

# Introduction to Machine Learning

Samrachana Adhikari  
Taman Upadhaya

December 22, 2018

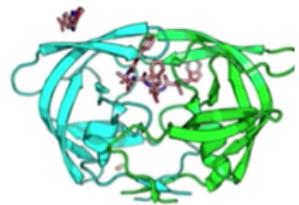
# State of the Art Applications of Machine Learning



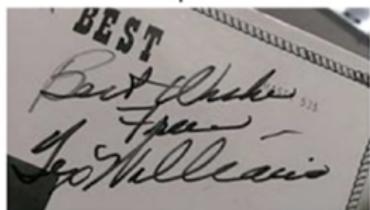
AlphaGo



Recommendation systems



Drug discovery



Character recognition



TWO SIGMA

Hedge fund stock predictions



Voice assistants



Assisted driving



Face detection/recognition



Cancer diagnosis

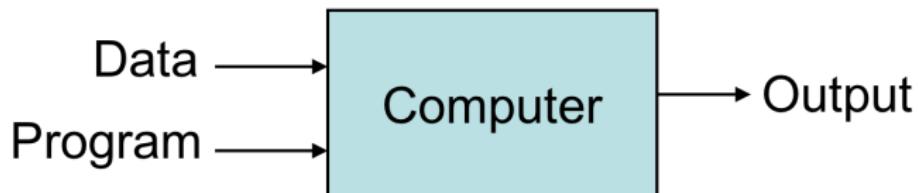
# What is in This Lecture Note

- Introduction
  - ▶ examples and application
- Supervised learning
  - ▶ Regression
  - ▶ classification
- Unsupervised learning
  - ▶ Clustering
  - ▶ ...
- Evaluation
  - ▶ Performance metrics
  - ▶ Bias-variance trade-off/ Over-fitting and under-fitting
  - ▶ Regularization methods
  - ▶ Cross-validation

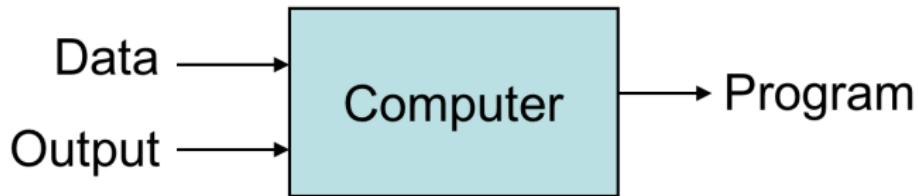
# Introduction

# What is Machine Learning

## Traditional Programming



## Machine Learning

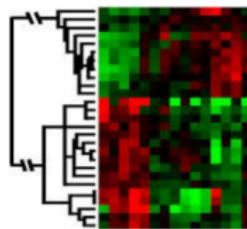


“Field of study that gives computers the ability to learn without being explicitly programmed” – Arthur Samuel 1959

# When Do We Use Machine Learning?

ML is used when:

- Human expertise does not exist (navigating on Mars)
- Models must be customized (personalized medicine)
- Models are based on huge amounts of data (genomics)
- ...



# Sample Applications

ML is used in:

- Web search
- Computational biology
- Finance
- E-commerce
- Space exploration
- Robotics
- Information extraction
- Social networks
- Debugging software
- Your favorite area ...

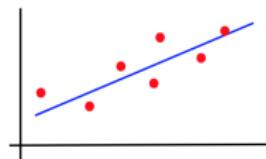
# Basic Paradigm

- **Training data:** set of example for observation
- **Testing data:** unseen set of example for prediction
- Variations on paradigm
  - ▶ Supervised: given a set of feature/label pairs, find a rules that predicts the label associated with a previously unseen input
  - ▶ Unsupervised: given a set of features vector (without label) group them into “natural cluster” (or create labels for groups)

# Three canonical learning problems

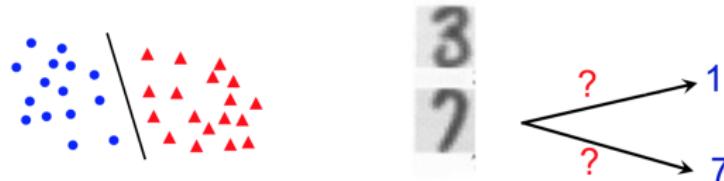
- ① Regression-supervised

estimate parameters, e.g. of weight vs height



- ② Classification-supervised

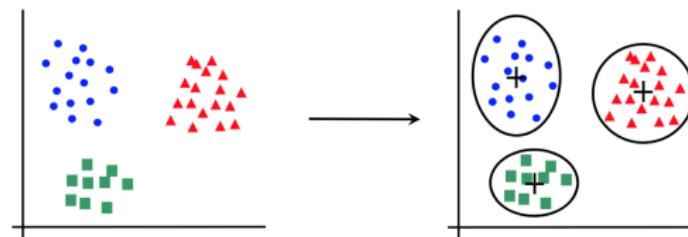
estimate class, e.g. handwritten digit classification



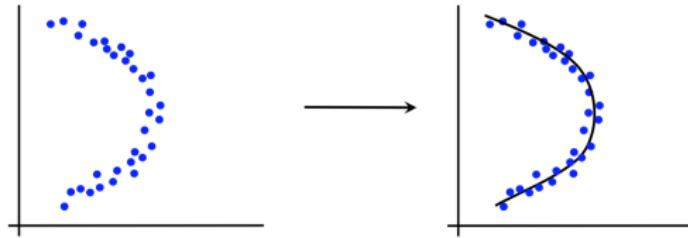
# Three canonical learning problems

## 3. Unsupervised learning model the data

- clustering



- dimensionality reduction



# Three components of ML

## ① Representation

- ▶ Regression model
- ▶ Neural networks
- ▶ Decision Tree
- ▶ Support Vector Machine, ...

## ② Evaluation

- ▶ Accuracy
- ▶ Squared error
- ▶ Precision and recall
- ▶ Posterior probability
- ▶ Entropy, ...

## ③ Optimization

- ▶ Combinatorial optimization, e.g. Greedy search
- ▶ Convex optimization, e.g Gradient descent
- ▶ Constrained optimization, e.g Linear programming

# Three components of ML

## ① Representation

- ▶ Regression model
- ▶ Neural networks
- ▶ Decision Tree
- ▶ Support Vector Machine, ...

## ② Evaluation

- ▶ Accuracy
- ▶ Squared error
- ▶ Precision and recall
- ▶ Posterior probability
- ▶ Entropy, ...

## ③ Optimization

- ▶ Combinatorial optimization, e.g. Greedy search
- ▶ Convex optimization, e.g Gradient descent
- ▶ Constrained optimization, e.g Linear programming

# Three components of ML

## ① Representation

- ▶ Regression model
- ▶ Neural networks
- ▶ Decision Tree
- ▶ Support Vector Machine, ...

## ② Evaluation

- ▶ Accuracy
- ▶ Squared error
- ▶ Precision and recall
- ▶ Posterior probability
- ▶ Entropy, ...

## ③ Optimization

- ▶ Combinatorial optimization, e.g. Greedy search
- ▶ Convex optimization, e.g Gradient descent
- ▶ Constrained optimization, e.g Linear programming

# Three components of ML

## ① Representation

- ▶ Regression model
- ▶ Neural networks
- ▶ Decision Tree
- ▶ Support Vector Machine, ...

## ② Evaluation

- ▶ Accuracy
- ▶ Squared error
- ▶ Precision and recall
- ▶ Posterior probability
- ▶ Entropy, ...

## ③ Optimization

- ▶ Combinatorial optimization, e.g. Greedy search
- ▶ Convex optimization, e.g Gradient descent
- ▶ Constrained optimization, e.g Linear programming

# Supervised learning

# Supervised learning: Motivation

- ① Predict weight from gender, height, age, ...
- ② Predict Google stock price today from Google, Yahoo, Microsoft prices yesterday
- ③ Predict each pixel intensity in robots current camera image, from previous image and previous action
- ④ **Discuss other examples in small group**

# Supervised Learning

## Formal motivation

Given a set of measurements  $(x_i, y_i)$  for,  $i = 1, \dots, N$  subjects, known as the **training data**, we want to construct a prediction rule for an arbitrary input  $x_0$ .

- $X$ s are often called **predictors or covariates**.
- $Y$  is the **outcome variable**; can be continuous, categorical or ordinal.
- $y$  is used to represent a realization of the random variable  $Y$ .

# Supervised learning: Linear Regression

Model continuous outcome as a **Linear** function of predictors

$$Y = \sum_{k=1}^p X^k \beta_k + \epsilon = X\beta + \epsilon$$
$$E(\epsilon) = 0$$

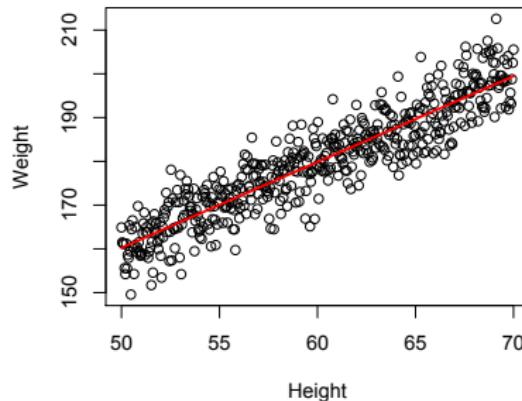
$Y$ : A **continuous outcome**

$X = [X_1, \dots, X_p]$ : matrix of  $p$

**predictors or features**

$\beta$ : a vector of **coefficients**

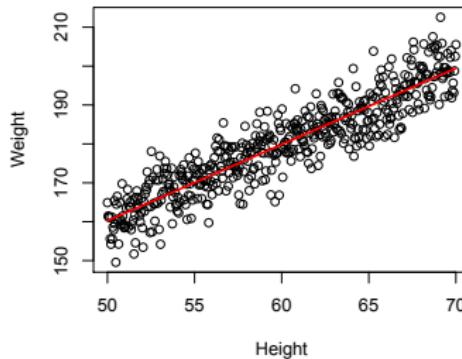
$\epsilon$  is sometimes called **noise or error term**



# How to predict $Y$ at an arbitrary $x_0$ ?

- ① **Training:** Define objective function to estimate  $\beta$  using training data
- ② **Prediction:** Predict outcome at  $x_0$  as

$$E(Y|x_0) = \hat{Y}|x_0 = \hat{\beta}x_0$$



# Linear Regression Estimation: Maximum Likelihood Estimate

MLE  $\beta$

$$\hat{\beta} = \operatorname{argmax}_{\beta} [\text{likelihood}(\beta)] \equiv \operatorname{argmin}_{\beta} [-\text{likelihood}(\beta)]$$

Let for each  $i$ ,  $\epsilon_i \sim \text{Normal}(0, \sigma^2)$

$$\Rightarrow Y_i | X_i, \beta \sim \text{Normal}(X_i \beta, \sigma^2)$$

$$\text{Likelihood}(\beta) = p(Y|X, \beta) = \prod_{i=1}^n p(Y_i|X_i, \beta)$$

- We are assuming each  $Y_i$  is independent and identically distributed
- Log of likelihood is used in practice

# Linear Regression Estimation: Maximum Likelihood Estimate

MLE  $\beta$

$$\hat{\beta} = \operatorname{argmax}_{\beta} [\text{likelihood}(\beta)] \equiv \operatorname{argmin}_{\beta} [-\text{likelihood}(\beta)]$$

Let for each  $i$ ,  $\epsilon_i \sim \text{Normal}(0, \sigma^2)$

$$\Rightarrow Y_i | X_i, \beta \sim \text{Normal}(X_i \beta, \sigma^2)$$

$$\text{Likelihood}(\beta) = p(Y|X, \beta) = \prod_{i=1}^n p(Y_i|X_i, \beta)$$

- We are assuming each  $Y_i$  is independent and identically distributed
- Log of likelihood is used in practice

# Linear Regression Estimation: Maximum Likelihood Estimate

MLE  $\beta$

$$\hat{\beta} = \operatorname{argmax}_{\beta} [\text{likelihood}(\beta)] \equiv \operatorname{argmin}_{\beta} [-\text{likelihood}(\beta)]$$

Let for each  $i$ ,  $\epsilon_i \sim \text{Normal}(0, \sigma^2)$

$$\Rightarrow Y_i | X_i, \beta \sim \text{Normal}(X_i \beta, \sigma^2)$$

$$\text{Likelihood}(\beta) = p(Y|X, \beta) = \prod_{i=1}^n p(Y_i|X_i, \beta)$$

- We are assuming each  $Y_i$  is independent and identically distributed
- Log of likelihood is used in practice

# MLE in Linear Regression

- $Y_1, \dots, Y_n$  are iid outcomes
- For each  $i$ ,  $\epsilon_i \sim \text{Normal}(0, \sigma^2)$

$$\Rightarrow Y_i | X_i, \beta \sim \text{Normal}(X_i \beta, \sigma^2)$$

$$\text{Likelihood}(\beta) = p(Y|X, \beta) = \prod_{i=1}^n p(Y_i|X_i, \beta)$$

$$\text{Log-Likelihood}(\beta) = \log p(Y|X, \beta) = \sum_{i=1}^n \log p(Y_i|X_i, \beta)$$

$$\bullet \log p(Y_i|X_i, \beta, \sigma^2) = \log \frac{1}{\sqrt{2\pi}\sigma} + \log \exp\left[\frac{-(Y_i - X_i \beta)^2}{2\sigma^2}\right]$$

## Discussion

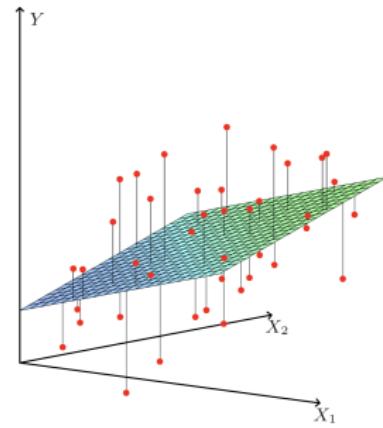
Derive MLE of  $\beta$  in linear regression.

# Linear Regression Estimation: MSE Loss

## Mean Square Error Loss

$$\begin{aligned}\hat{\beta} &= \operatorname{argmin}_{\beta} \|Y - X\beta\|_2^2 = \operatorname{argmin}_{\beta} \sum_{i=1}^n (Y_i - X_i\beta)^2 \\ \hat{\beta} &= (X^T X)^{-1} X^T Y\end{aligned}$$

- When  $\epsilon \sim \text{Normal}(\mu, \sigma^2)$ , minimizing squared error loss is equivalent to maximizing likelihood

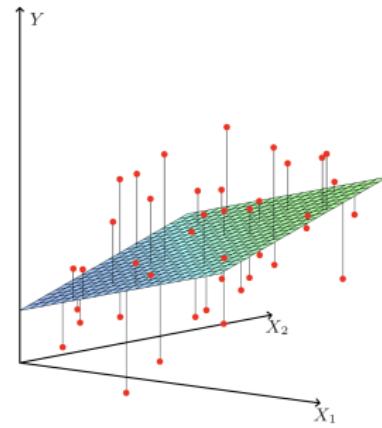


# Linear Regression Estimation: MSE Loss

## Mean Square Error Loss

$$\begin{aligned}\hat{\beta} &= \operatorname{argmin}_{\beta} \|Y - X\beta\|_2^2 = \operatorname{argmin}_{\beta} \sum_{i=1}^n (Y_i - X_i\beta)^2 \\ \hat{\beta} &= (X^T X)^{-1} X^T Y\end{aligned}$$

- When  $\epsilon \sim \text{Normal}(\mu, \sigma^2)$ , minimizing squared error loss is equivalent to maximizing likelihood



# Linear Regression Discussion

Pros:

- Interpretability and simplicity
- Solution exists (in most cases)
- Favorable large-sample statistical properties
- Numerous software packages to fit the model

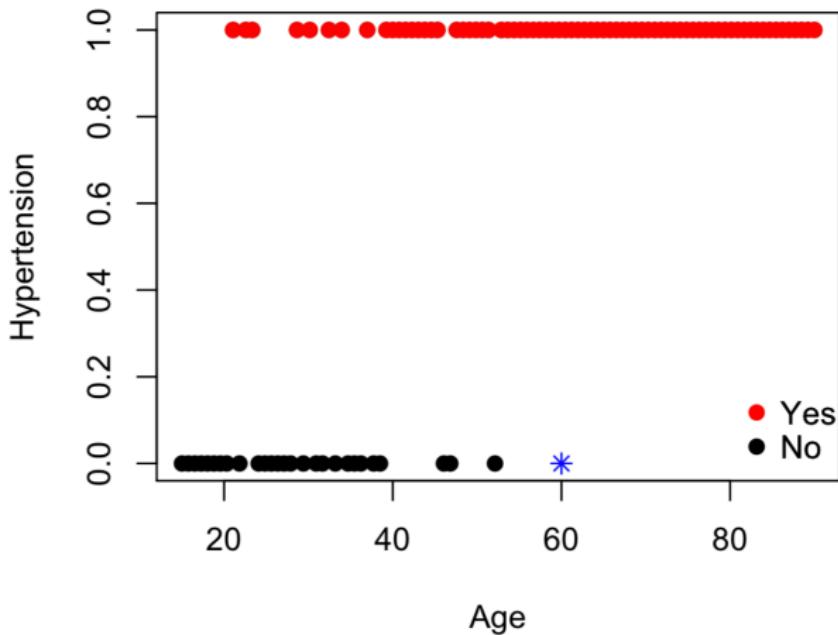
Cons:

- Real world is never that simple!!!
  - ▶ **Too restrictive:** Assumes linear relationship between features and outcomes
  - ▶ **Big data:** Not suitable when feature space is large ( $p \gg n$ )

# Supervised learning: Classification models

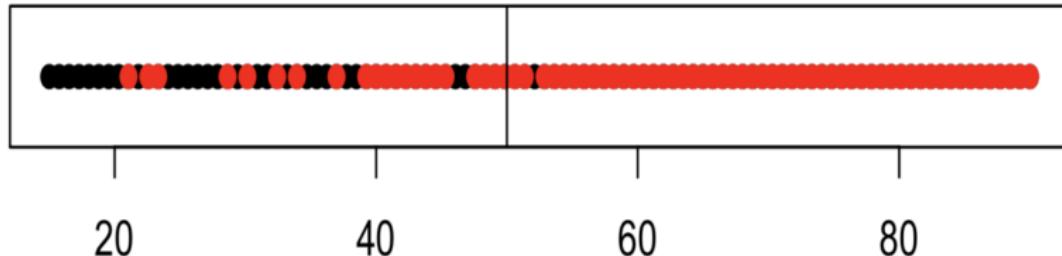
- Outcome is discrete
  - ▶ Predict weather, rainy versus sunny, based on past weather.
  - ▶ Predict whether a patient will survive after a major surgery using vitals prior to the surgery.
  - ▶ Classifying patients into high or low risk of hypertension based on their age.
- Goal: find a rule to correctly classify subjects into classes by learning from the existing data.

# Classification: Motivation

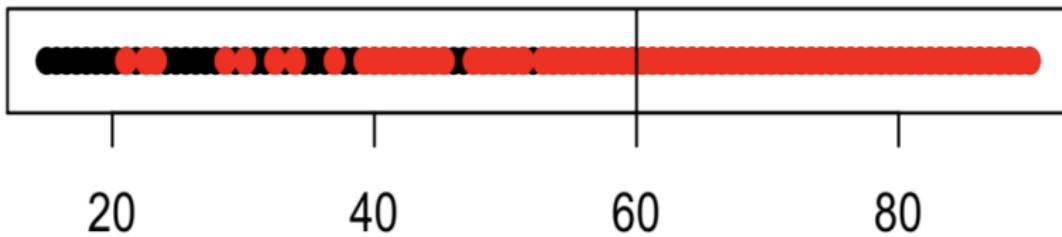


# Classification: Motivation

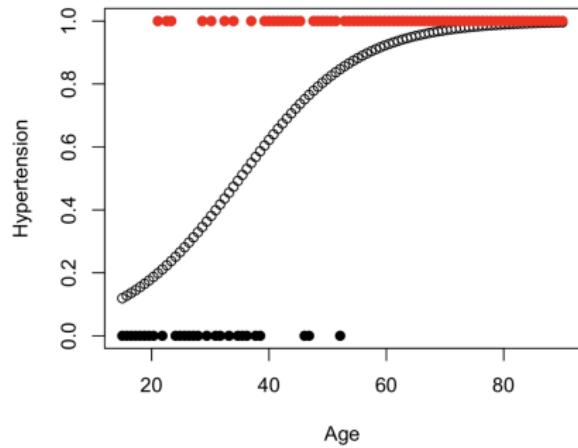
Rule 1



Rule 2



# Logistic Regression



$$\log \frac{Pr(Y = 1)}{Pr(Y = 0)} = X\beta$$

Prediction:

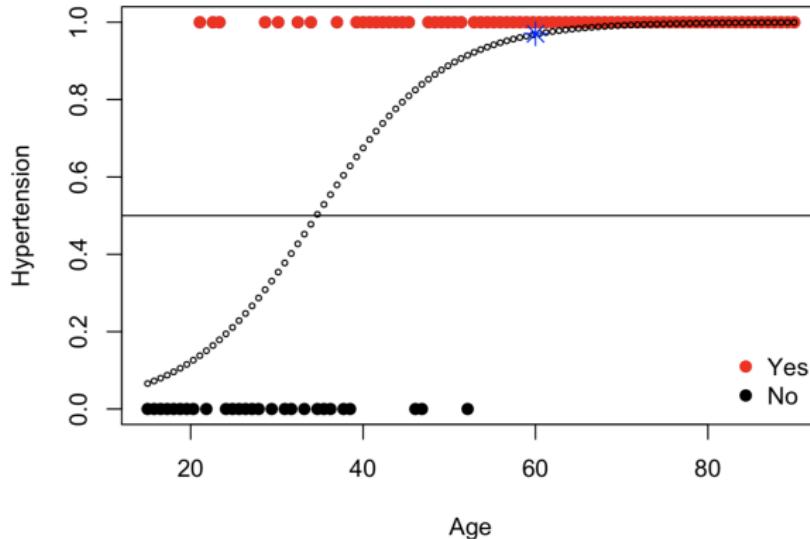
$$\begin{aligned}\hat{p} &= \hat{Pr}(Y = 1)|x_o &= \frac{1}{1 + \text{Exp}(-x_0\hat{\beta})} \\ \hat{Y} &= \mathbf{I}(\hat{p} > 0.5)\end{aligned}$$

# Decision Boundary: Linear Classifier

- A convenient property of this form is that it leads to a simple **linear** expression for classification.
- To classify any given  $Y$  we generally want to assign the value  $k$  that maximizes  $Pr(Y = k|X)$ , for  $k = \{0, 1\}$
- Classify  $Y = 1$  if

$$\begin{aligned} Pr(Y = 1|X) &> Pr(Y = 0|X) \\ \equiv 1 &> \text{Exp}(\beta X) \equiv 0 > \beta X \end{aligned}$$

- **Classification rule is linear in  $X$**



For 2-classes case,  
classify  $\hat{Y}|x_0 = 1$  if  
 $\hat{Pr}(Y = 1|x_0) > 0.5$ .

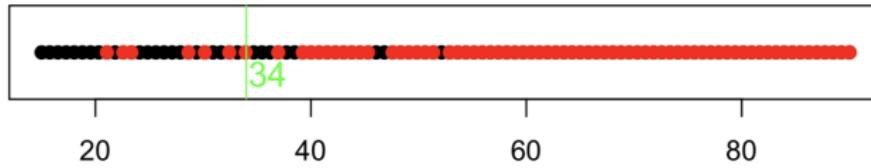


Figure: Linear Classifier

# Supervised learning: Regression Models

## A general regression model

$$E(Y|X) = f(X)$$

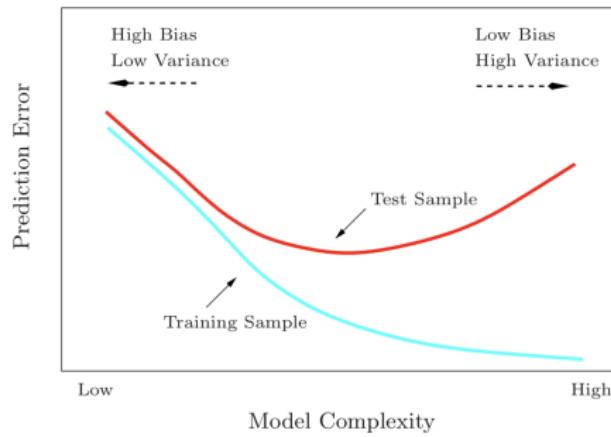
- Many other functional forms on  $f$  is possible
- $f$  is often referred to as a **link function** in statistical learning
- Parametric versus non-parametric models: Restrict  $f$  to rely on few parameters versus using  $f$  in a bigger class of family with no restrictions
- How to decide which  $f$  to use?

# General Estimation: Minimization of Expected Risk or Loss Function

$$\hat{f} = \arg \min_f E(\text{Prediction loss}(f))$$

# General Estimation: Minimization of Expected Risk or Loss Function

- Optimization methods are used to do the minimization.
- Challenging if the loss function is not convex.
- Desired to build a most parsimonious model that minimizes the prediction error.



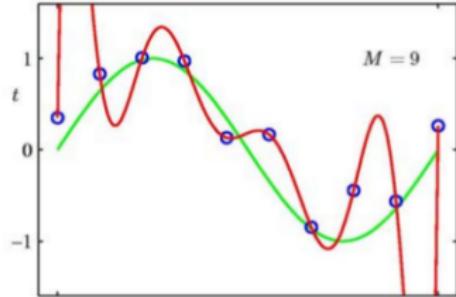
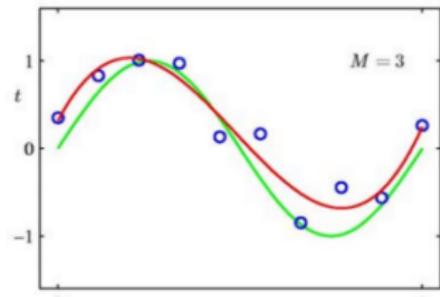
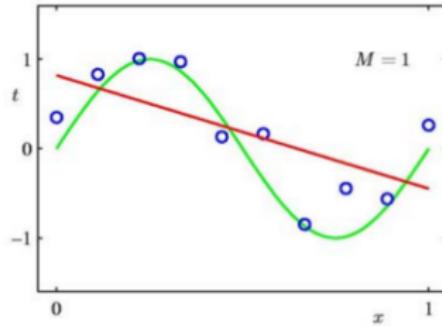
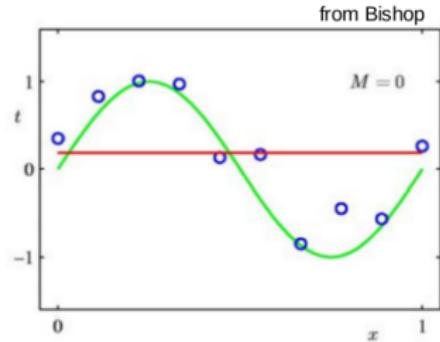
# Performance metrics

Truth table (confusion matrix)

		Predicted class	
		<i>P</i>	<i>N</i>
<b>Actual Class</b>	<i>P</i>	True Positives (TP)	False Negatives (FN)
	<i>N</i>	False Positives (FP)	True Negatives (TN)

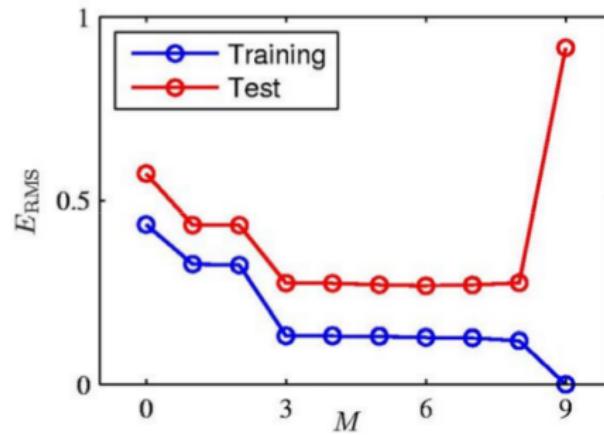
- Accuracy =  $\frac{TP+TN}{P+N}$
- Sensitivity =  $\frac{TP}{\text{Actual Positives}}$
- Specificity =  $\frac{TN}{\text{Actual Negatives}}$
- ...

# Some fits to the data: which is best?



# Bias-variance trade-off

- test data: a different sample from the same true function



RootMeanSquare (RMS) Error

- training error goes to zero, but test error increases with M

$$\text{RMSE} = \sqrt{\left(\frac{(Y_i - \hat{Y}_i)^2}{n}\right)}$$

# Unsupervised learning

# Whats is Data Clustering

- The organization of unlabeled data into similarity groups called clusters.
- A cluster is a collection of data items which are similar between them, and dissimilar to data items in other clusters.

# Whats is Data Clustering

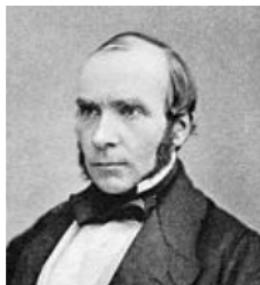
- The organization of unlabeled data into similarity groups called clusters.
- A cluster is a collection of data items which are similar between them, and dissimilar to data items in other clusters.



Can be done in many ways

# Historic application of clustering

- John Snow, a London physician plotted the location of cholera deaths on a map during an outbreak in the 1855.
- The location indicate that cases were clustered around certain intersections where there were polluted wells – thus exposing both the problem and solution.



John Snow and famous cholera map

# What do we need for clustering?

## ① Proximity measure:

- ▶ similarity measure  $S(X_i, X_k)$ : large if  $X_i, X_k$  are similar
- ▶ dissimilarity (or distance) measure  $D(X_i, X_k)$ : small if  $X_i, X_k$  are similar



## ② Criterion function to evaluate a clustering

## ③ Algorithm to computer clustering

- ▶ For example, by optimizing the criterion function

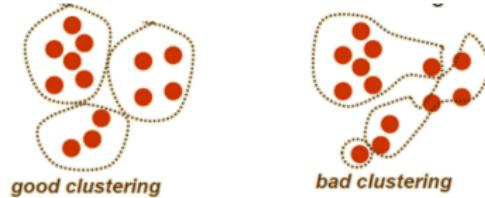
# What do we need for clustering?

## ① Proximity measure:

- ▶ similarity measure  $S(X_i, X_k)$ : large if  $X_i, X_k$  are similar
- ▶ dissimilarity (or distance) measure  $D(X_i, X_k)$ : small if  $X_i, X_k$  are similar



## ② Criterion function to evaluate a clustering



## ③ Algorithm to computer clustering

- ▶ For example, by optimizing the criterion function

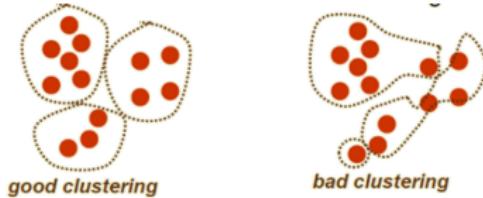
# What do we need for clustering?

## ① Proximity measure:

- ▶ similarity measure  $S(X_i, X_k)$ : large if  $X_i, X_k$  are similar
- ▶ dissimilarity (or distance) measure  $D(X_i, X_k)$ : small if  $X_i, X_k$  are similar



## ② Criterion function to evaluate a clustering



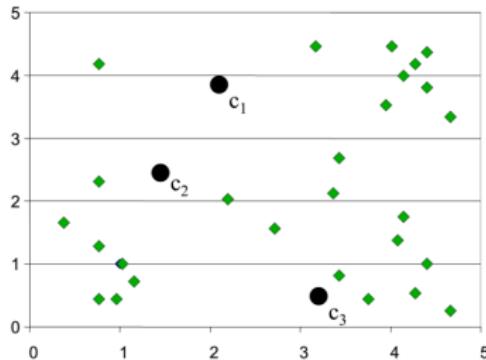
## ③ Algorithm to computer clustering

- ▶ For example, by optimizing the criterion function

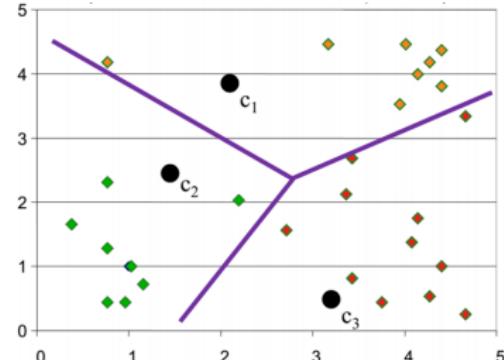
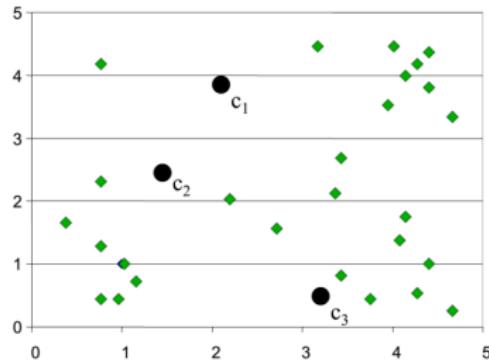
# K-Means Clustering

- K-means (MacQueen, 1967) is a partitional clustering algorithm
- Let the set of data points D be  $x_1, x_2, \dots, x_n$ 
  - ▶ where  $x_i = (x_{i1}, x_{i2}, \dots, x_{ir})$  is a vector in  $X \subseteq R^n$  and n is the number of dimensions
- The k-means algorithm partitions the given data into K clusters:
  - ▶ Each cluster has a cluster center C, called centroid.
  - ▶ K is specified by the user
- convergence (stopping) criterion
  - ▶ minimum decrease in the sum of squared error (SSE)
$$SSE = \sum_{j=1}^k \sum_{x \in C_j} d(x, m_j)^2$$

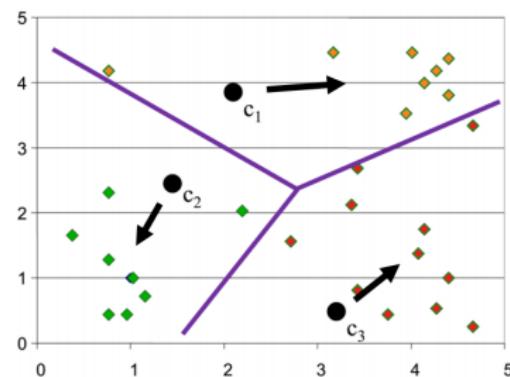
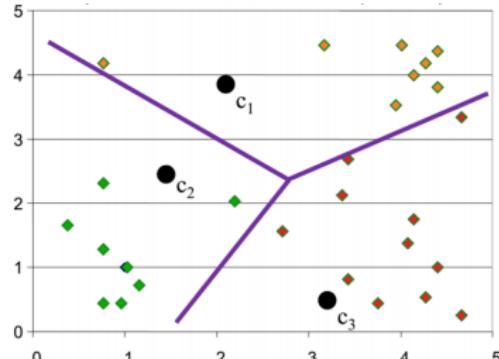
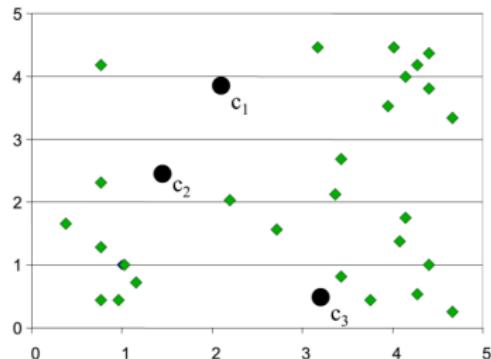
## K-Means Clustering: working example



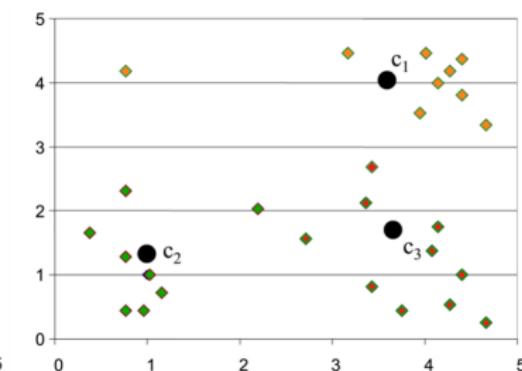
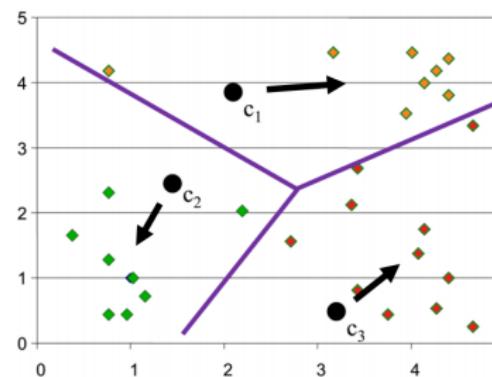
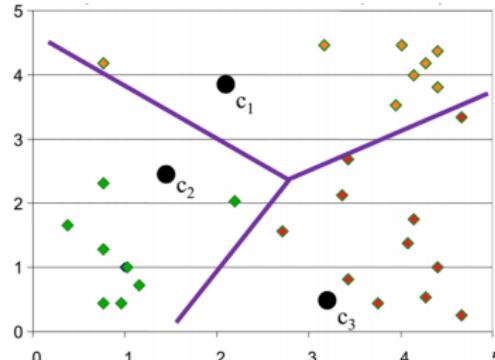
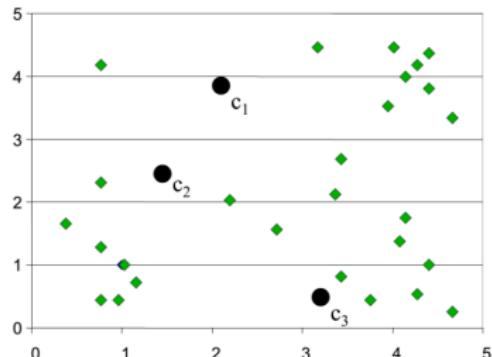
# K-Means Clustering: working example



# K-Means Clustering: working example

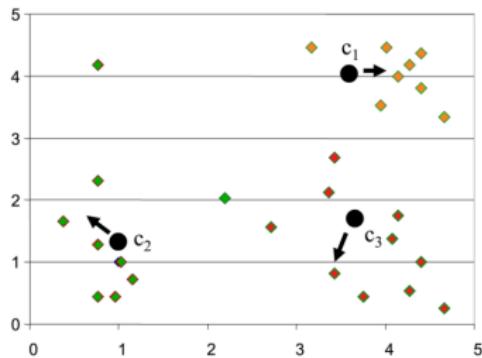


# K-Means Clustering: working example

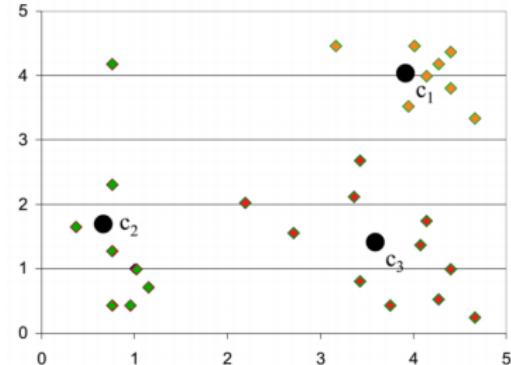
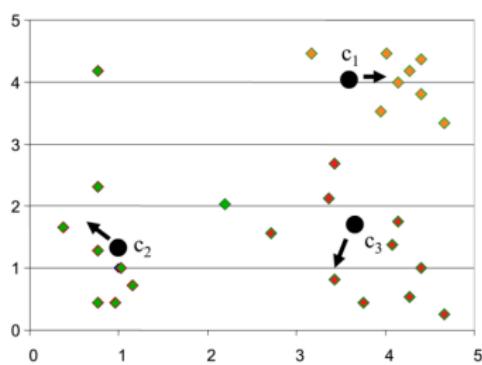


Result of first iteration

# K-Means Clustering: working example



# K-Means Clustering: working example



Result of second iteration

# K-Means Clustering: lab working example

- Demo class work:  
[http://stanford.edu/class/ee103/visualizations/kmeans/  
kmeans.html](http://stanford.edu/class/ee103/visualizations/kmeans/kmeans.html)
- LAB:
  - ① Data science: any working example (Samarachana)
  - ② Medical imaging: brain tumor segmentation (Taman)

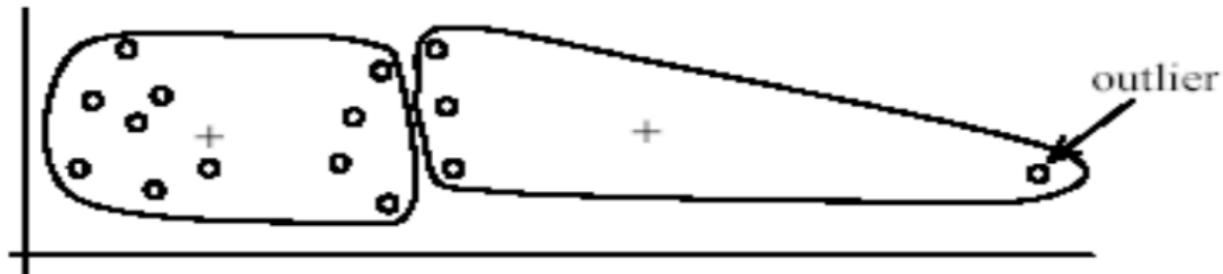
# Why use K-means?

- Strengths:
  - ▶ Simple: easy to understand and to implement
  - ▶ Efficient: Time complexity:  $O(tkn)$ , where  $n$  is the number of data points,  $k$  is the number of clusters, and  $t$  is the number of iterations.
  - ▶ Since both  $k$  and  $t$  are small. k-means is considered a linear algorithm.
- K-means is the most popular clustering algorithm.
- Note that: it terminates at a local optimum if SSE is used. The global optimum is hard to find due to complexity.

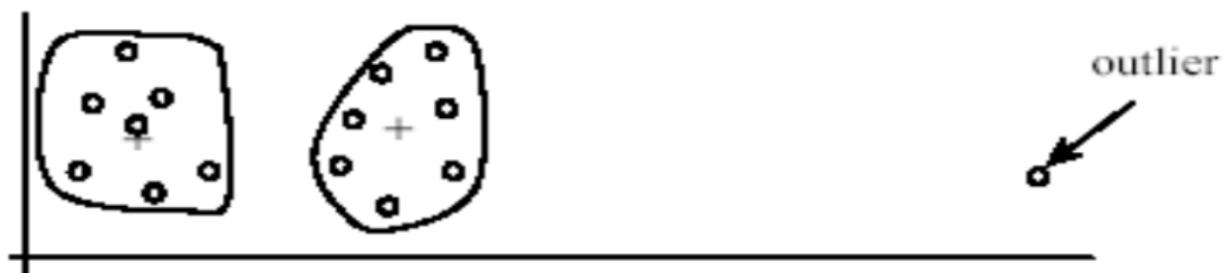
# Weaknesses of K-means

- The algorithm is only applicable if the mean is defined.
  - ▶ For categorical data, k-mode - the centroid is represented by most frequent values.
- The user needs to specify k.
- The algorithm is sensitive to outliers
  - ▶ Outliers are data points that are very far away from other data points.
  - ▶ Outliers could be errors in the data recording or some special data points with very different values.

# Outliers



(A): Undesirable clusters

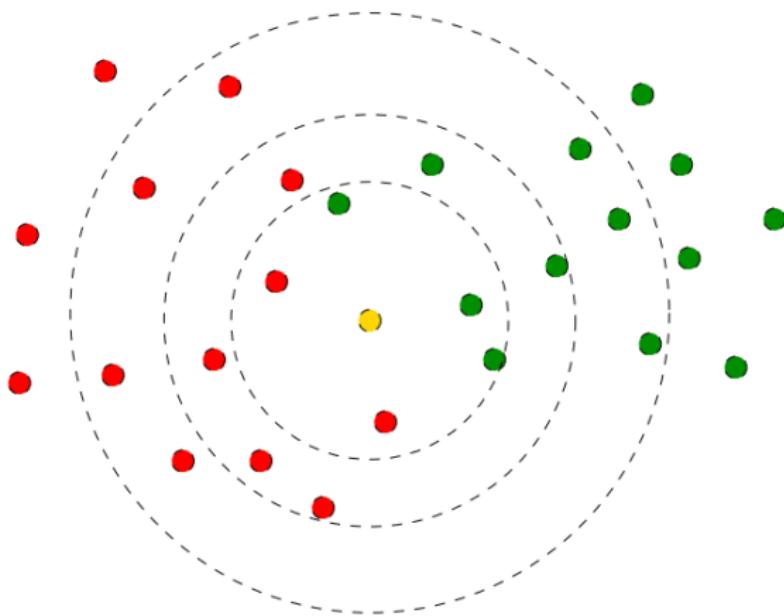


(B): Ideal clusters

# K Nearest Neighbors

- **K-nearest neighbours** uses the local neighborhood to obtain a prediction
- The K memorized examples more similar to the one that is being classified are retrieved
- A distance function is needed to compare the examples similarity
  - ▶ Euclidean distance  $d(x_j, x_k) = \sqrt{\sum_i (x_{j,i} - x_{k,i})^2}$
  - ▶ Manhattan distance  $d(x_j, x_k) = \sqrt{\sum_i |x_{j,i} - x_{k,i}|^2}$
- This means that if we change the distance function, we change how examples are classified

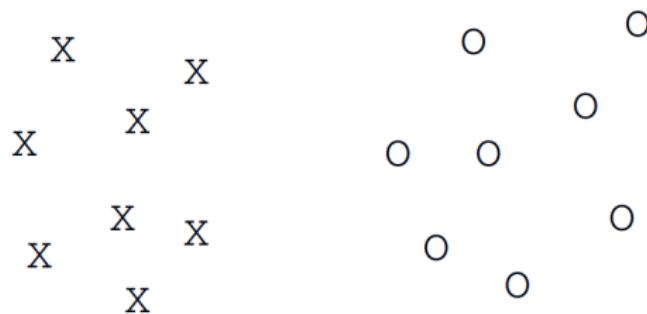
# K Nearest Neighbors



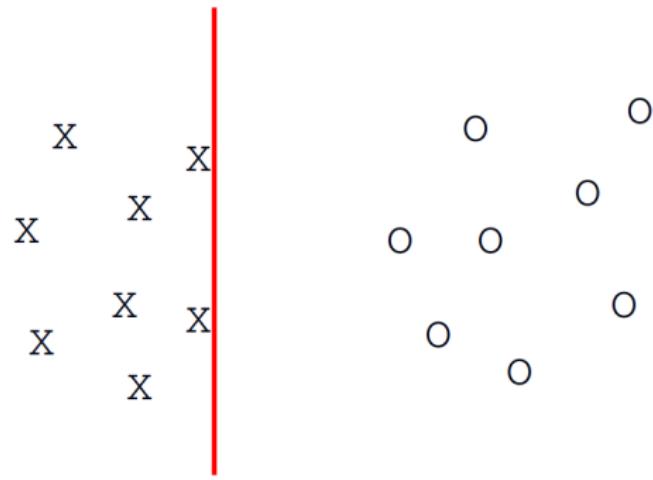
# Support Vector Machine

Support Vector Machine

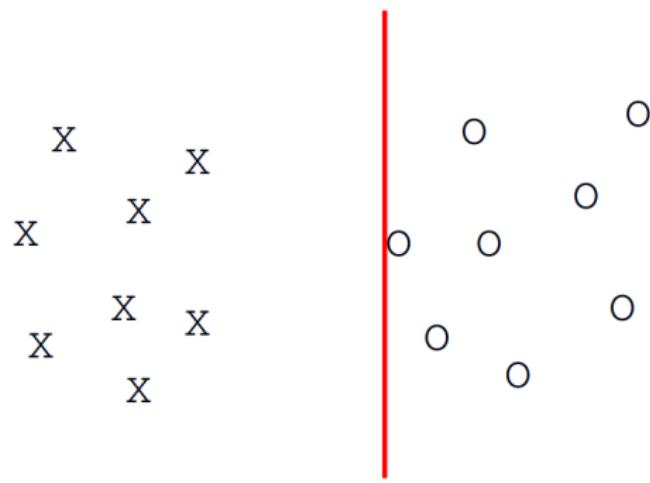
# Intuitions



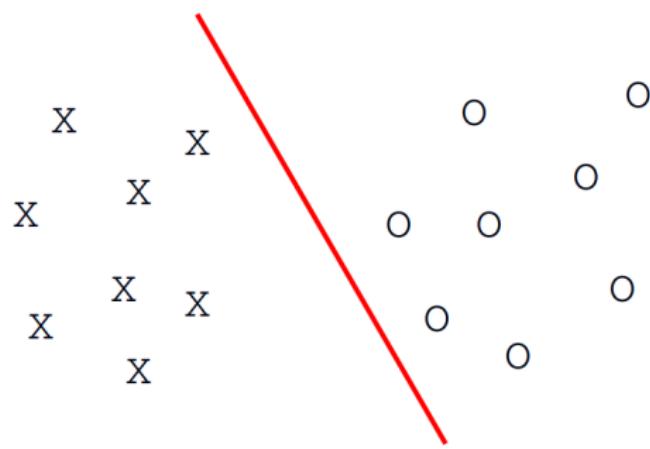
# Intuitions



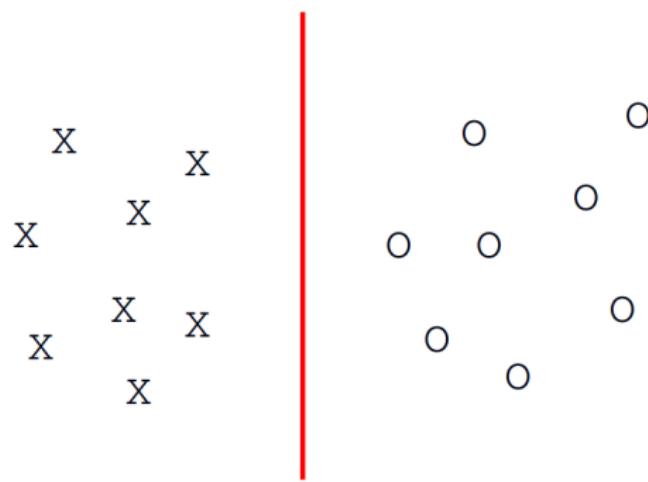
# Intuitions



# Intuitions

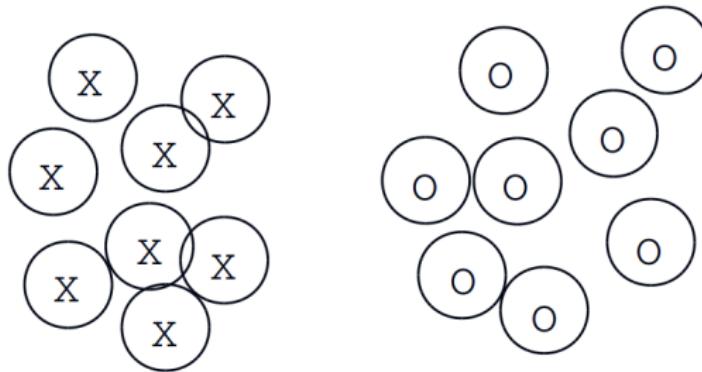


# Intuitions

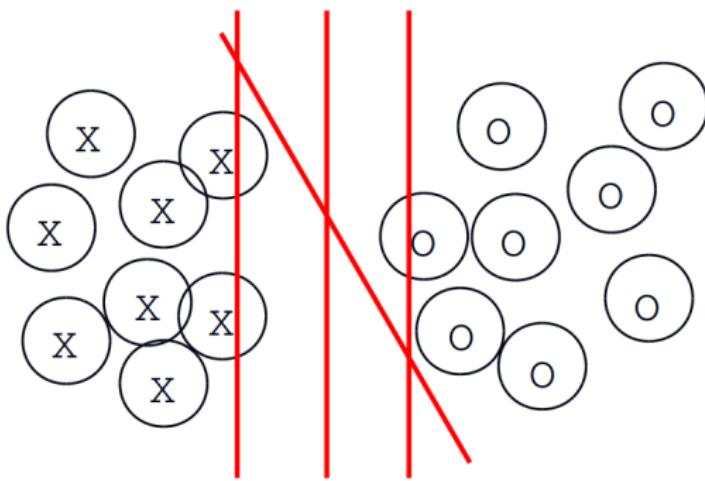


A "Good" Separator Noise in the Observations

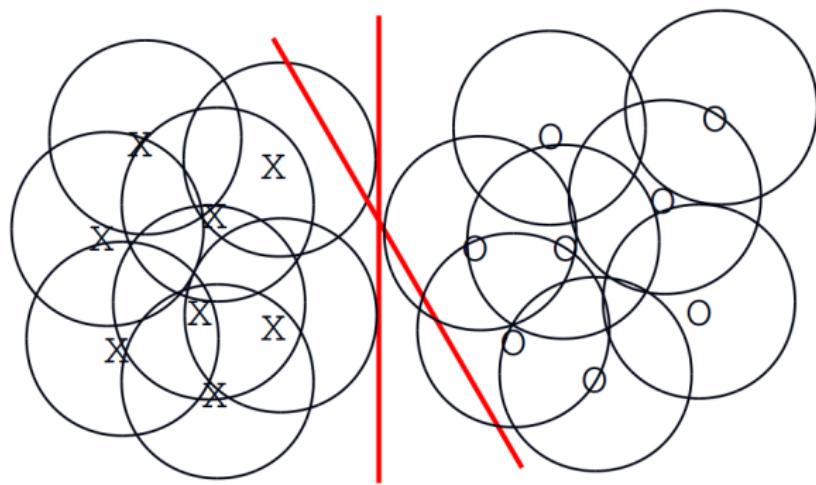
# Intuitions



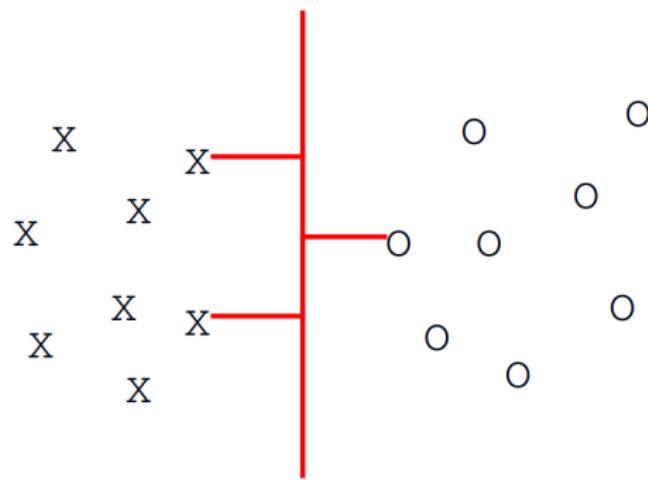
# Intuitions



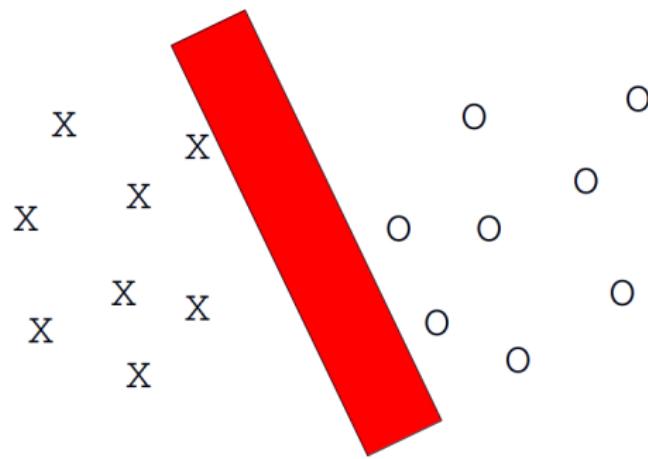
# Intuitions



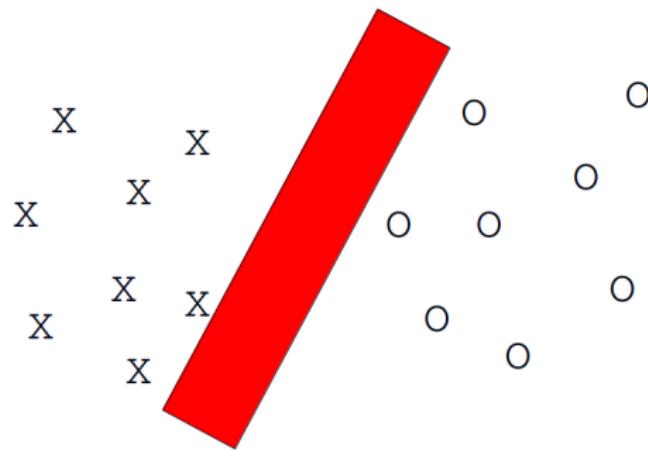
# Intuitions



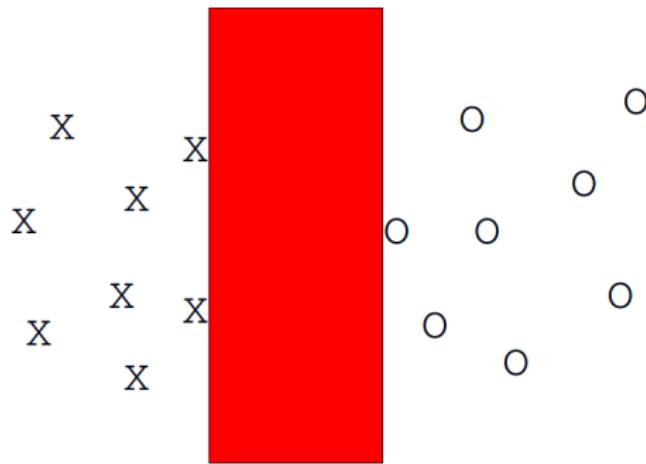
# Intuitions



# Intuitions



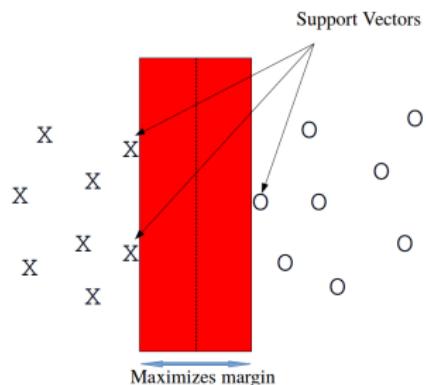
# Intuitions



Idea: to select the separating hyperplane that maximizes the margin

# Support Vector machine

- SVMs maximize the margin around the separating hyperplane.
  - ▶ A.k.a. large margin classifiers
- The decision function is fully specified by a subset of training samples, the support vectors
- Solving SVMs is a quadratic programming problem



# Why might predictions be wrong?

- True non-determinism
- Partial observability
  - ▶ hard, soft
- Representational bias
- Algorithmic bias
- Bounded resources

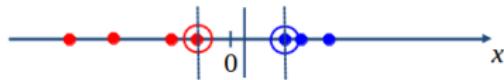
# Support Vector machine

- Datasets that are linearly separable works great

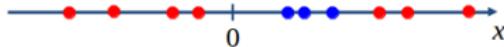


# Support Vector machine

- Datasets that are linearly separable works great

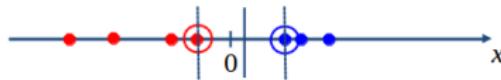


- But what are we going to do if the dataset is just too hard?

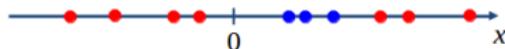


# Support Vector machine

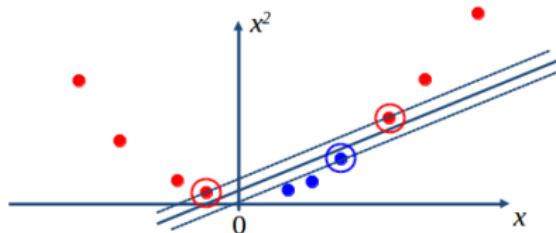
- Datasets that are linearly separable works great



- But what are we going to do if the dataset is just too hard?

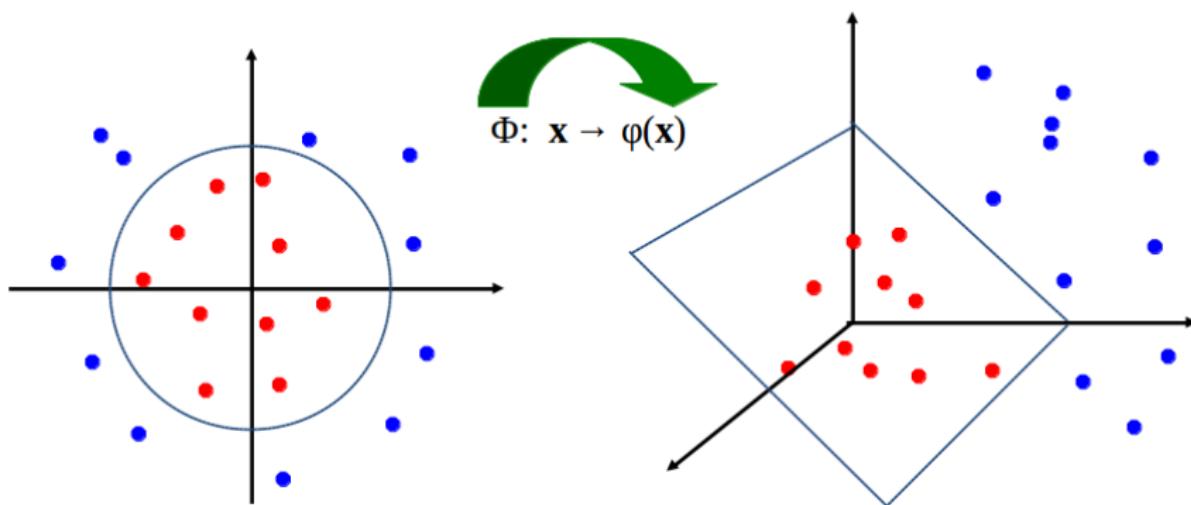


- How about mapping data to a higher-dimensional space



# Non-linear SVMs: Feature spaces

- General idea: the original feature space can always be mapped to some higher-dimensional feature space where the training set is separable:



<https://www.youtube.com/watch?v=9NrALgHFwTo>

# The “Kernel Trick”

- The linear classifier relies on an inner product between vectors  
 $K(x_i, x_j) = (x_i)^T (x_j)$
- If every datapoint is mapped into high-dimensional space via some transformation  $\Phi : x \rightarrow \varphi(x)$ , the inner product becomes:  
 $K(x_i, x_j) = \varphi(x_i)^T \varphi(x_j)$
- A kernel function is some function that corresponds to an inner product in some expanded feature space
  - ▶ Make non-separable problem separable.
  - ▶ Map data into better representational space
- Common Kernels
  - ▶ Linear
  - ▶ polynomial  $K(x, z) = (1 + x^T z)^d$
  - ▶ Radial basis function (infinite dimensional space)  
 $K(x_i, x_j) = e^{(-\|x_i - x_j\|^2)/2\sigma^2}$

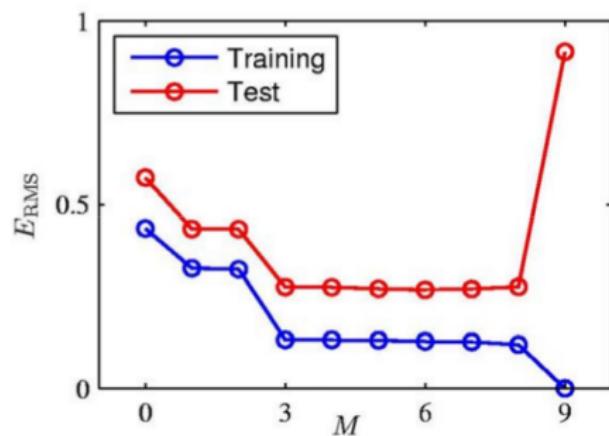
# SVM “Demo class work”

<https://cs.stanford.edu/~karpathy/svmjs/demo/>

# Regularization and Evaluation

# Bias - Variance Tradeoff

$$\text{MSE} = \text{Bias}^2 + \text{Variance}$$



# Regularization

In regression, when  $p \gg n$ , non-convex optimization. Regularization is introduced as a convex relaxation of the regression model.

$$\operatorname{argmin}_{\beta} \text{loss} + \text{penalty}(\beta)$$

# Regularization - $L_1$ and $L_2$ norms

- ① Lasso Penalty

$$\operatorname{argmin}_{\beta} \text{loss} + \lambda |\beta|$$

- ② Ridge Penalty

$$\operatorname{argmin}_{\beta} \text{loss} + \lambda |\beta|_2$$

How to choose  $\lambda$ ?

# Regularization - $L_1$ and $L_2$ norms

① Lasso Penalty

$$\operatorname{argmin}_{\beta} \text{loss} + \lambda |\beta|$$

② Ridge Penalty

$$\operatorname{argmin}_{\beta} \text{loss} + \lambda |\beta|_2$$

How to choose  $\lambda$ ?

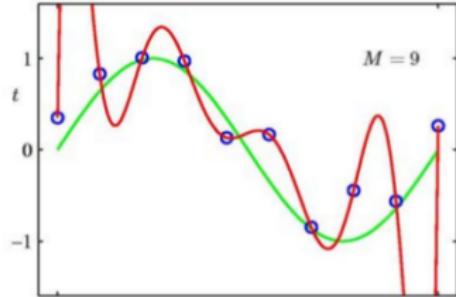
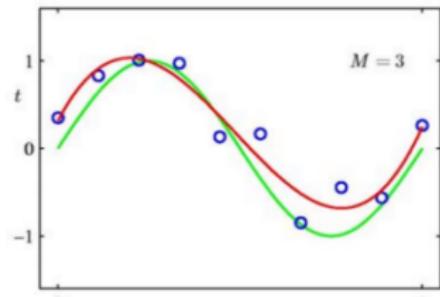
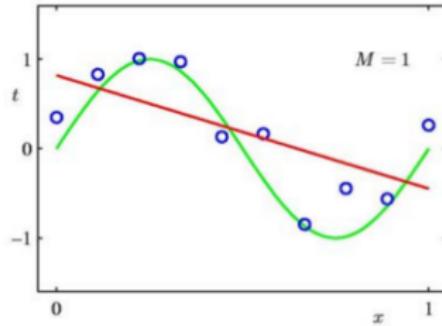
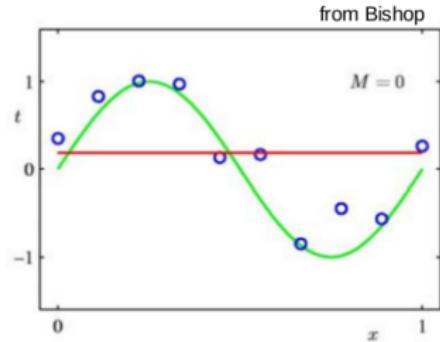
# Performance metrics

Truth table (confusion matrix)

		Predicted class	
		<i>P</i>	<i>N</i>
<b>Actual Class</b>	<i>P</i>	True Positives (TP)	False Negatives (FN)
	<i>N</i>	False Positives (FP)	True Negatives (TN)

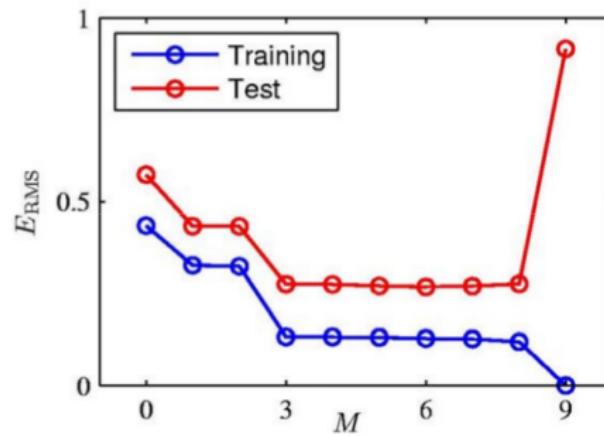
- Accuracy =  $\frac{TP+TN}{P+N}$
- Sensitivity =  $\frac{TP}{\text{Actual Positives}}$
- Specificity =  $\frac{TN}{\text{Actual Negatives}}$
- ...

# Some fits to the data: which is best?



# Bias-variance trade-off

- test data: a different sample from the same true function

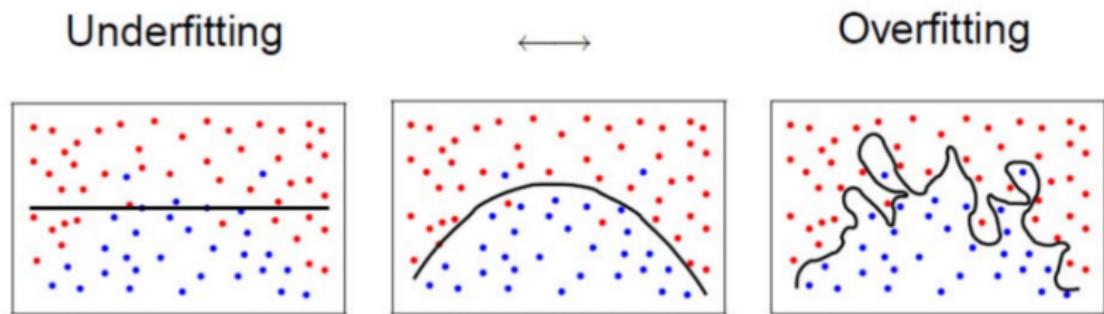


RootMeanSquare (RMS) Error

- training error goes to zero, but test error increases with M

# Generalization Problem in Regression/Classification

- Best fit will be somewhere in between under-fitting and over-fitting



# Cross-validation

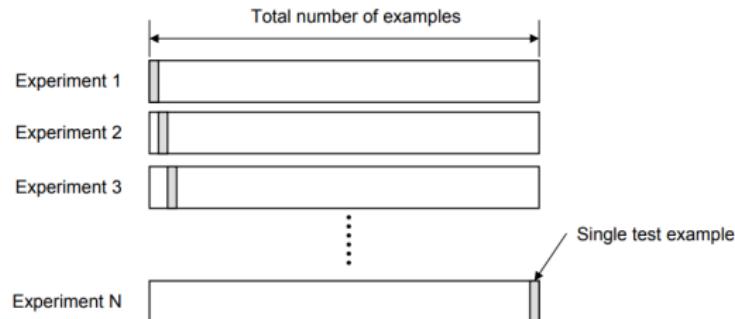


Figure: Leave-one-out cross-validation

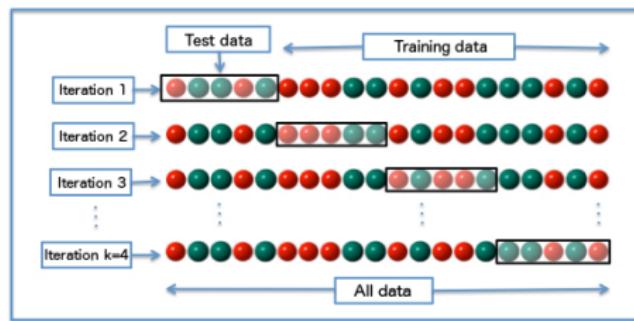


Figure: 5-fold cross-validation

# Cross-validation

- 10-fold cross-validation ?

# Cross-validation

- 10-fold cross-validation ?

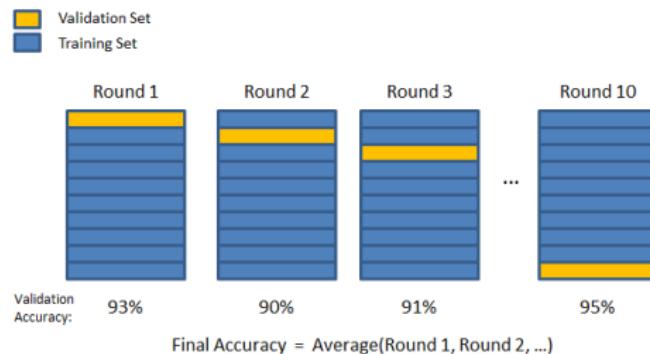


Figure: 10-fold cross-validation

# Cross-validation

- 10-fold cross-validation ?

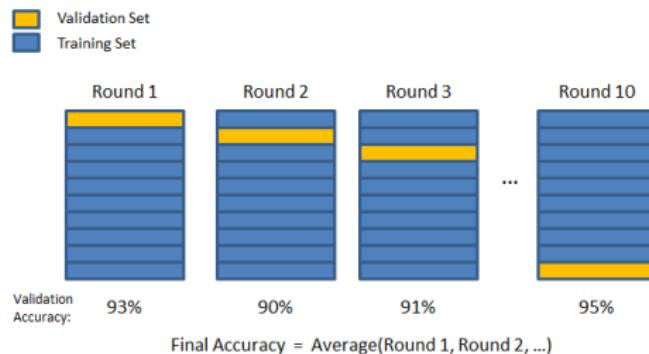
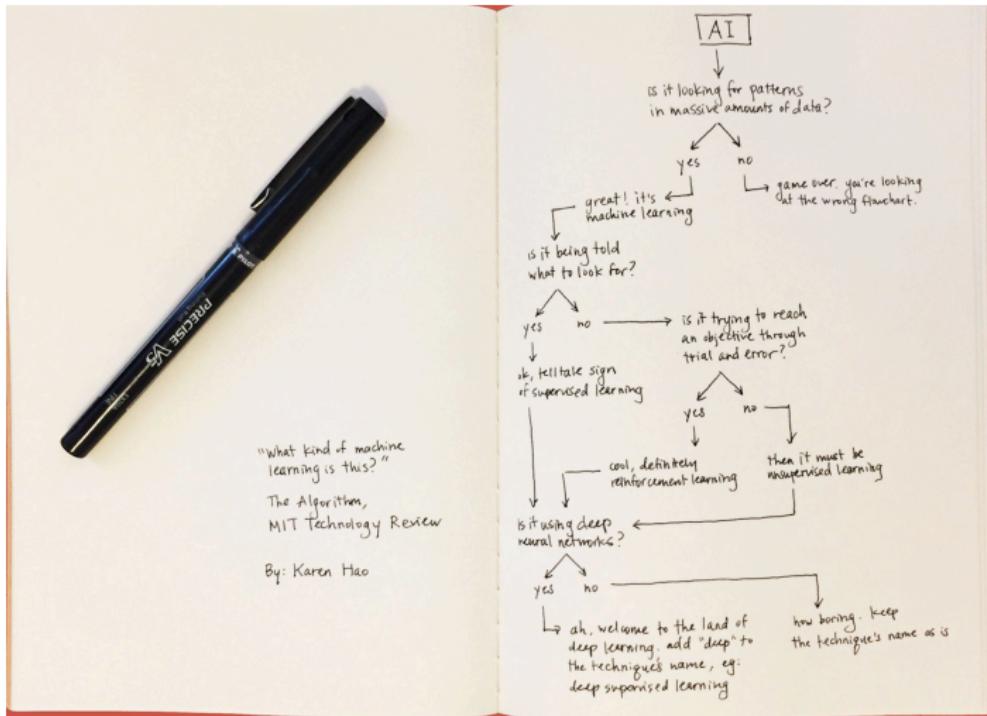


Figure: 10-fold cross-validation

- 5/10-repeated 5/10-fold cross-validation ?

The only thing that differs is that you randomly divide the data into 5/10 folds a different way each time



# References

COMP24111 Machine Learning

Machine Learning 10-701, Tom Mitchell

The Elements of Statistical Learning, (Friedman, Hastie and Tibshirani)

# Slide credit

Pedro Domingos

E. Alpaydin

Sebastian Thruns multimedia website

Geoffrey Hinton

Tim Oates