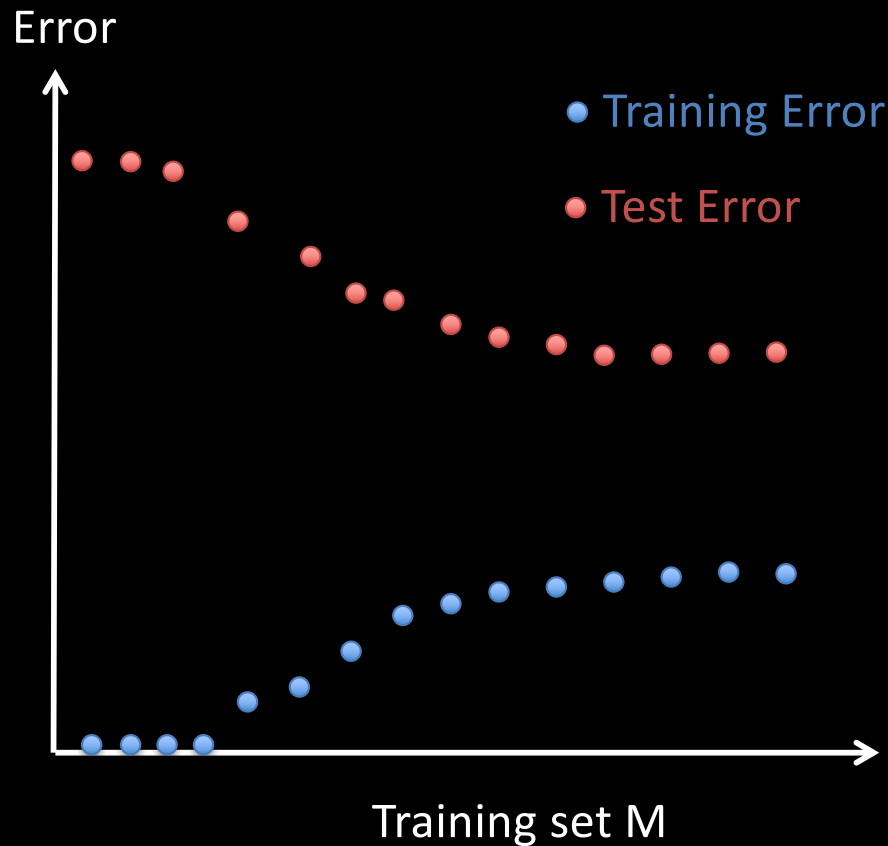
Clicker:

Test error

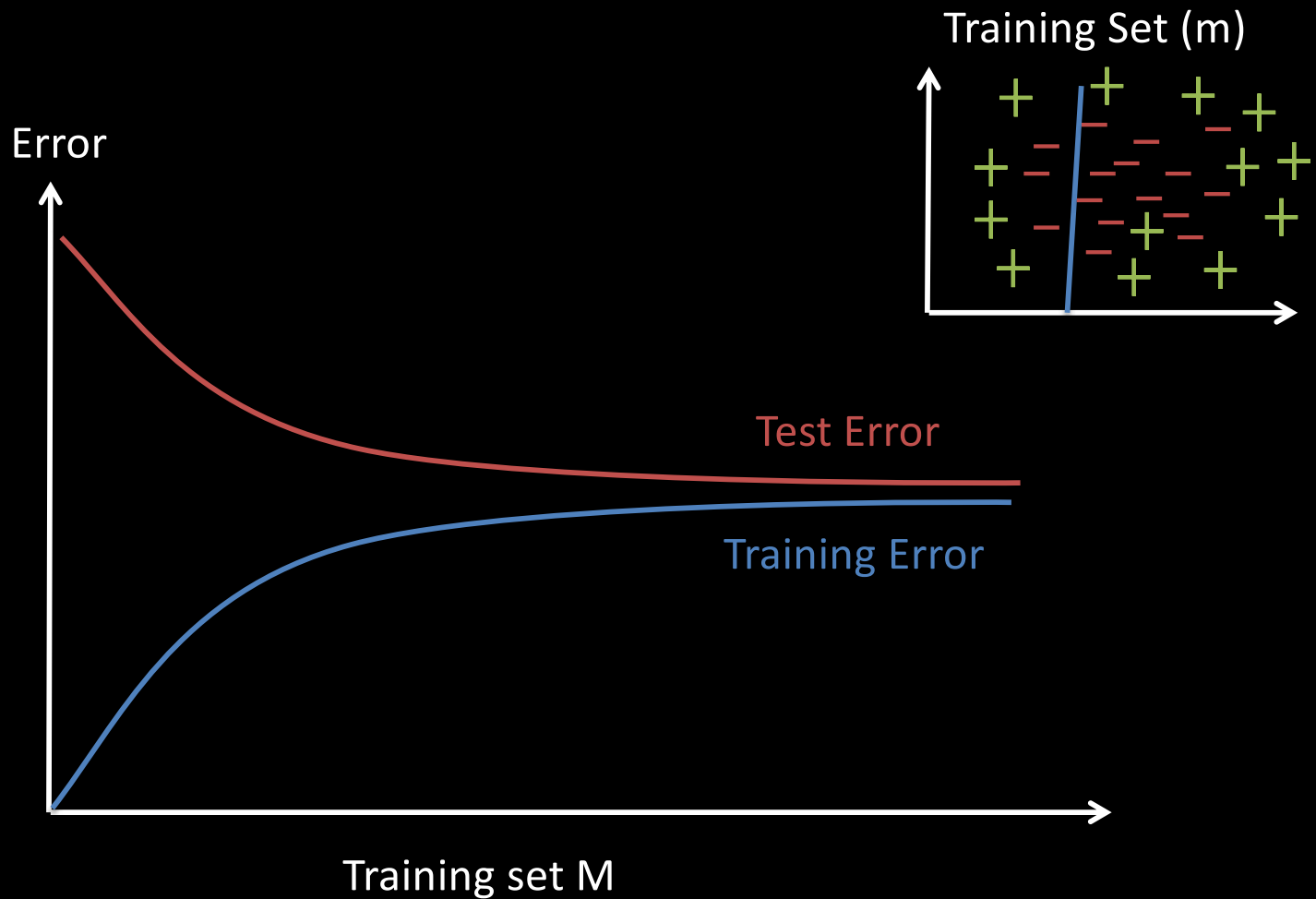
- a) decreases with M
- b) increases with M
- c) stays constant



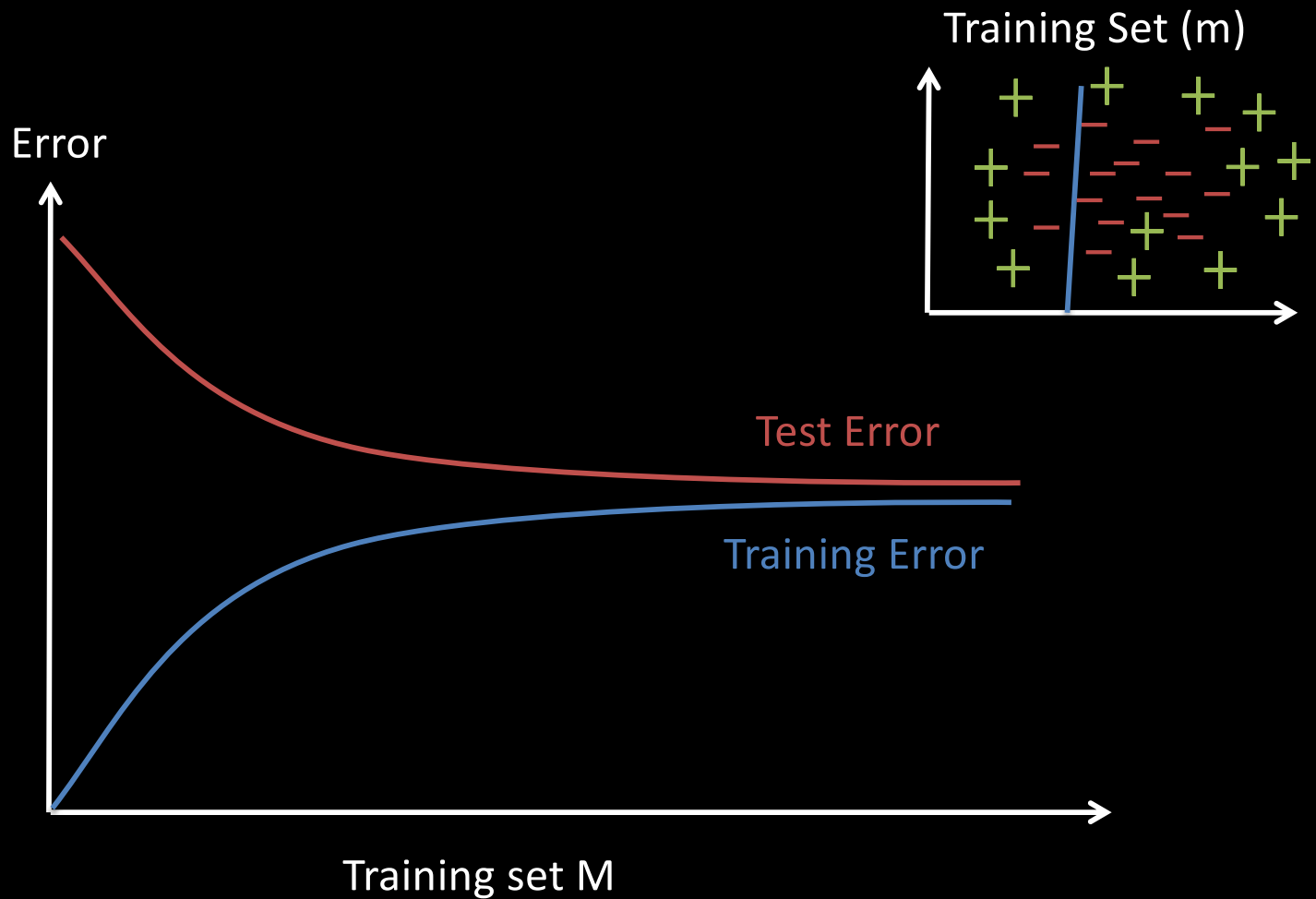
Bias and Variance



High Bias



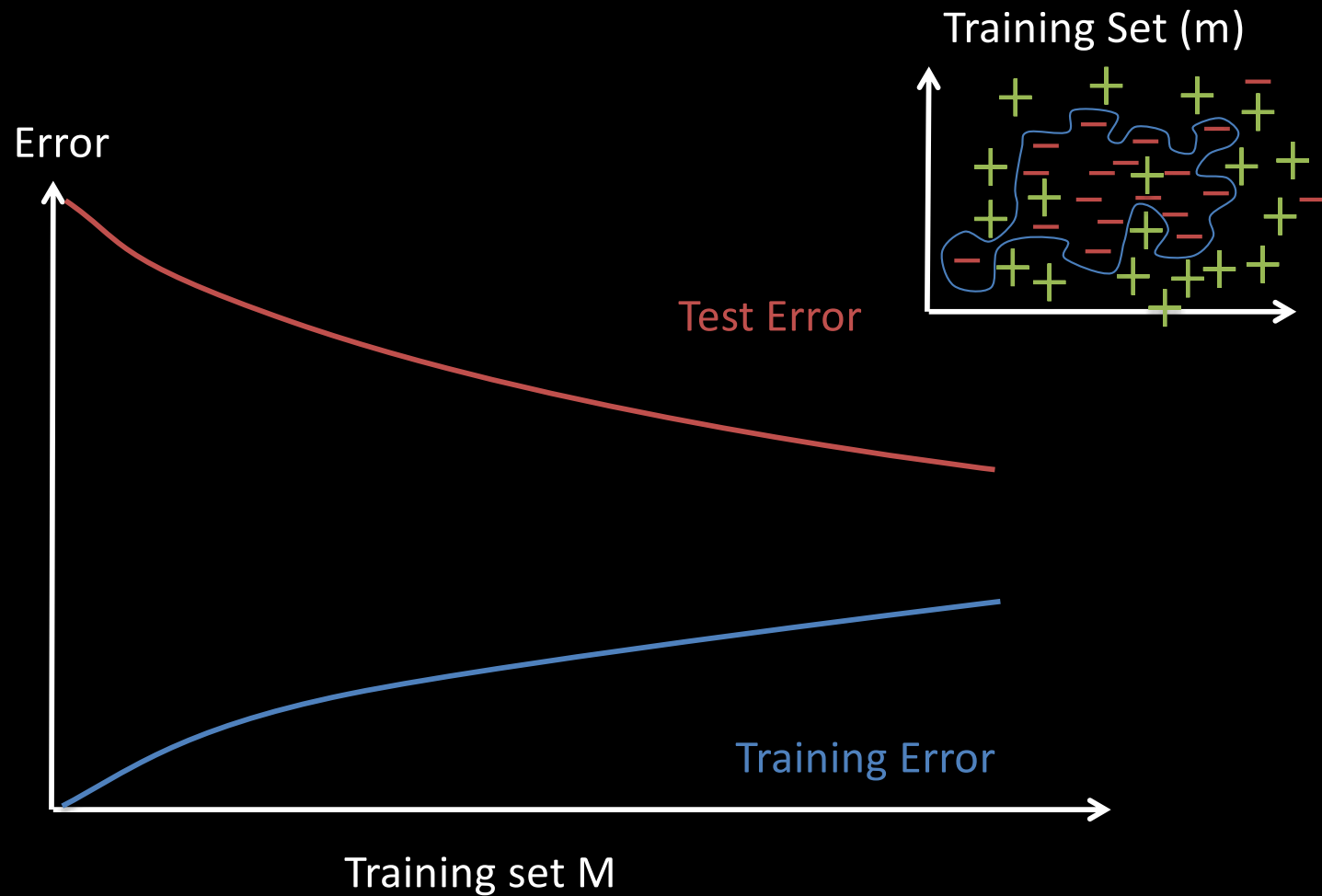
High Bias



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

- a) No
- b) Yes

High Variance



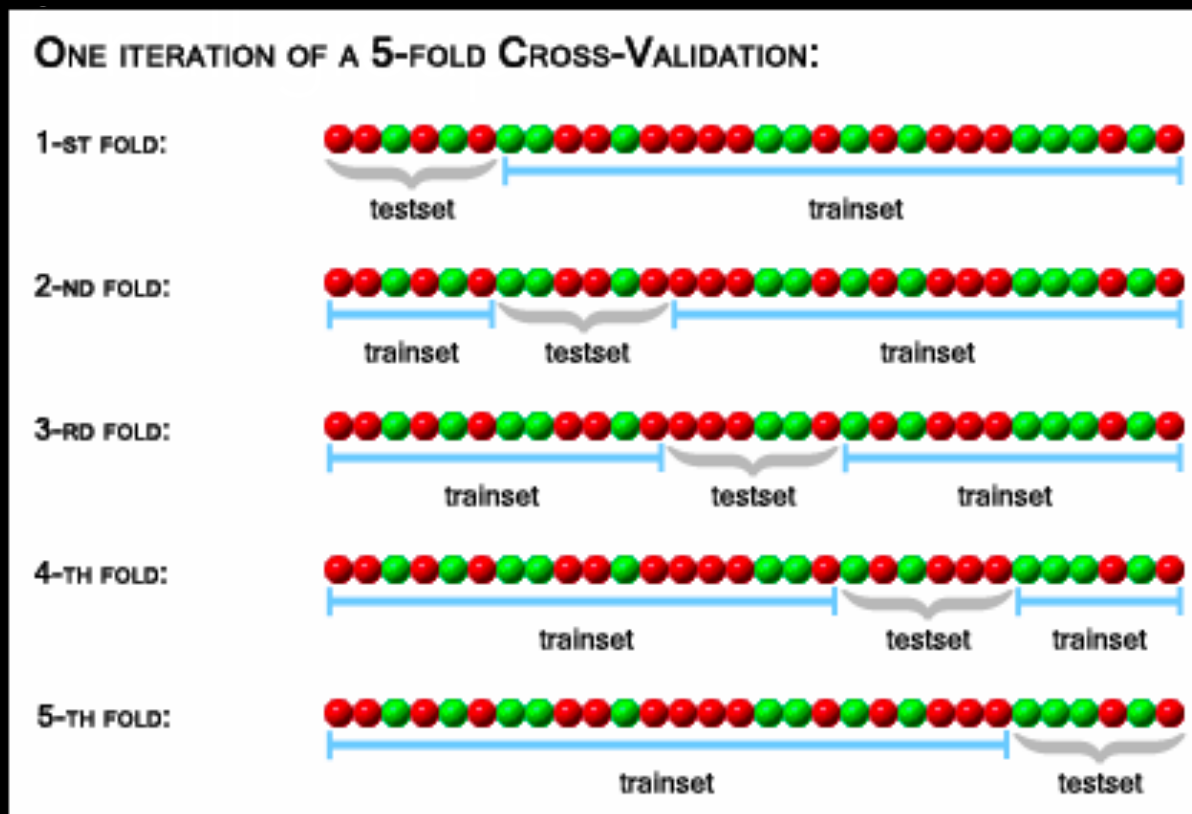
Clicker: If you have high-bias, does more data help?

- a) No
- b) Yes

Cross-validation

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

Repeat



Parameter Tuning

Grid Search

