

# Machine Learning

Lecture 1: Introduction to ML, Statistical learning: concepts,  
Bayes classifier, LDA, QDA

**Asst. Prof. Dr. Santitham Prom-on**

Department of Computer Engineering, Faculty of Engineering  
King Mongkut's University of Technology Thonburi

# Teaching Team



Aj. Yai



Aj. Por



Aj. Joe



P' Noina

Lecturers

TA

# Course Learning Outcome

CLO1: Demonstrate mastery in concepts and details of supervised learning algorithms.

CLO2: Demonstrate mastery in concepts and details of unsupervised and reinforcement learning algorithms.

CLO3: Implement machine learning models and workflows.

# Grading

Midterm Exam 30%

Final Exam 35%

Homework and assignments 35%

A:  $\geq 85$ , B+: [84,80], B: [79,75], C+: [74,65],  
C: [64,55], D+: [54,50], D: [49,45], F: < 45

The instructor reserves the right to change the grading policy as deemed appropriate.

# References

- G. James, et.al., Introduction to Statistical Learning with Application with Python, 2023
- A. Geron, Hands-On Machine Learning with Scikit-Learn & Tensorflow, O'Reilly, 2019

# Class Policies

- Assignments are due in one week before class in LEB2.
- Late submissions are only accepted under reasonable excuses and explicit permission from the instructor. No submission is accepted after the solution has been posted.
- Posted solutions will be brief and does not show routine works. You should attempt to work out detailed solutions on your own.
- Academic integrity is strictly enforced.

# Class Schedule

Date	Topic
Aug 7	Introduction to ML, Statistical learning: concepts, Bayes classifier, LDA, QDA
Aug 14	Training models: Direct (OLS) and iterative (gradient descent) approaches
Aug 21	Support vector machines: Linear and nonlinear SVM
Aug 28	Regression, Constraints, Generalized linear models
Sep 4	Tree-based models: Decision tree and Regression tree
Sep 9 - 13	1st University Exam Period - No class
Sep 18	Ensemble models
Sep 25	MLOps
Oct 2	Neural networks and deep learning
Oct 9	Convolutional neural networks
Oct 16	Dimensionality reduction
Oct 21 - 29	2nd University Exam Period - No class
Nov 6	Clustering (1): Dissimilarity measures, K-means clustering, Cluster evaluation
Nov 13	Clustering (2): Gaussian mixture model, Hierarchical clustering, Density-based clustering
Nov 20	Anomaly detection
Nov 27	Reinforcement learning
Dec 2 - 13	3rd University Exam Period - No class

# Topics

- Introduction to ML
- Statistical learning concepts
- Bayes' classifier
- Discriminant analysis

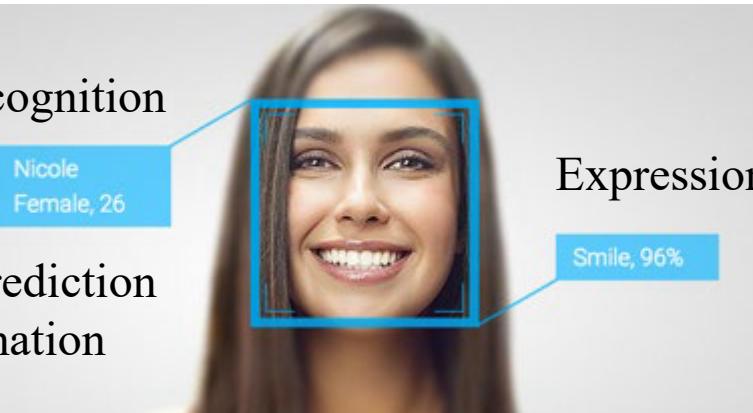
# Introduction to ML Statistical Learning Concepts

Section 1

# Consumer ML

Face Detection

Face Recognition



Expression Recognition

Smile, 96%

Gender Prediction

Age Estimation



Voice Interface



Smart Home

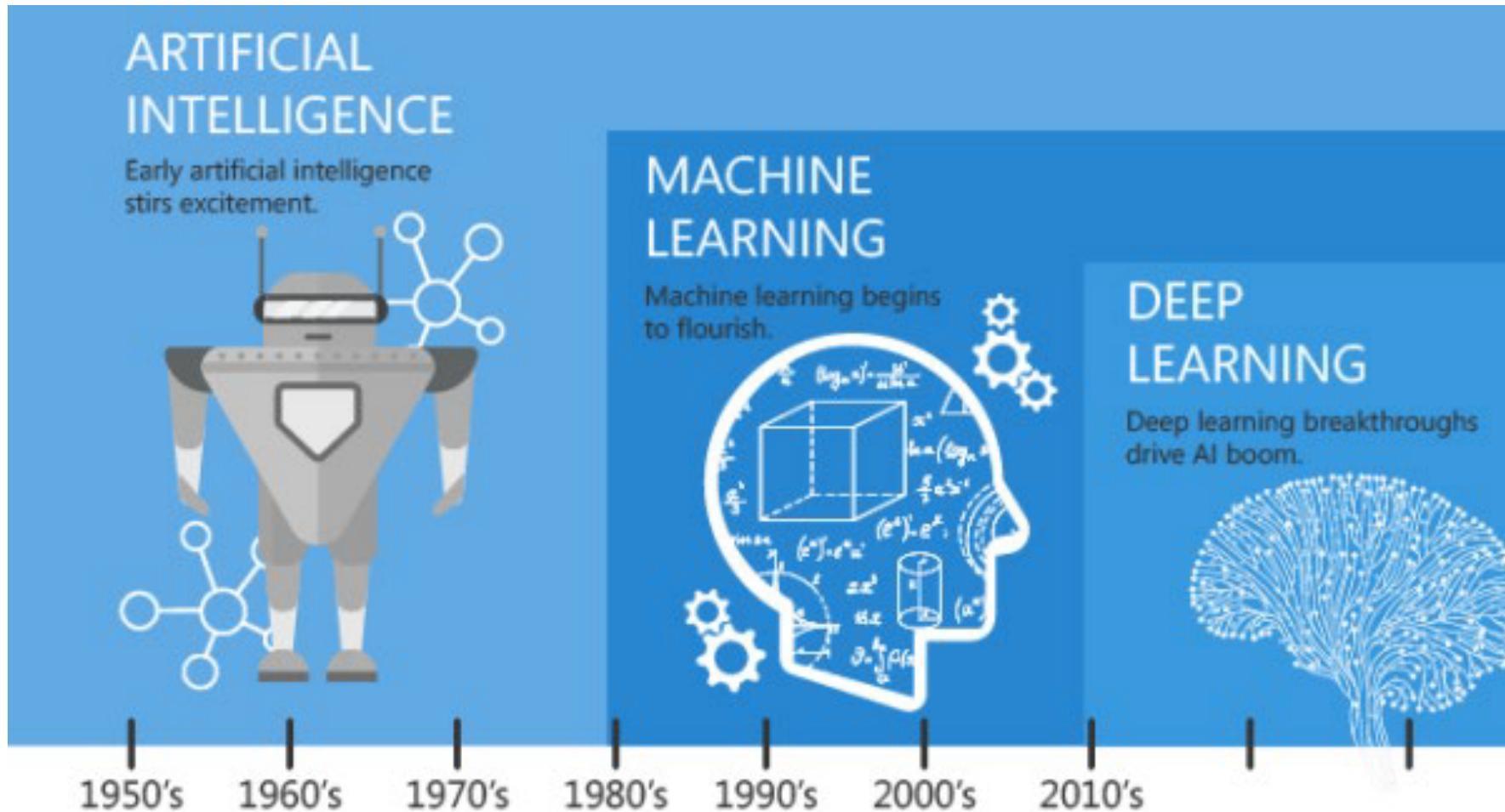
# What exactly is machine learning?

- Machine learning (ML) is a method of data analysis that automates analytical model building.
- Using algorithms that iteratively learn from data, machine learning allows computers to find hidden insights without being explicitly programmed where to look.

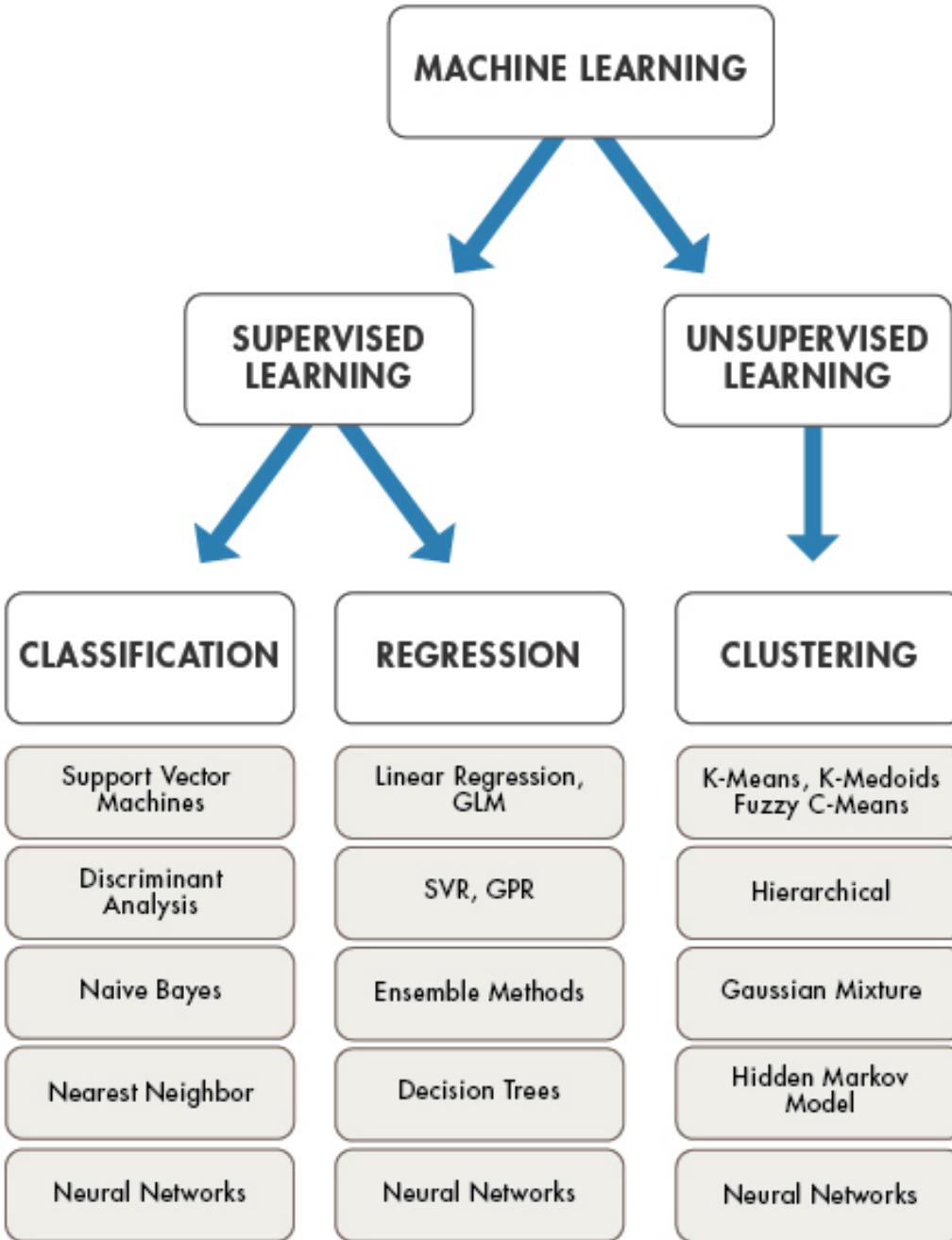
# Why machine learning matter

- Automatic: Train it once and it can be run automatically
- Fast: With big data, work faster than human
- Accurate: Can predict groups more accurately than manual methods
- Scale: Able to handle large data

# Timeline

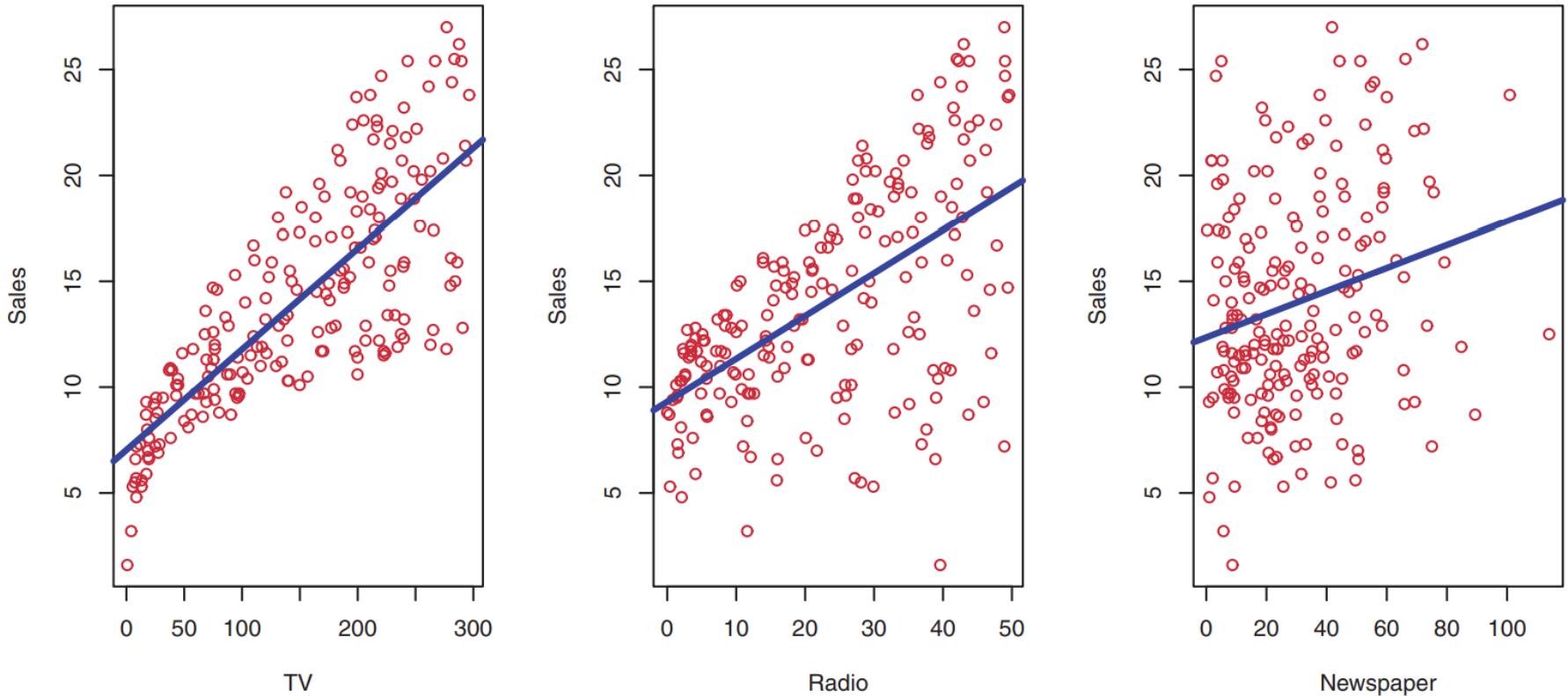


# Types of ML



# Style of Learning

<ul style="list-style-type: none"><li>• Data has <b>known labels</b> or output</li></ul>	<ul style="list-style-type: none"><li>• Labels or output unknown</li><li>• Focus on <b>finding patterns and gaining insight</b> from the data</li></ul>	<ul style="list-style-type: none"><li>• Labels or output known for <b>a subset of data</b></li><li>• A blend of supervised and unsupervised learning</li></ul>	<ul style="list-style-type: none"><li>• Focus on <b>making decisions</b> based on previous experience</li><li>• Policy-making with feedback</li></ul>
<ul style="list-style-type: none"><li>• Insurance underwriting</li><li>• Fraud detection</li></ul>	<ul style="list-style-type: none"><li>• Customer clustering</li><li>• Association rule mining</li></ul>	<ul style="list-style-type: none"><li>• Medical predictions (where tests and expert diagnoses are expensive, and only part of the population is tested)</li></ul>	<ul style="list-style-type: none"><li>• Game AI</li><li>• Complex decision problems</li><li>• Reward systems</li></ul>

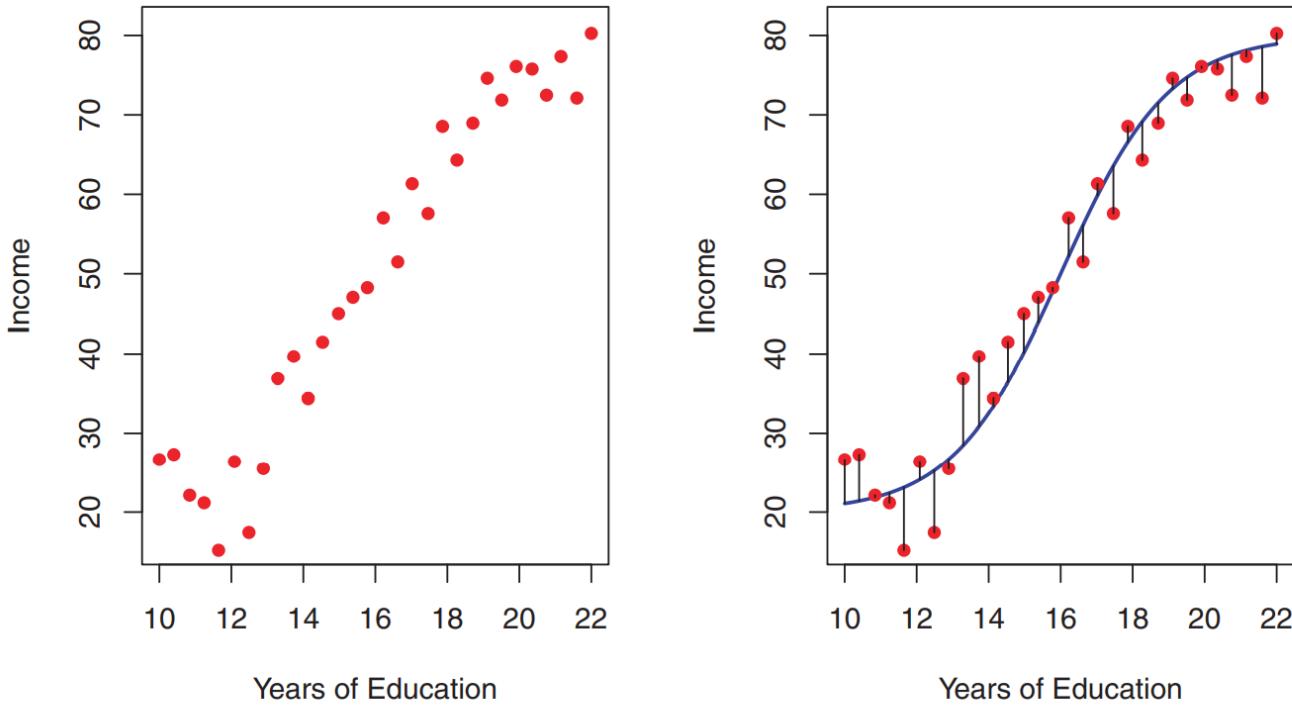


**FIGURE 2.1.** The Advertising data set. The plot displays **sales**, in thousands of units, as a function of **TV**, **radio**, and **newspaper** budgets, in thousands of dollars, for 200 different markets. In each plot we show the simple least squares fit of **sales** to that variable, as described in Chapter 3. In other words, each blue line represents a simple model that can be used to predict **sales** using **TV**, **radio**, and **newspaper**, respectively.

# Representing relationship

- Suppose that we observe a quantitative response  $Y$  and  $p$  different predictors,  $X_1, X_2, \dots, X_p$ .
- Assume that there is some relationship between  $Y$  and  $X = (X_1, X_2, \dots, X_p)$ ,

$$Y = f(X) + \epsilon.$$



**FIGURE 2.2.** The `Income` data set. Left: The red dots are the observed values of `income` (in tens of thousands of dollars) and `years of education` for 30 individuals. Right: The blue curve represents the true underlying relationship between `income` and `years of education`, which is generally unknown (but is known in this case because the data were simulated). The black lines represent the error associated with each observation. Note that some errors are positive (if an observation lies above the blue curve) and some are negative (if an observation lies below the curve). Overall, these errors have approximately mean zero.

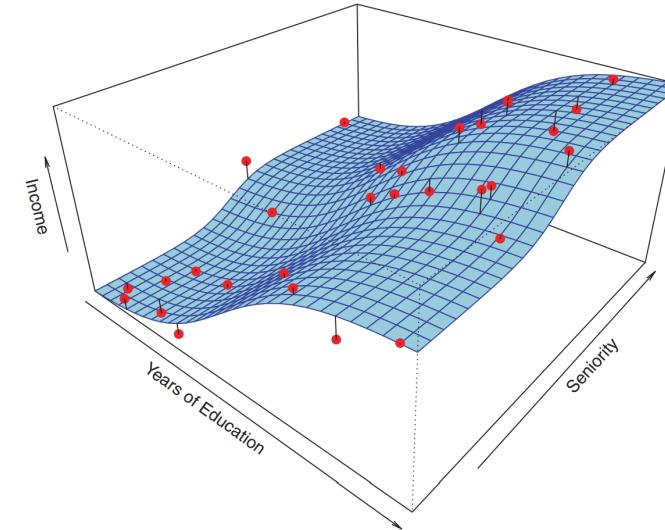
# Why estimate $f$ ?

- Prediction

$$\hat{Y} = \hat{f}(X),$$

- Inference

- Which predictors are associated with the response?
- What is the relationship between the response and each predictor?
- Can the relationship between  $Y$  and each predictor be adequately summarized using a linear equation, or is the relationship more complicated?



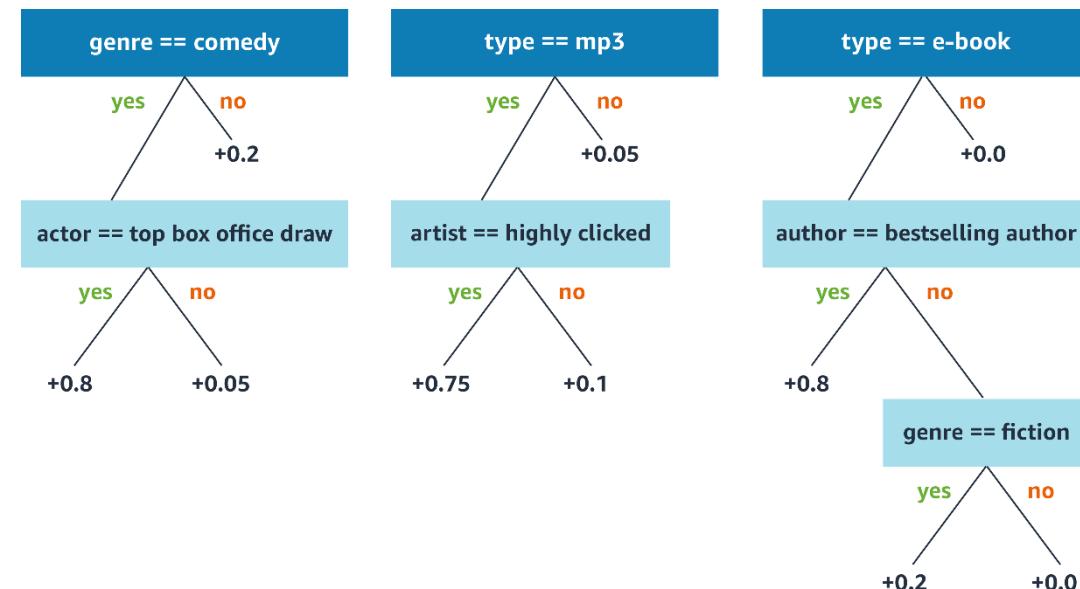
**FIGURE 2.3.** The plot displays income as a function of years of education and seniority in the Income data set. The blue surface represents the true underlying relationship between income and years of education and seniority, which is known since the data are simulated. The red dots indicate the observed values of these quantities for 30 individuals.

# How do we estimate $f$ ?

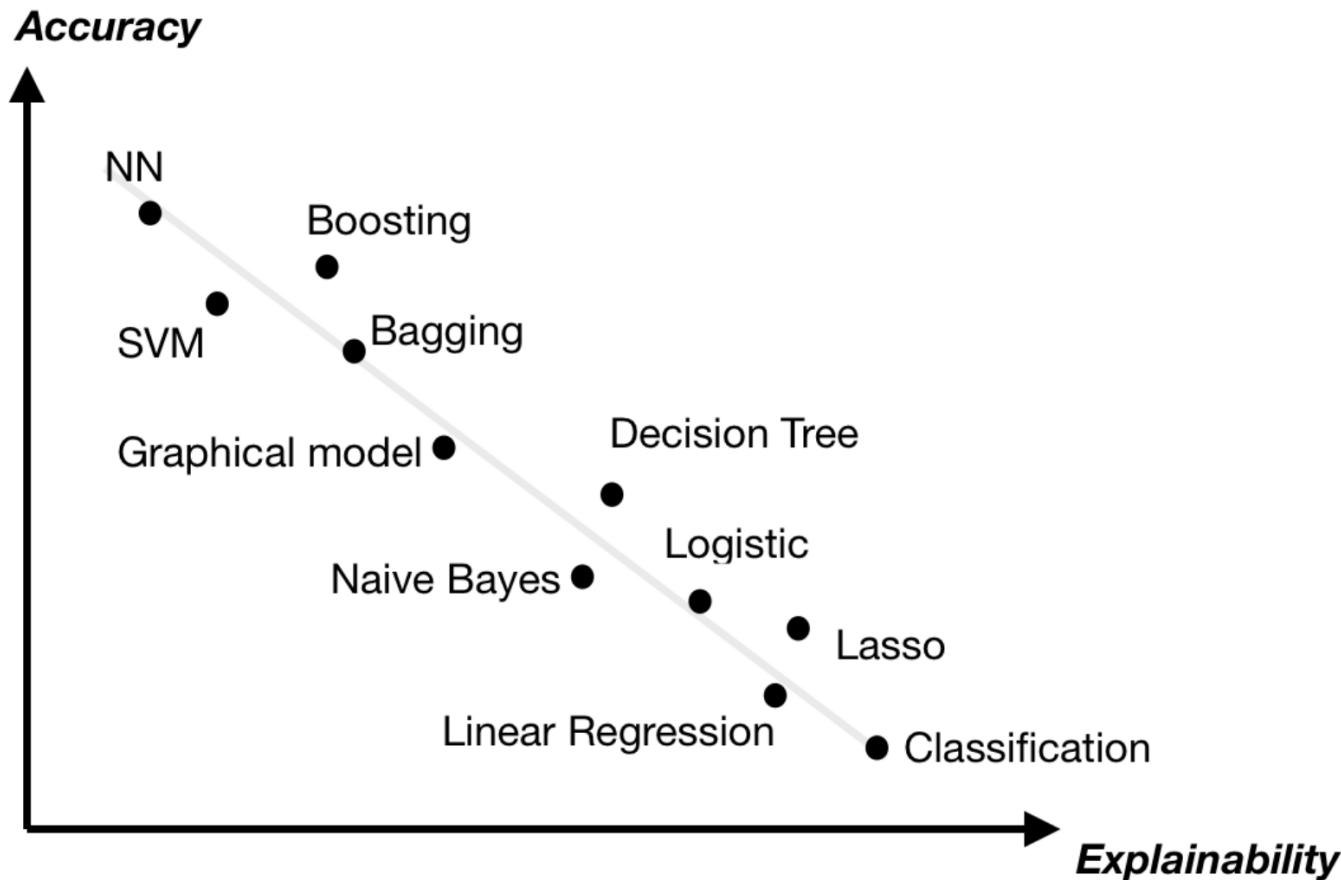
- Linear models

$$f(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p.$$

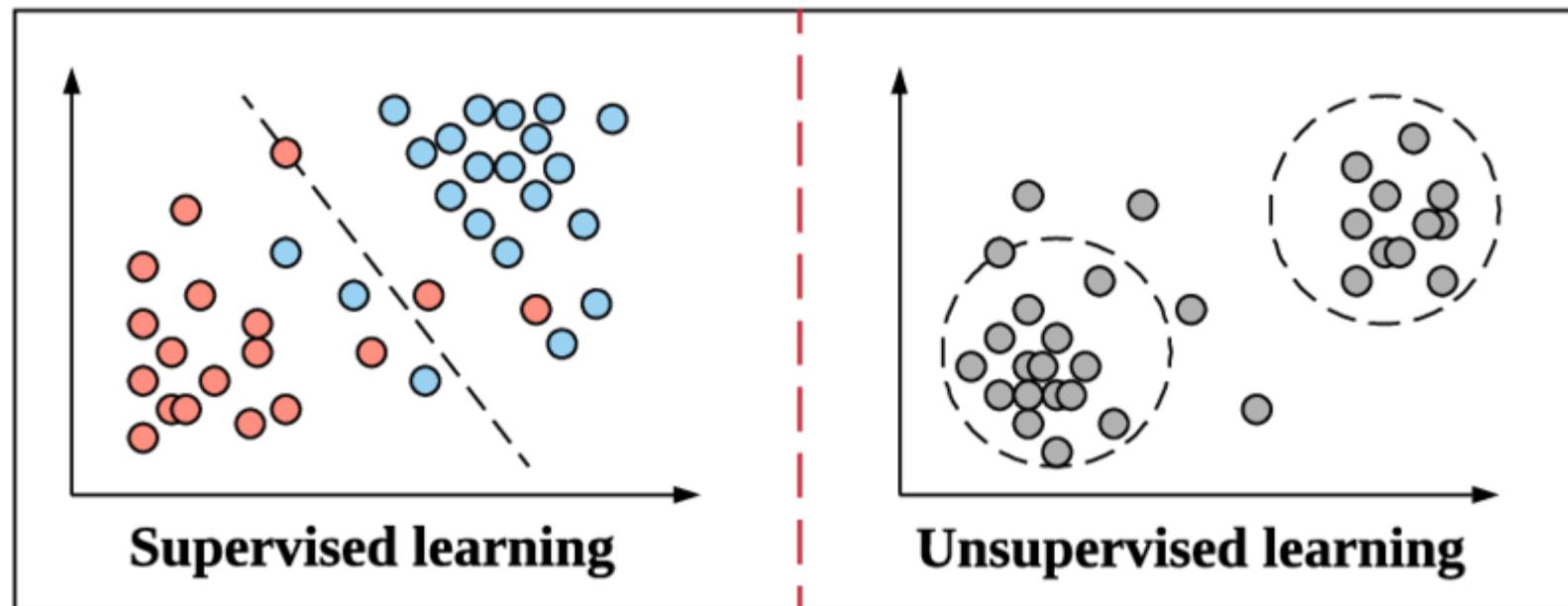
- Non-linear models



# Explainability vs Accuracy



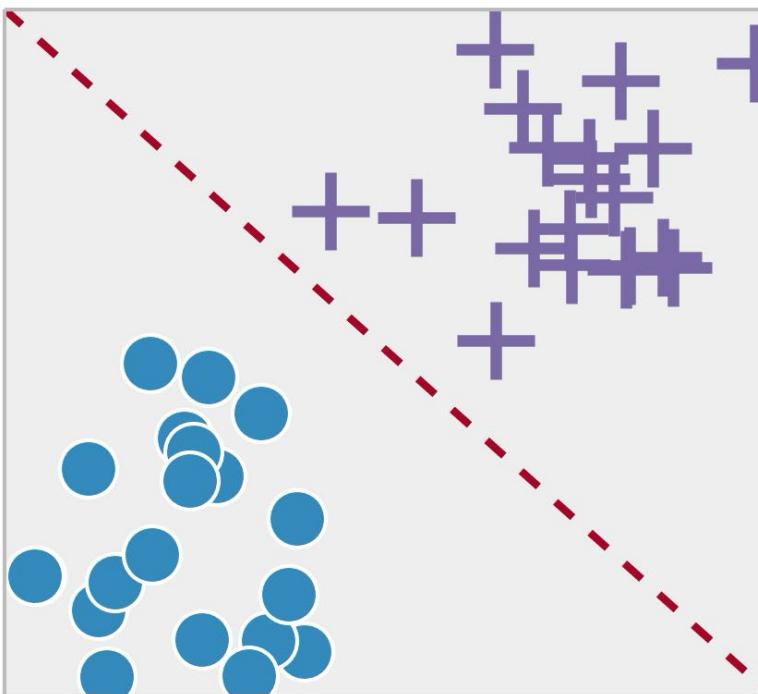
# Supervised vs Unsupervised



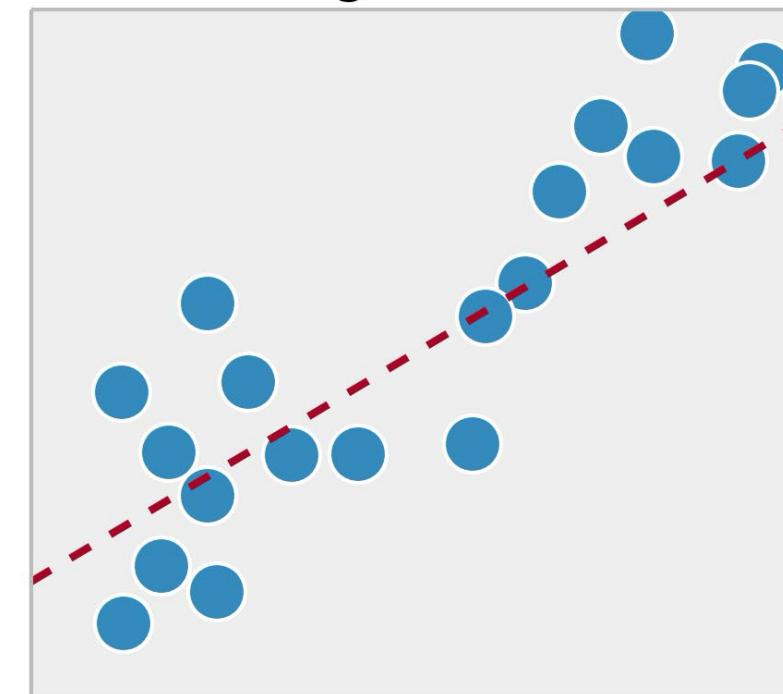
Qian, Bin & Su, Jie & Wen, Zhenyu & Yang, Renyu & Zomaya, Albert & Rana, Omer. (2019). Orchestrating the Development Lifecycle of Machine Learning-Based IoT Applications: A Taxonomy and Survey.

# Regression vs Classification

Classification

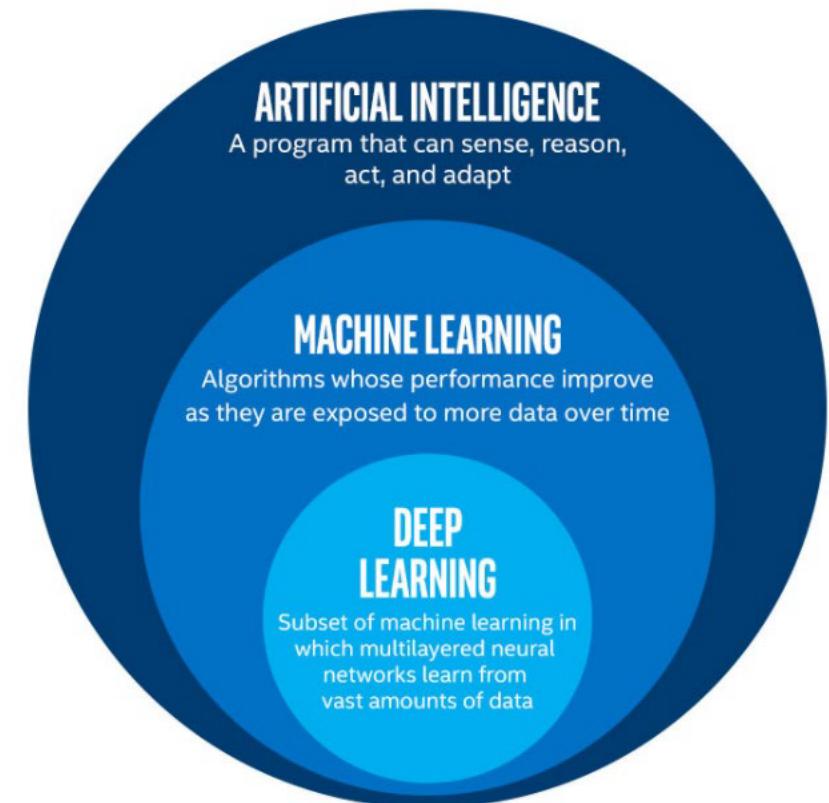


Regression



# Machine Learning Basics

- Deep learning is a specific kind of machine learning.
- To understand deep learning well, one must have a solid understanding of the basic principles of machine learning.



# Learning Algorithm - Definition

“A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .” (Mitchell, 1997)

$$\begin{array}{ccc} T & \longrightarrow & P \\ \uparrow & & \\ E & & \end{array}$$

# Tasks, T

## Classification

- In this type of task, the computer program is asked to specify which of  $k$  categories some input belongs to.
- To solve this task, the learning algorithm is usually asked to produce a function  $f: \mathbb{R}^n \rightarrow \{1, \dots, k\}$ .
- When  $y = f(x)$ , the model assigns an input described by vector  $x$  to a category identified by numeric code  $y$ .

# Tasks, T

## Regression

- In this type of task, the computer program is asked to predict a numerical value given some input.
- To solve this task, the learning algorithm is asked to output a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ .
- This type of task is similar to classification, except that the format of output is a value instead of a class.

# Tasks, T

## Clustering

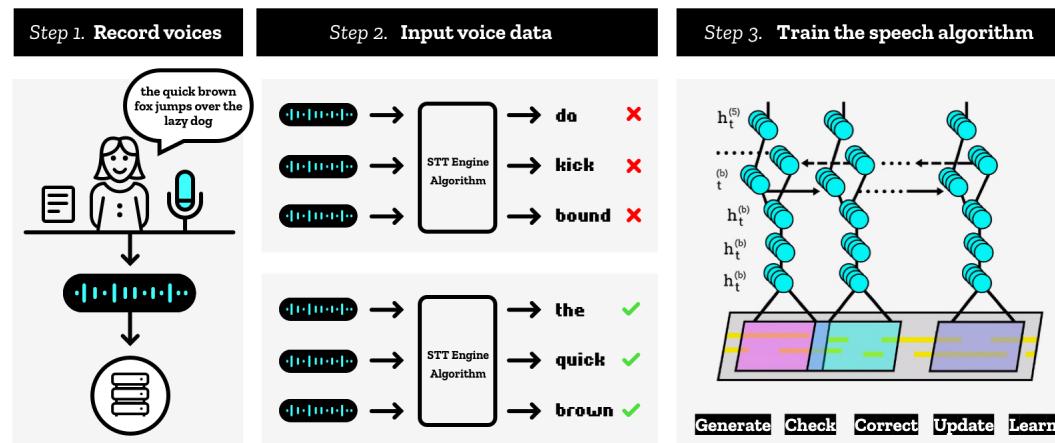
- In this type of task, the computer program is asked to divide data into groups.
- To solve this task, the learning algorithm is asked to output a function  $f: \mathbb{R}^n \rightarrow \{1, \dots, k\}$ , but without the label.
- The goal is usually to divide data into groups based on specific features, e.g. Recency-Frequency-Monetary customer segmentation

# Tasks, T

## Transcription

- In this type of task, the machine learning system is asked to observe a relatively unstructured representation of some kind of data and transcribe the information into discrete textual form.

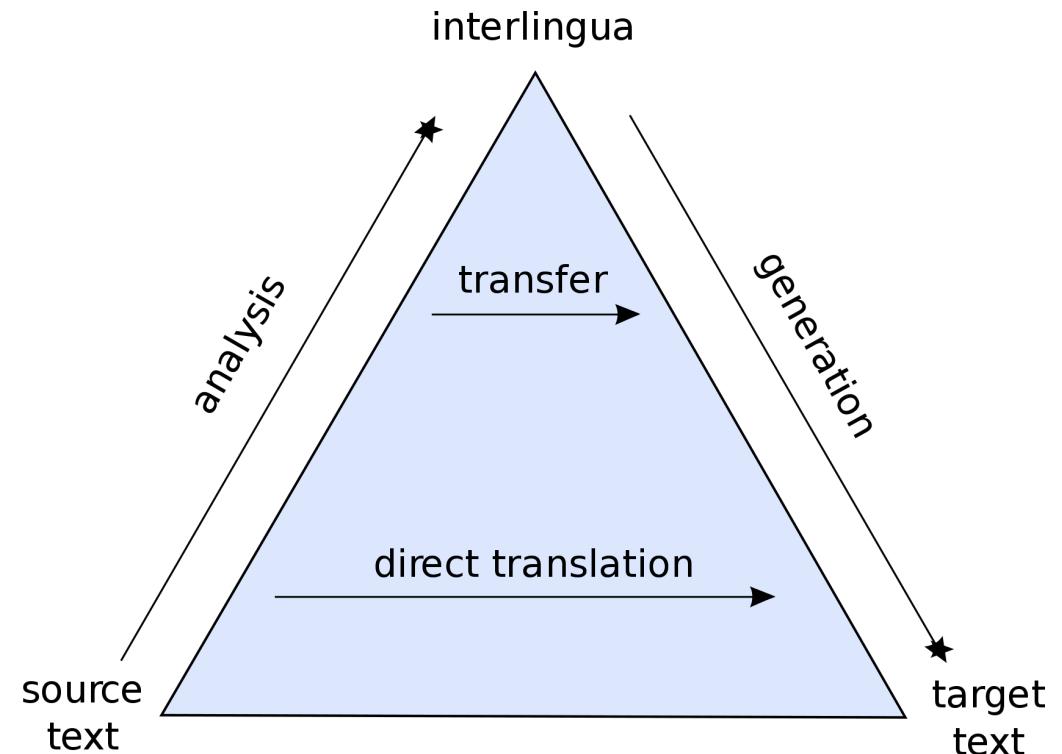
### How a Speech Application Learns



# Tasks, T

## Machine Translation

- In a machine translation task, the input already consists of a sequence of symbols in some language, and the computer program must convert this into a sequence of symbols in another language.



# Tasks, T

## Anomaly Detection

- In this type of task, the computer program sifts through a set of events or objects and flags some of them as being unusual or atypical.



# Tasks, T

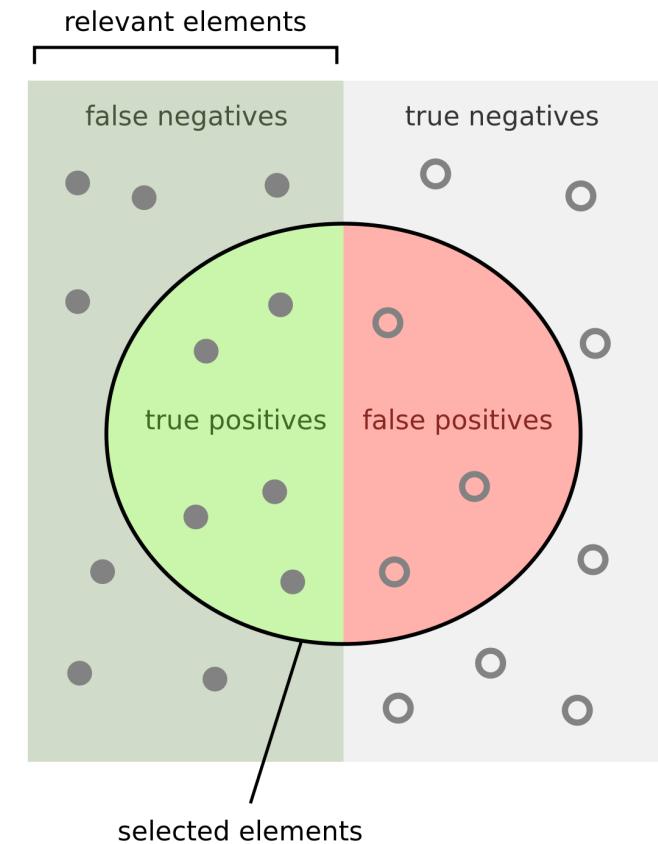
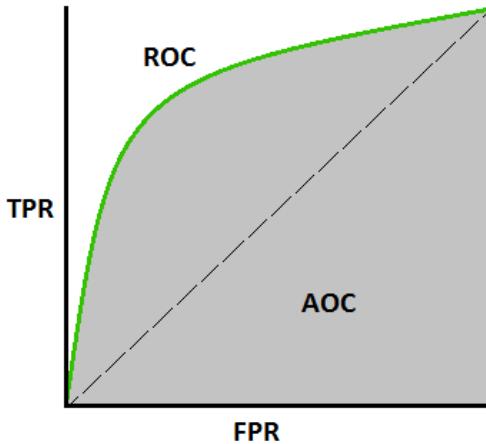
## Synthesis and sampling

- In this type of task, the machine learning algorithm is asked to generate new examples that are similar to those in the training data.



# Performance, P

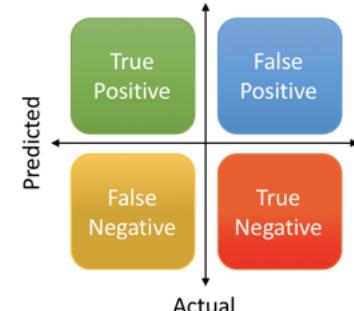
- Classification



$$\text{Precision} = \frac{\text{True Positive}}{\text{Actual Results}} \quad \text{or} \quad \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

$$\text{Recall} = \frac{\text{True Positive}}{\text{Predicted Results}} \quad \text{or} \quad \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

$$\text{Accuracy} = \frac{\text{True Positive} + \text{True Negative}}{\text{Total}}$$



How many selected items are relevant?

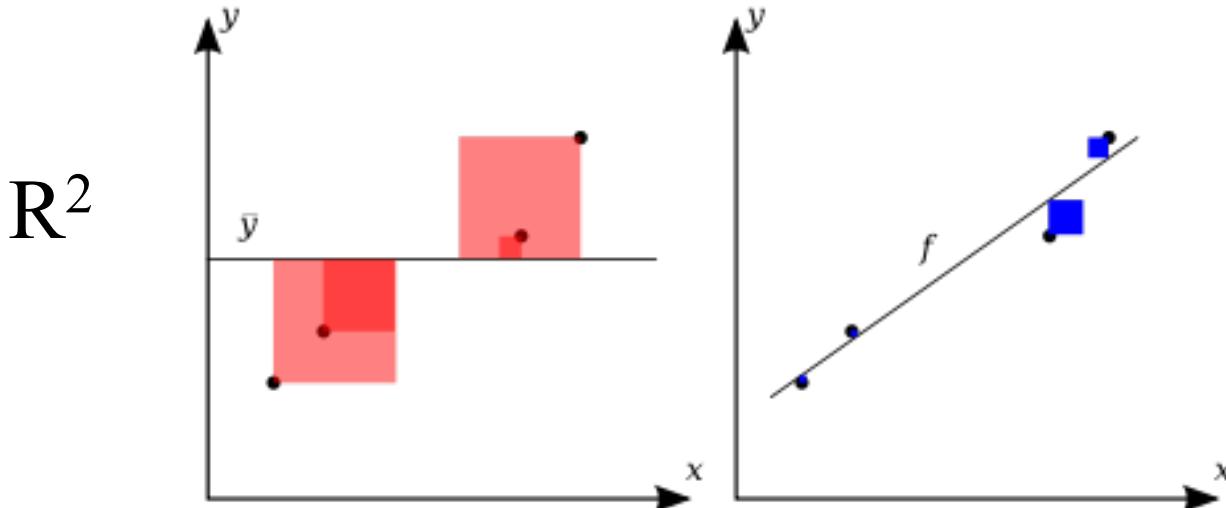
$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

# Performance, P

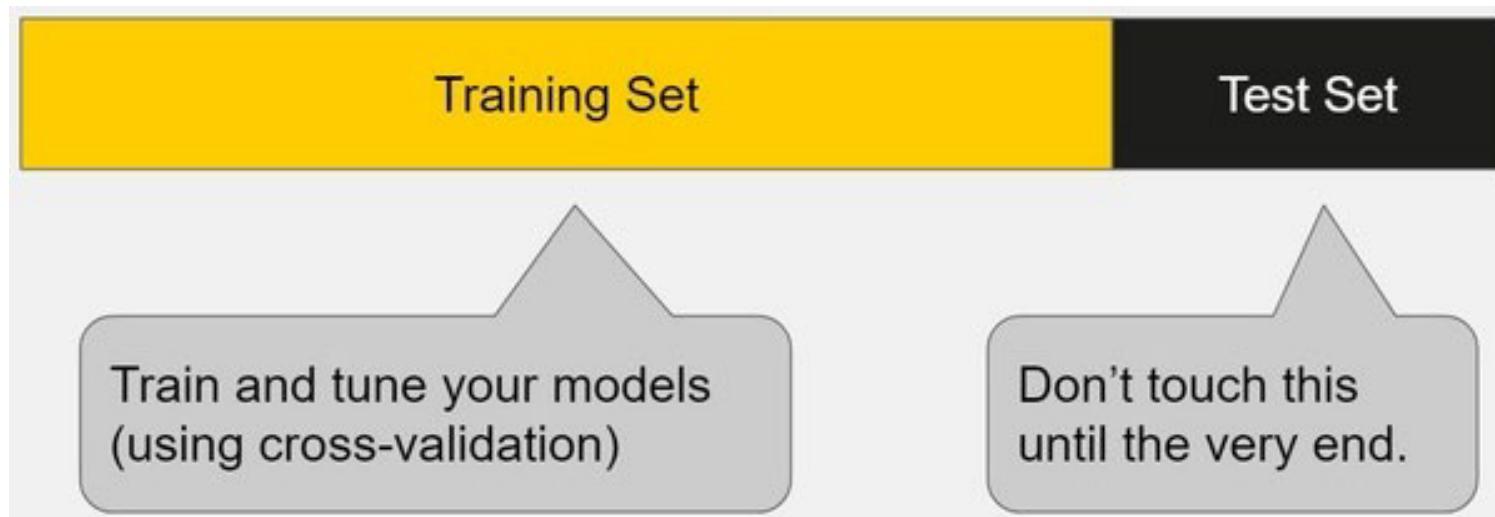
- Regression



$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |x_i - \hat{x}|$$

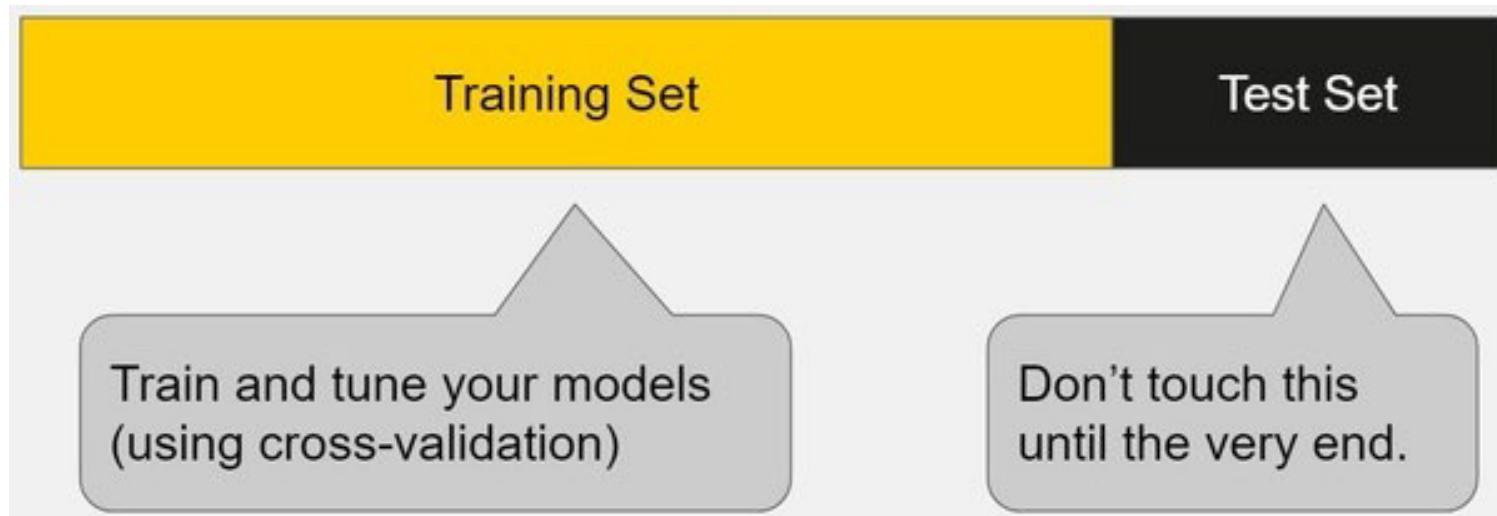
# Performance, P

- Train set and Test set

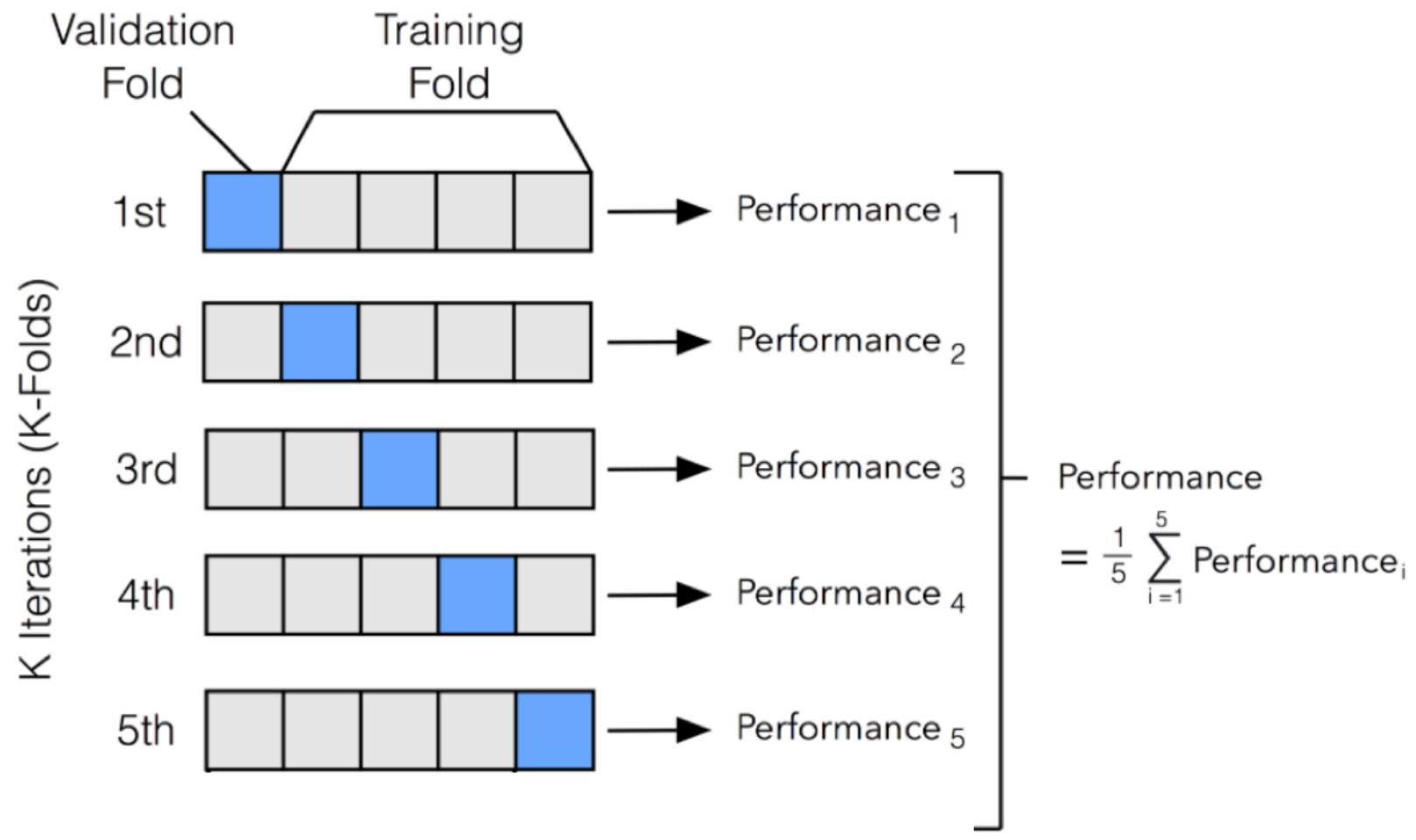


# Experience, E

- Train set and Test set



# Cross-validation



# Example: Linear Regression

- Linear regression solves a regression problem.
- In other words, the goal is to build a system that can take a vector  $\mathbf{x} \in \mathbb{R}^n$  as input and predict the value of a scalar  $y \in \mathbb{R}$  as its output.
- The output of linear regression is a linear function of the input.
- Let  $\hat{y}$  be the value that our model predicts  $y$  should take on. We define the output to be

$$\hat{y} = \mathbf{w}^\top \mathbf{x},$$

where  $\mathbf{w} \in \mathbb{R}^n$  is a vector of parameters

# Linear regression: Task, T

To predict  $y$  from  $x$  by outputting

$$\hat{y} = \mathbf{w}^\top \mathbf{x},$$

# Performance measure, P

Mean square error

$$\text{MSE}_{\text{test}} = \frac{1}{m} \sum_i (\hat{\mathbf{y}}^{(\text{test})} - \mathbf{y}^{(\text{test})})_i^2.$$

# Experience, E

$$\nabla_{\mathbf{w}} \text{MSE}_{\text{train}} = 0$$

$$\Rightarrow \nabla_{\mathbf{w}} \frac{1}{m} \|\hat{\mathbf{y}}^{(\text{train})} - \mathbf{y}^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \frac{1}{m} \nabla_{\mathbf{w}} \|\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}\|_2^2 = 0$$

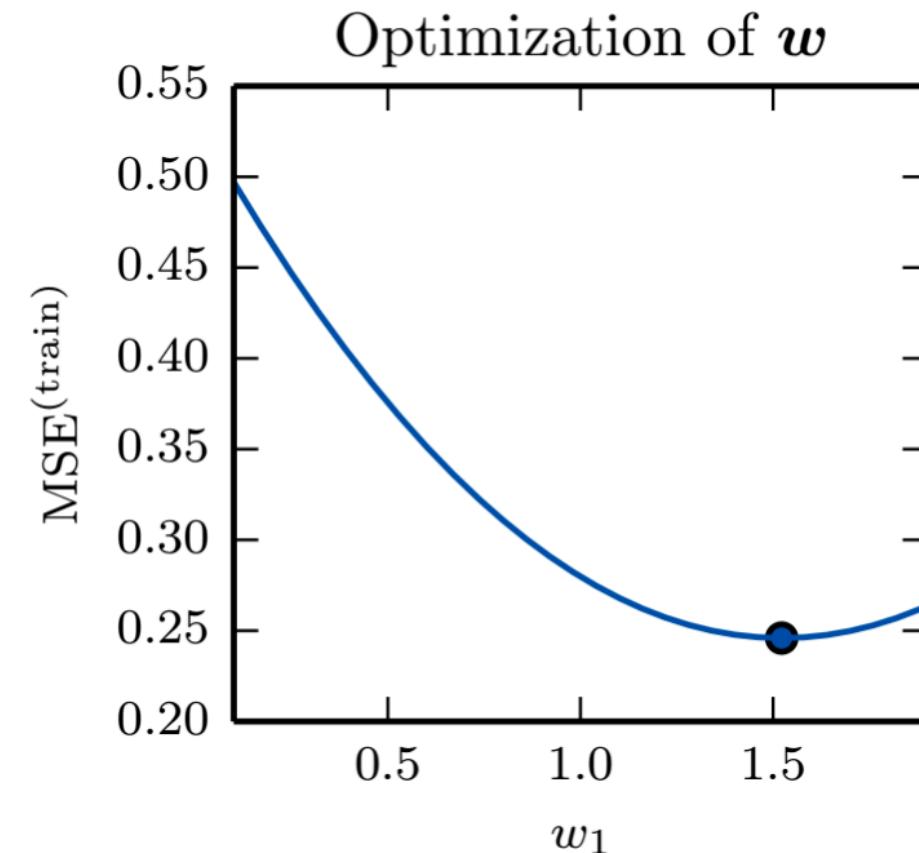
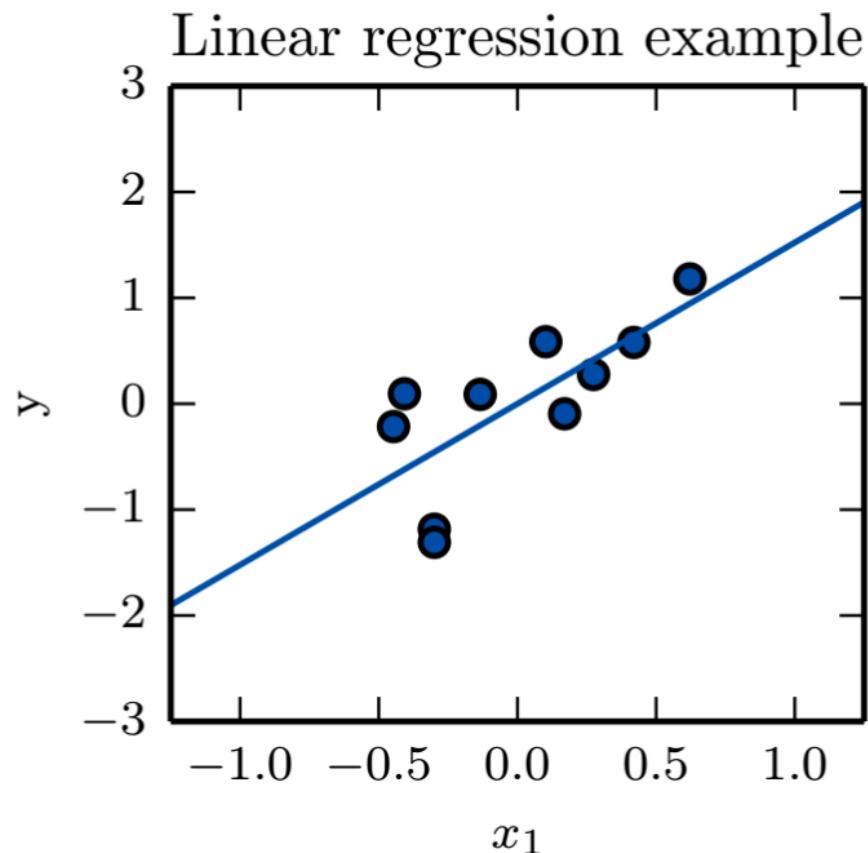
$$\Rightarrow \nabla_{\mathbf{w}} \left( \mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right)^\top \left( \mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right) = 0$$

$$\Rightarrow \nabla_{\mathbf{w}} \left( \mathbf{w}^\top \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \mathbf{w} - 2\mathbf{w}^\top \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} + \mathbf{y}^{(\text{train})\top} \mathbf{y}^{(\text{train})} \right) = 0$$

$$\Rightarrow 2\mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \mathbf{w} - 2\mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} = 0$$

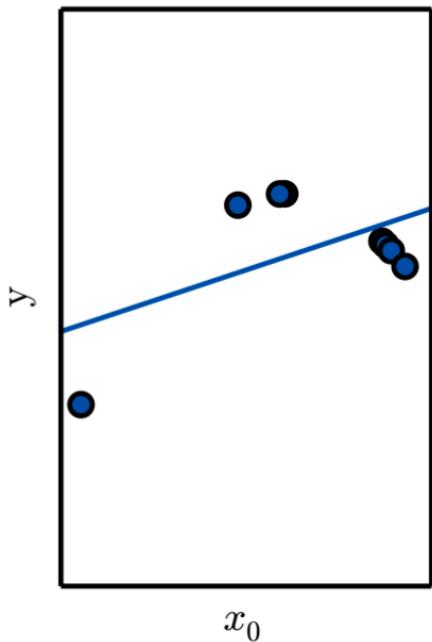
$$\Rightarrow \mathbf{w} = \left( \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \right)^{-1} \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})}$$

# Linear regression – learning weight



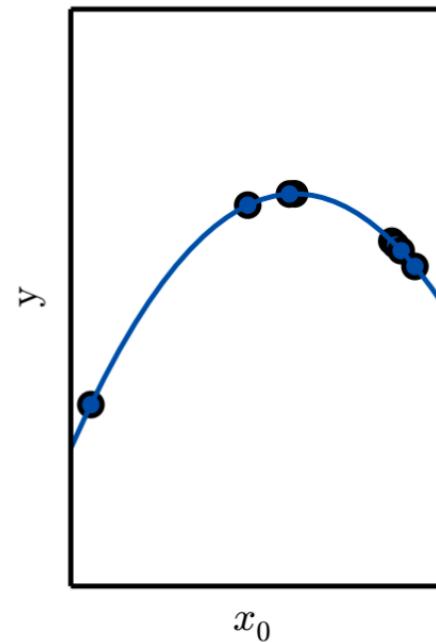
# Regression, underfitting and overfitting

Underfitting



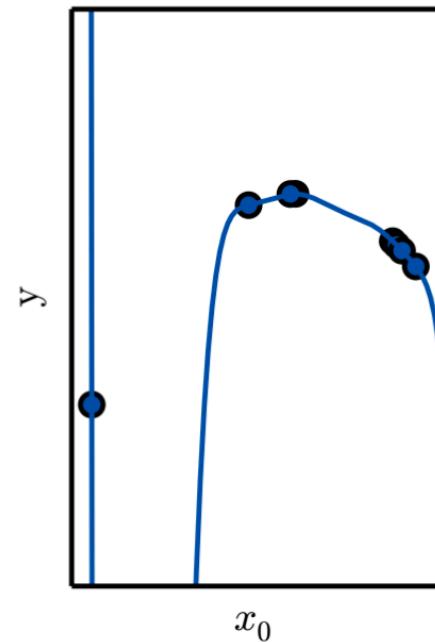
$$\hat{y} = b + wx.$$

Appropriate capacity



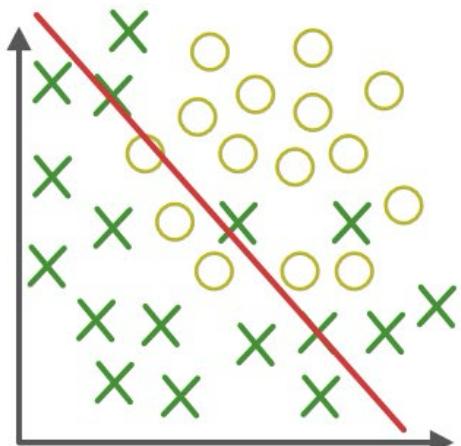
$$\hat{y} = b + w_1x + w_2x^2.$$

Overfitting



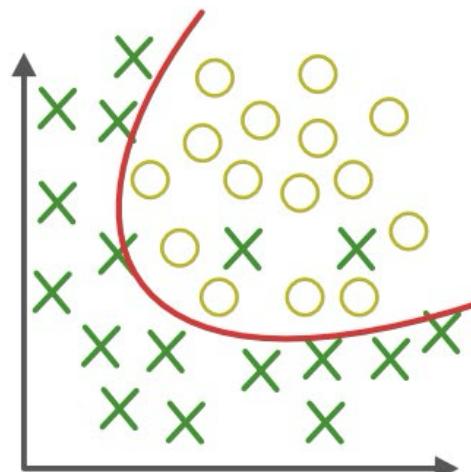
$$\hat{y} = b + \sum_{i=1}^9 w_i x^i.$$

# Classification, underfitting and overfitting

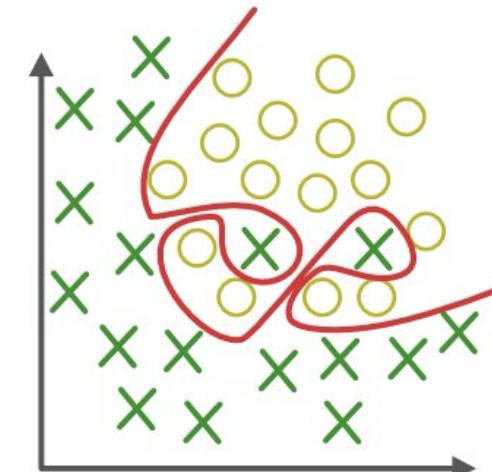


**Under-fitting**

(too simple to explain the variance)



**Appropriate-fitting**



**Over-fitting**

(forcefitting--too good to be true)

# Bayes' Classifier

Section 2

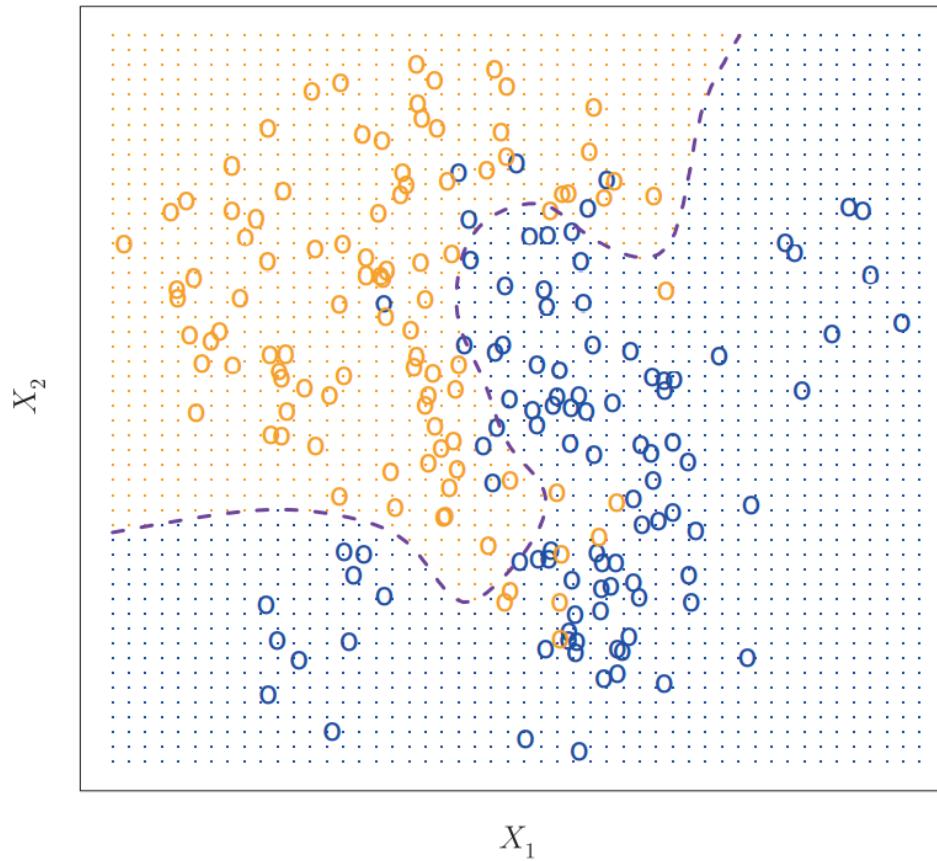
# Bayes' Classifier

- assigns each observation to the most likely class, given its predictor values.
- In other words, we should simply assign a test observation with predictor vector  $x_0$  to the class  $j$  for which

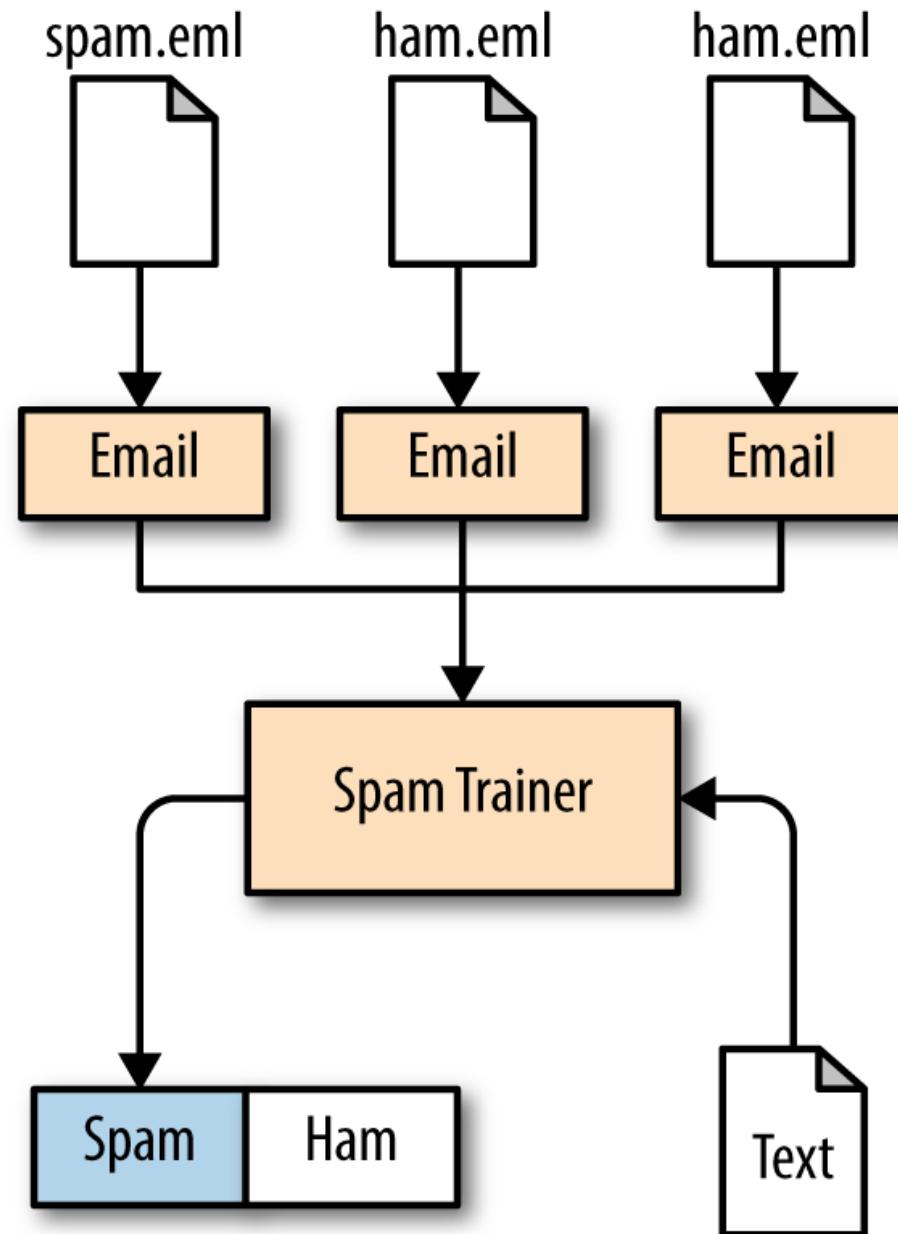
$$\Pr(Y = j \mid X = x_0)$$

is the largest

- This is a conditional probability, the probability conditional that  $Y = j$ , given the observed predictor vector  $x_0$ .



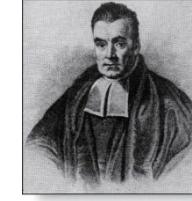
**FIGURE 2.13.** A simulated data set consisting of 100 observations in each of two groups, indicated in blue and in orange. The purple dashed line represents the Bayes decision boundary. The orange background grid indicates the region in which a test observation will be assigned to the orange class, and the blue background grid indicates the region in which a test observation will be assigned to the blue class.



# Motivation for Bayes' Theorem

- Bayes' theorem allows us to use probability to answer questions such as the following:
  - Given that someone tests positive for having a particular disease, what is the probability that they actually do have the disease?
  - Given that someone tests negative for the disease, what is the probability, that in fact they do have the disease?
- Bayes' theorem has applications to medicine, law, artificial intelligence, engineering, and many diverse other areas.

# Bayes' Theorem



Thomas Bayes  
(1702-1761)

**Bayes' Theorem:** Suppose that  $E$  and  $F$  are events from a sample space  $S$  such that  $p(E) \neq 0$  and  $p(F) \neq 0$ . Then:

$$p(F|E) = \frac{p(E|F)p(F)}{p(E|F)p(F) + p(E|\bar{F})p(\bar{F})}$$

$p(E)$

A blue wavy line underlines the term  $p(E)$  in the denominator of the equation.

# Derivation of Bayes' Theorem

- Recall the definition of the conditional probability  $p(E|F)$ :

$$p(E|F) = \frac{p(E \cap F)}{p(F)}$$

- From this definition, it follows that:

$$p(E|F) = \frac{p(E \cap F)}{p(F)}, \quad p(F|E) = \frac{p(E \cap F)}{p(E)}$$

# Derivation of Bayes' Theorem

$$p(E|F)p(F) = p(E \cap F), \quad p(F|E)p(E) = p(E \cap F)$$

Equating the two formulas for  $p(E|F)p(F)$  shows that

$$p(E|F)p(F) = p(F|E)p(E)$$

Solving for  $p(E|F)$  and for  $p(F|E)$  tells us that

$$p(E|F) = \frac{p(F|E)p(E)}{p(F)}, \quad p(F|E) = \frac{p(E|F)p(F)}{p(E)}$$

# Bayes' Theorem

$$p(F|E) = \frac{p(E|F)p(F)}{p(E)}$$

The diagram illustrates the components of Bayes' Theorem. The formula  $p(F|E) = \frac{p(E|F)p(F)}{p(E)}$  is centered. Arrows point from each term to its corresponding label: 'likelihood' points to  $p(E|F)$ , 'Prior probability of F' points to  $p(F)$ , 'Posterior probability' points to  $p(F|E)$ , and 'Prior probability of E' points to  $p(E)$ .

# Derivation of Bayes' Theorem

$$p(F|E) = \frac{p(E|F)p(F)}{p(E)}$$

Note that

$$p(E) = p(E|F)p(F) + p(E|\bar{F})p(\bar{F})$$

since  $p(E) = p(E \cap F) + p(E \cap \bar{F})$

$$p(E) = p(E \cap F) + p(E \cap \bar{F}) = p(E|F)p(F) + p(E|\bar{F})p(\bar{F})$$

Hence,

$$p(F|E) = \frac{p(E|F)p(F)}{p(E|F)p(F) + p(E|\bar{F})p(\bar{F})}$$

# Applying Bayes' Theorem

**Example:** Suppose that one person in 100,000 has a particular disease. There is a test for the disease that gives a positive result 99% of the time when given to someone with the disease. When given to someone without the disease, 99.5% of the time it gives a negative result. Find

- a) the probability that a person who test positive has the disease.
- b) the probability that a person who test negative does not have the disease.
- Should someone who tests positive be worried?

$$P(D) = 10^{-5}$$

$$P(+|D) = 0.99$$

$$P(-|\bar{D}) = 0.995$$

$$P(\bar{D}) = 0.99999$$

$$P(+|\bar{D}) = 0.01$$

$$P(-|D) = 0.005$$

# Applying Bayes' Theorem

**Solution:** Let  $D$  be the event that the person has the disease, and  $E$  be the event that this person tests positive. We need to compute  $p(D|E)$  from  $p(D)$ ,  $p(E|D)$ ,  $p(E|\bar{D})$ ,  $p(\bar{D})$ .

$$\overline{D} \quad \overline{D}$$

$$p(D) = 1/100,000 = 0.00001 \quad p(\overline{D}) = 1 - 0.00001 = 0.99999$$

$$p(E|D) = .99 \quad p(\overline{E}|D) = .01 \quad p(E|\overline{D}) = .005 \quad p(\overline{E}|\overline{D}) = .995$$

$$\begin{aligned} p(D|E) &= \frac{p(E|D)p(D)}{p(E|D)p(D) + p(E|\overline{D})p(\overline{D})} \\ &= \frac{(0.99)(0.00001)}{(0.99)(0.00001) + (0.005)(0.99999)} \end{aligned}$$

Can you use this formula to explain why the resulting probability is surprisingly small?

$$\approx 0.002$$

So, don't worry too much, if your test for this disease comes back positive.

# Applying Bayes' Theorem

- What if the result is negative?

So, the probability you have the disease if you test negative is

$$\begin{aligned} p(D|\bar{E}) &\approx 1 - 0.9999999 \\ &= 0.0000001. \end{aligned}$$

$$\begin{aligned} p(\bar{D}|\bar{E}) &= \frac{p(\bar{E}|\bar{D})p(\bar{D})}{p(\bar{E}|\bar{D})p(\bar{D}) + p(\bar{E}|D)p(D)} \\ &= \frac{(0.995)(0.99999)}{(0.995)(0.99999) + (0.01)(0.00001)} \\ &\approx 0.9999999 \end{aligned}$$

- So, it is extremely unlikely you have the disease if you test negative.

# Generalized Bayes' Theorem

**Generalized Bayes' Theorem:** Suppose that  $E$  is an event from a sample space  $S$  and that  $F_1, F_2, \dots, F_n$  are mutually exclusive events such that

$$\bigcup_i^n F_i = S.$$

Assume that  $p(E) \neq 0$  for  $i = 1, 2, \dots, n$ . Then

$$p(F_j|E) = \frac{p(E|F_j)p(F_j)}{\sum_{i=1}^n p(E|F_i)p(F_i)}.$$

# Bayesian Spam Filters

- How do we develop a tool for determining whether an email is likely to be spam?
- If we have an initial set  $B$  of spam messages and set  $G$  of non-spam messages. We can use this information along with Bayes' law to predict the probability that a new email message is spam.
- We look at a particular word  $w$ , and count the number of times that it occurs in  $B$  and in  $G$ ;  $n_B(w)$  and  $n_G(w)$ .
  - Estimated probability that an email containing  $w$  is spam:  
 $p(w) = n_B(w)/|B|$
  - Estimated probability that an email containing  $w$  is not spam:  
 $q(w) = n_G(w)/|G|$

# Bayesian Spam Filters

- Let  $S$  be the event that the message is spam, and  $E$  be the event that the message contains the word  $w$ .
- Using Bayes' rule

$$p(S|E) = \frac{p(E|S)p(S)}{p(E|S)p(S) + p(E|\bar{S})p(\bar{S})}$$

$$p(S|E) = \frac{p(E|S)}{p(E|S) + p(E|\bar{S})}$$

Assuming that it is equally likely that an arbitrary message is spam and is not spam; i.e.,  $p(S) = \frac{1}{2}$ .

# Bayesian Spam Filters

Using our empirical estimates of  $p(E | S)$  and  $p(E | \bar{S})$ .

$$r(w) = \frac{p(w)}{p(w) + q(w)}$$

$r(w)$  estimates the probability that the message is spam. We can class the message as spam if  $r(w)$  is above a **threshold**.

# Bayesian Spam Filters

**Example:** We find that the word “Rolex” occurs in 250 out of 2000 spam messages and occurs in 5 out of 1000 non-spam messages. Estimate the probability that an incoming message is spam. Suppose our threshold for rejecting the email is 0.9.

**Solution:**  $p(\text{Rolex}) = 250/2000 = .125$  and  $q(\text{Rolex}) = 5/1000 = 0.005$ .

$$r(\text{Rolex}) = \frac{p(\text{Rolex})}{p(\text{Rolex}) + q(\text{Rolex})} = \frac{0.125}{0.125 + .005} = \frac{0.125}{0.125 + .005} \approx 0.962$$

We class the message as spam  
and reject the email!

# Bayesian Spam Filters using Multiple Words

- Accuracy can be improved by considering more than one word as evidence.
- Consider the case where  $E_1$  and  $E_2$  denote the events that the message contains the words  $w_1$  and  $w_2$  respectively.
- We make the simplifying assumption that the events are independent. And again we assume that  $p(S) = \frac{1}{2}$ .

$$p(S|E_1 \cap E_2) = \frac{p(E_1|S)p(E_2|S)}{p(E_1|S)p(E_2|S) + p(E_1|\bar{S})p(E_2|\bar{S})}$$

$$r(w_1, w_2) = \frac{p(w_1)p(w_2)}{p(w_1)p(w_2) + q(w_1)q(w_2)}$$

# Bayesian Spam Filters using Multiple Words

**Example:** We have 2000 spam messages and 1000 non-spam messages.

- The word “stock” occurs 400 times in the spam messages and 60 times in the non-spam.
- The word “undervalued” occurs in 200 spam messages and 25 non-spam.

# Bayesian Spam Filters using Multiple Words

**Solution:**

$$p(\text{stock}) = 400/2000 = .2$$

$$q(\text{stock}) = 60/1000 = .06$$

$$p(\text{undervalued}) = 200/2000 = .1$$

$$q(\text{undervalued}) = 25/1000 = .025$$

$$\begin{aligned} r(\text{stock, undervalued}) &= \frac{p(\text{stock}) p(\text{undervalued})}{p(\text{stock}) p(\text{undervalued}) + q(\text{stock}) q(\text{undervalued})} \\ &= \frac{(0.2)(0.1)}{(0.2)(0.1) + (0.06)(0.025)} \approx 0.930 \end{aligned}$$

If our threshold is .9, we class the message as spam and reject it.

# Bayesian Spam Filters using Multiple Words

- In general, the more words we consider, the more accurate the spam filter. With the independence assumption if we consider  $k$  words:

$$p(S | \bigcap_{i=1}^k E_i) = \frac{\prod_{i=1}^k p(E_i | S)}{\prod_{i=1}^k p(E_i | S) + \prod_{i=1}^k p(E_i | \bar{S})}$$

$$r(w_1, w_2, \dots, w_n) = \frac{\prod_{i=1}^k p(w_i)}{\prod_{i=1}^k p(w_i) + \prod_{i=1}^k q(w_i)}$$

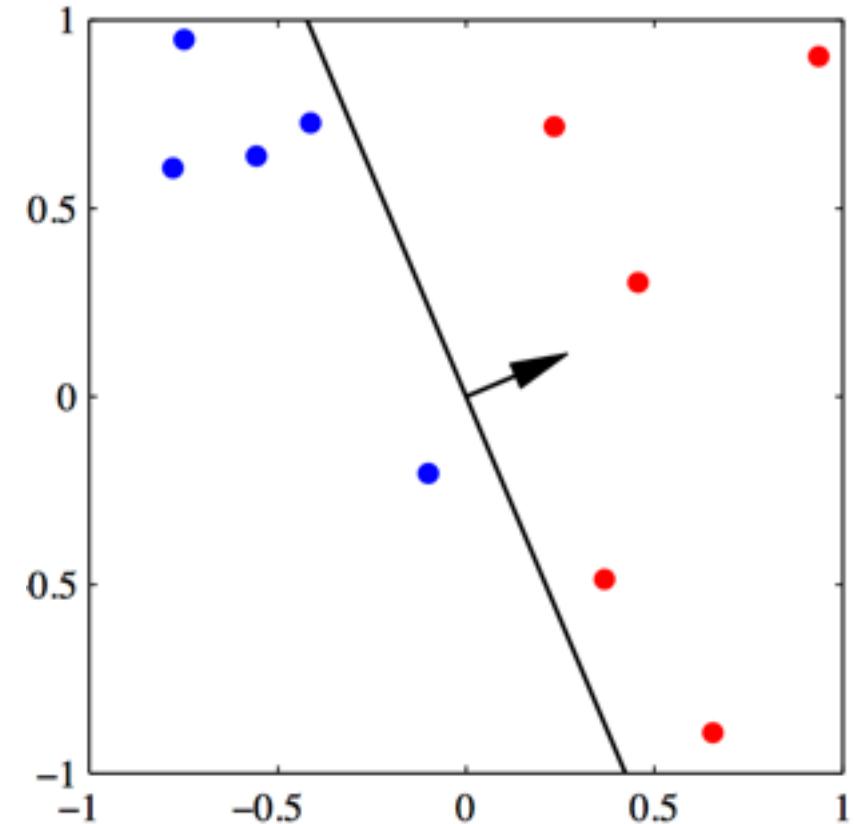
We can further improve the filter by considering pairs of words as a single block or certain types of strings.

# Discriminant analysis

Section 3

# Linear classification

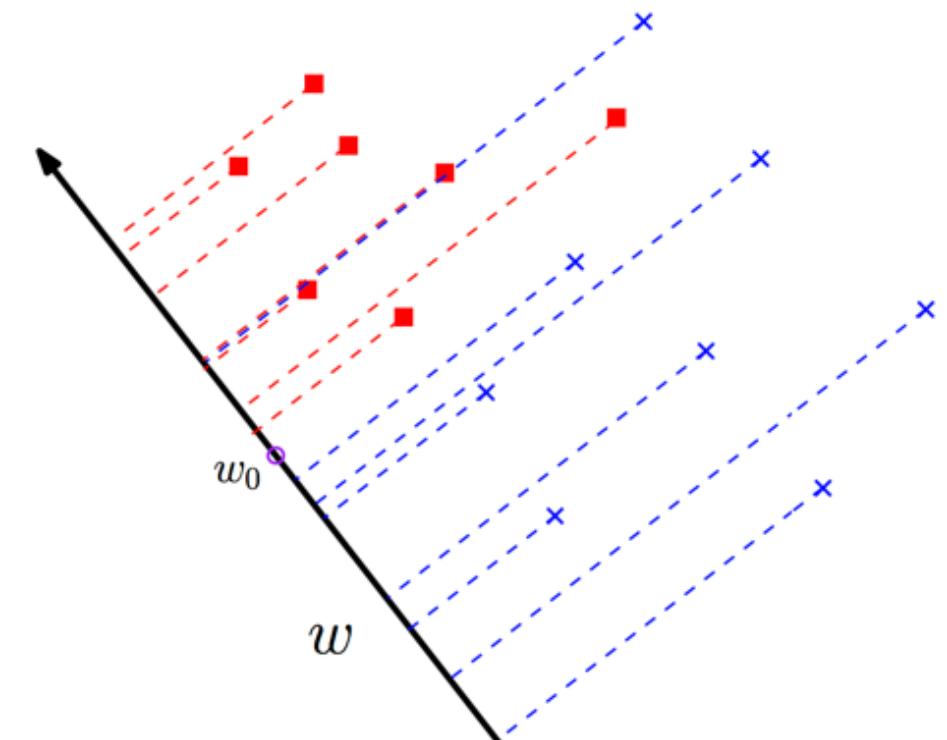
- Focus on linear classification model, i.e., the decision boundary is a linear function of  $x$ 
  - Defined by  $(D - 1)$ -dimensional hyperplane
  - If the data can be separated exactly by linear decision surfaces, they are called **linearly separable**
  - Implicit assumption: Classes can be modeled well by Gaussians
  - Simply speaking, treat **classification as a projection** problem



From PRML (Bishop, 2006)

# Projection

- Assume we know the basic vector  $w$ , we can compute the projection,  $y$ , of any points,  $x$ .
- Threshold  $w_0$ , such that we decide on  $C_1$  if  $y \geq w_0$  and  $C_2$  otherwise.



ដីបាប្រាប់លុក 2 ខ្លួចឱនទំនើនកំណែ

$$7, 14, 21, 28, \dots, 70$$

$$7, 7n \approx 1 \text{ ឬ } 7n \approx 17 \text{ ឬ } 7n \approx 24 \text{ ឬ } 7n \approx 27$$

$$\text{ខ្លួចឱនទំនើនកំណែ } 7: 7, 14, 21, 28, \dots \rightarrow 7n = 3600$$

$$\text{ខ្លួចឱនទំនើនកំណែ } 2: 7, 17, 27, 37, \dots$$

$$\left. \begin{array}{l} 6 \text{ ឬ } 7 \approx 7 \\ 60 \text{ ឬ } 7 \approx 70 \end{array} \right\} 89$$

ក្រោមរំលែក 1 ស៊ូរិយៈ 3600 គីឡូតី

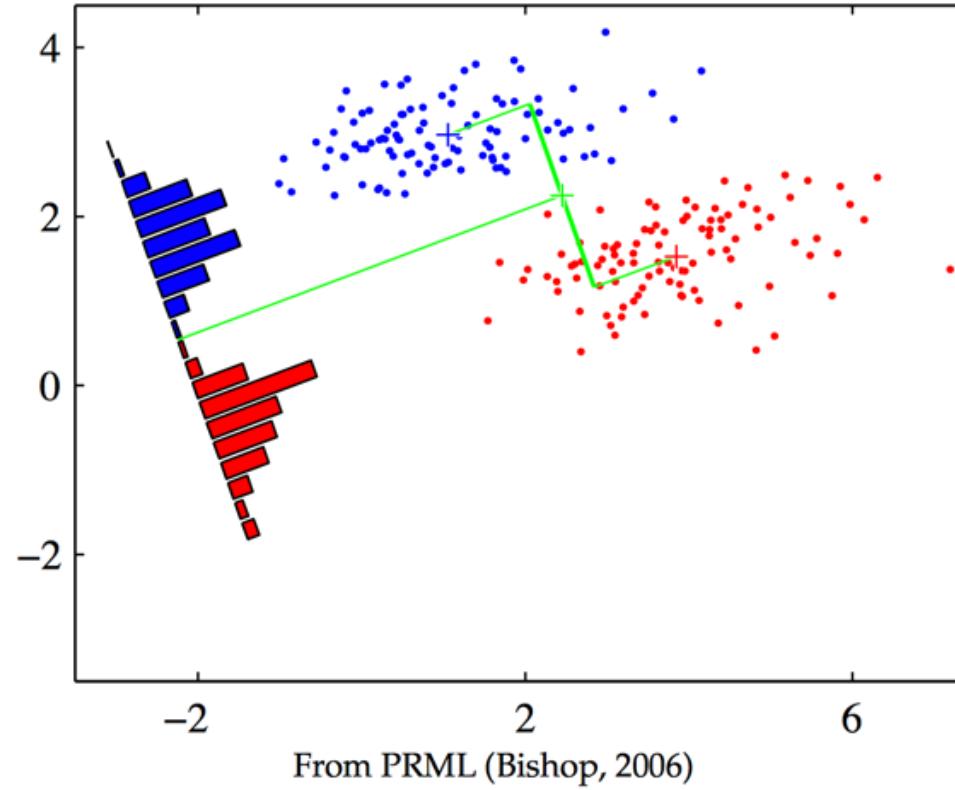
$$\begin{array}{ccccccc} 3597 & & 3587 & & 3598 & & 3577 \\ & & & & \underbrace{10\textcircled{1} + 7}_{3600} & & \\ & & & & \overbrace{\quad\quad\quad}^{\text{គីឡូតី}} & & \end{array}$$

$$\left. \begin{array}{l} \text{ខ្លួចឱនទំនើនកំណែ } 514 \text{ រឿង} \\ \text{ខ្លួចឱនទំនើនកំណែ } 359 \text{ រឿង} \end{array} \right\} 873 - 51 = 822$$

$$70n + 7 \approx 3600$$

$$70n \approx 3593$$

## LDA



- Separate samples of distinct groups by projecting them onto a space that
  - Maximize their between-class separability while
  - Minimize their within-class variability

# Linear discriminant analysis

- We model the distribution of the predictors  $X$  separately in each of the response classes (i.e. given  $Y$ ), and then use Bayes' theorem to flip these around into estimates for  $\Pr(Y = k | X = x)$
- Assume that we have 1 predictor
- Suppose we assume that  $f_k(x)$  is normal or Gaussian

$$f_k(x) = \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{1}{2\sigma_k^2}(x - \mu_k)^2\right)$$

# Bayes' theorem for classification

- Let  $\pi_k$  be the overall or prior probability of the  $k^{\text{th}}$  class;
- This is the probability that a given observation is associated with the  $k^{\text{th}}$  category of the response variable  $Y$ .
- Let  $f_k(x) \equiv \Pr(X = x|Y = k)$  denote the density function of  $X$  for an observation that comes from the  $k^{\text{th}}$  class

$$\Pr(Y = k|X = x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)}.$$

# Bays' theorem for classification

$$p_k(x) = \frac{\pi_k \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x - \mu_k)^2\right)}{\sum_{l=1}^K \pi_l \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x - \mu_l)^2\right)}.$$

- Taking a log, we get

$$\delta_k(x) = x \cdot \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k)$$

- We will assign  $x$  to class  $k$  if  $\delta_k(x)$  is largest.

# Deriving Bayes' classifier solution

$$p_k(x) = \frac{\pi_k \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x - \mu_k)^2\right)}{p(X = x)}$$

$$\log(p_k(x)) = \log(\pi_k) + \log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) + \log\left(\exp\left(-\frac{1}{2\sigma^2}(x - \mu_k)^2\right)\right) - \log(p(X = x))$$

$$\log(p_k(x)) = \log(\pi_k) + \log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) + \left(-\frac{1}{2\sigma^2}(x - \mu_k)^2\right) - \log(p(X = x))$$

$$\log(p_k(x)) = \log(\pi_k) + \log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) - \frac{x^2}{2\sigma^2} + \frac{x\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} - \log(p(X = x))$$

$$\delta_k(x) = \log(\pi_k) + \frac{x\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2}$$

# Linear discriminant analysis (LDA)

- The linear discriminant analysis method approximates the Bayes classifier by plugging estimates for  $\pi_k$ ,  $\mu_k$ , and  $\sigma_2$ :

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} x_i$$

$$\hat{\sigma}^2 = \frac{1}{n - K} \sum_{k=1}^K \sum_{i:y_i=k} (x_i - \hat{\mu}_k)^2$$

$$\hat{\pi}_k = n_k/n.$$

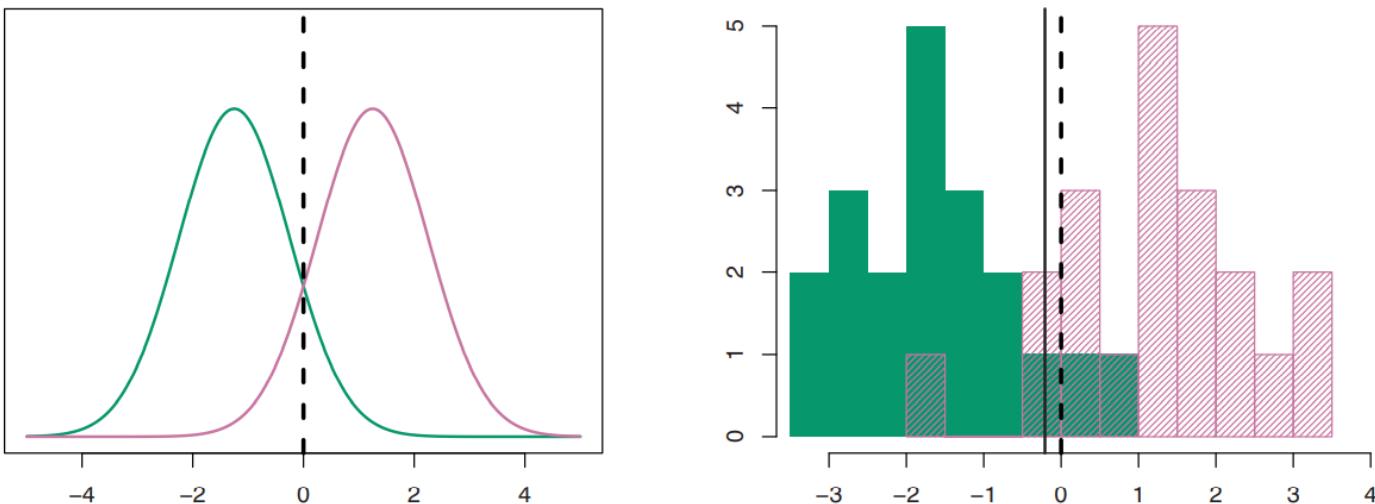
# LDA discriminant function

- Assign  $x$  to class  $k$  if

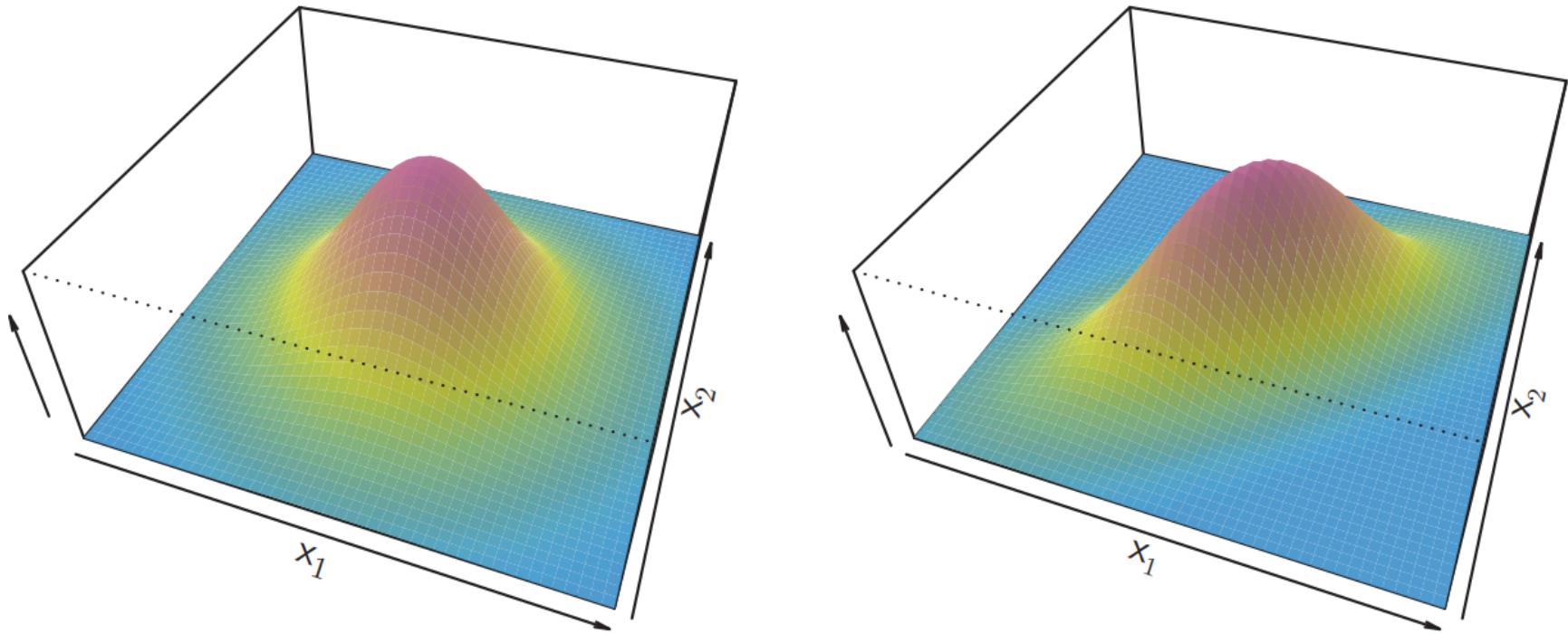
$$\hat{\delta}_k(x) = x \cdot \frac{\hat{\mu}_k}{\hat{\sigma}^2} - \frac{\hat{\mu}_k^2}{2\hat{\sigma}^2} + \log(\hat{\pi}_k)$$

is the largest

# Bayes decision boundary and LDA

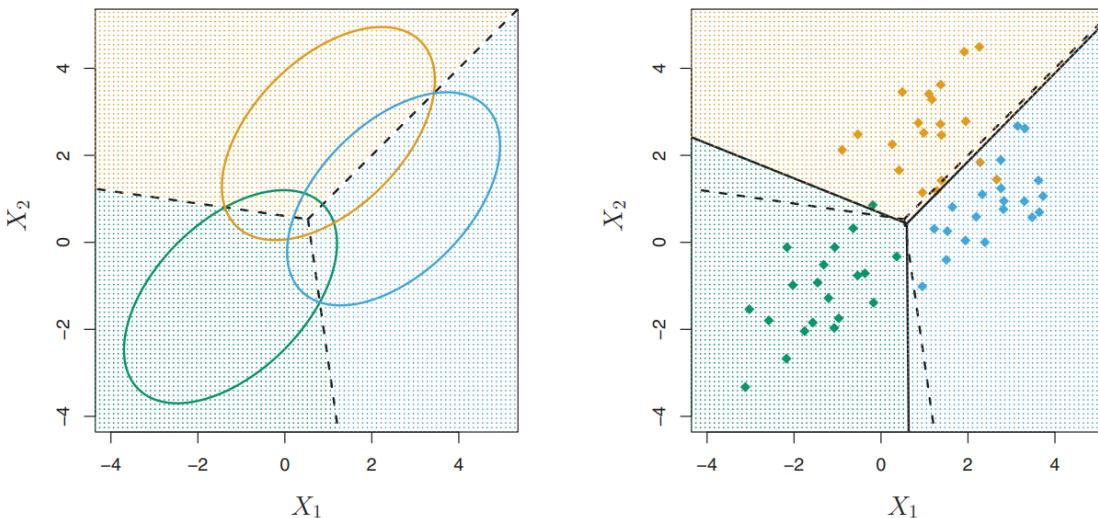


**FIGURE 4.4.** Left: Two one-dimensional normal density functions are shown. The dashed vertical line represents the Bayes decision boundary. Right: 20 observations were drawn from each of the two classes, and are shown as histograms. The Bayes decision boundary is again shown as a dashed vertical line. The solid vertical line represents the LDA decision boundary estimated from the training data.



**FIGURE 4.5.** Two multivariate Gaussian density functions are shown, with  $p = 2$ . Left: The two predictors are uncorrelated. Right: The two variables have a correlation of 0.7.

# Discriminant analysis with more than one predictors



**FIGURE 4.6.** An example with three classes. The observations from each class are drawn from a multivariate Gaussian distribution with  $p = 2$ , with a class-specific mean vector and a common covariance matrix. Left: Ellipses that contain 95 % of the probability for each of the three classes are shown. The dashed lines are the Bayes decision boundaries. Right: 20 observations were generated from each class, and the corresponding LDA decision boundaries are indicated using solid black lines. The Bayes decision boundaries are once again shown as dashed lines.

# Quadratic Discriminant Analysis

A generic class of discriminant analysis

From the distribution function, if we take log of the function

$$\begin{aligned}\log P(y = k|x) &= \log P(x|y = k) + \log P(y = k) + Cst \\ &= -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2}(x - \mu_k)^t \Sigma_k^{-1} (x - \mu_k) + \log P(y = k) + Cst,\end{aligned}$$

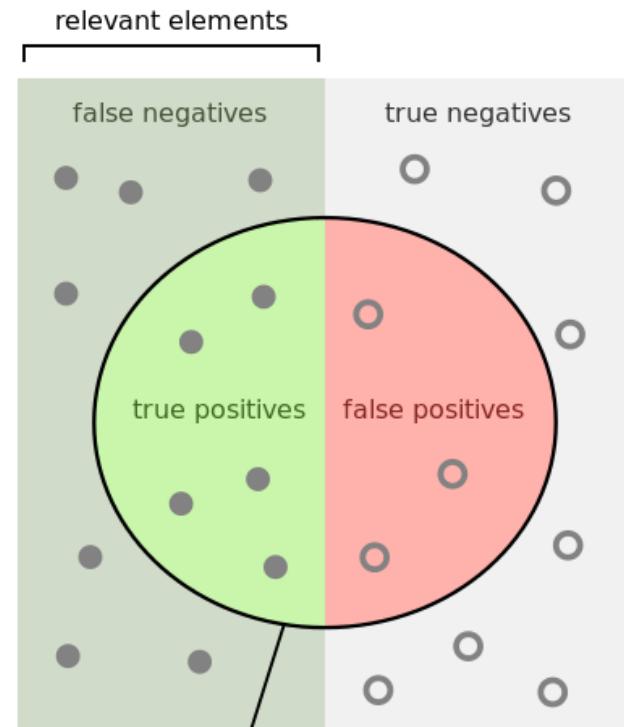
where the constant term  $Cst$  corresponds to the denominator  $P(x)$ , in addition to other constant terms from the Gaussian. The predicted class is the one that maximises this log-posterior.

LDA is a special case of QDA when each class assume shared covariance.

# Accuracy, Precision and Recall

	Actual Positive (p)	Actual Negative (n)
The model says “Yes” = positive (y)	True positives	False positives
The model says “No” = not positive (n)	False negatives	True negatives

- Accuracy =  $(TP + TN) / (TP + FP + TN + FN)$
- Recall (Completeness) = true positive rate =  $TP / (TP + FN)$
- Precision (Exactness) = the accuracy over the cases predicted to be positive,  $TP / (TP + FP)$
- F-measure = the harmonic mean of precision and recall
  - = the balance between recall and precision
  - =  $2 \cdot \frac{precision * recall}{precision + recall}$



How many selected items are relevant?

Precision =  $\frac{\text{true positives}}{\text{true positives} + \text{false positives}}$

How many relevant items are selected?

Recall =  $\frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$

# Receiver operating characteristics Area under the ROC curve

**True Positive Rate (TPR)** is a synonym for recall and is therefore defined as follows:

$$TPR = \frac{TP}{TP + FN}$$

**False Positive Rate (FPR)** is defined as follows:

$$FPR = \frac{FP}{FP + TN}$$

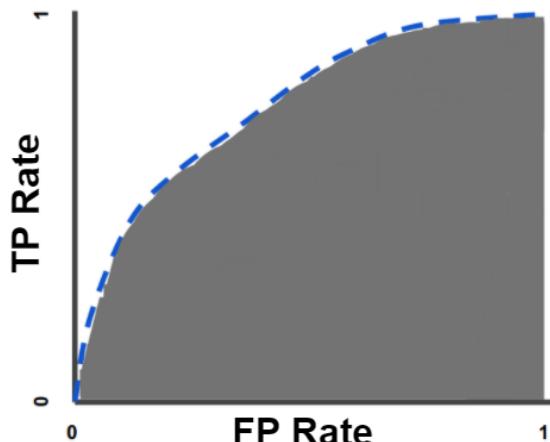


Figure 5. AUC (Area under the ROC Curve).

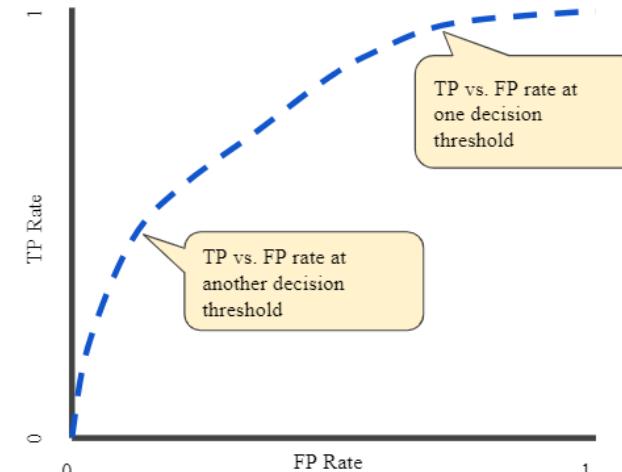
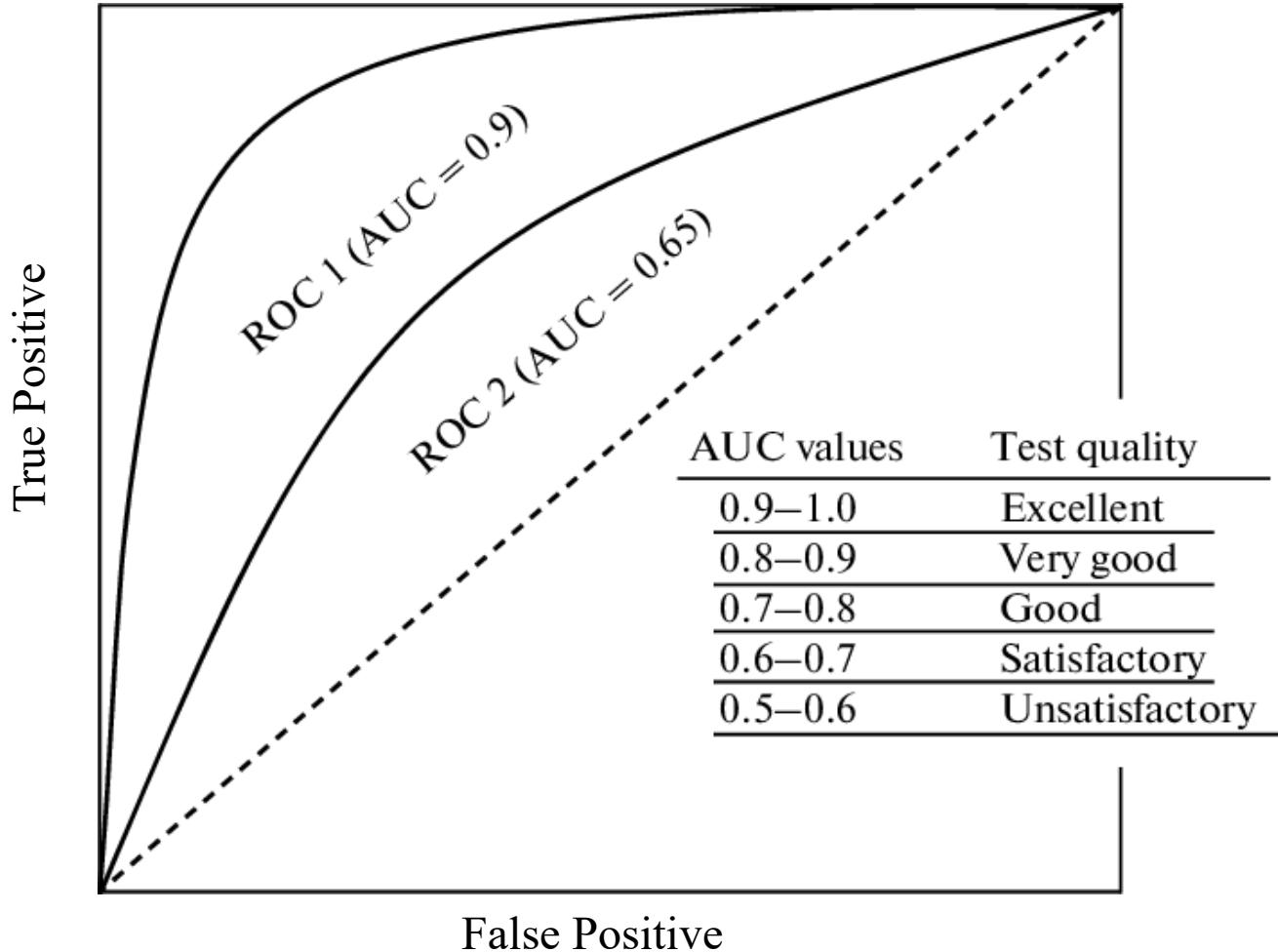


Figure 4. TP vs. FP rate at different classification thresholds.

# AUC - ROC



# Coding Practice

- [Implementing linear discriminant analysis \(LDA\) in Python - IBM Developer](#)
- [Linear Discriminant Analysis Made Simple & How To Tutorial \(spotintelligence.com\)](#)

# End of Lecture 1

## Question?