

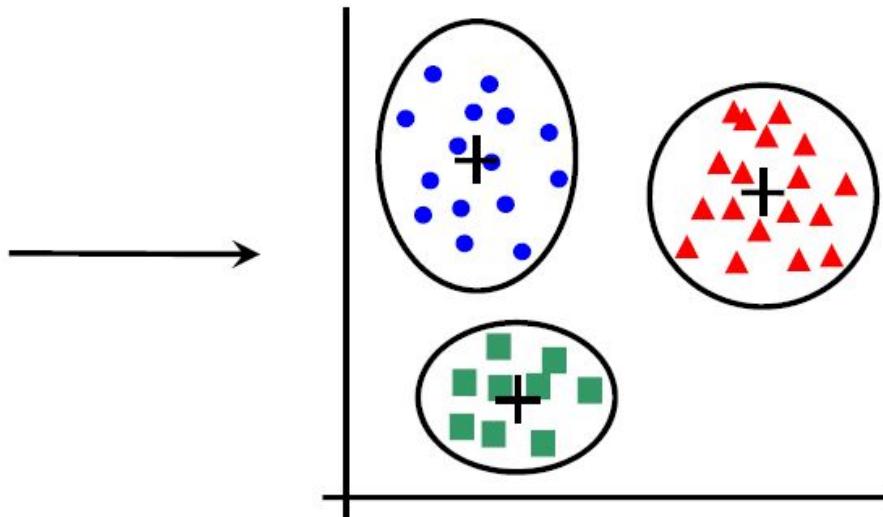
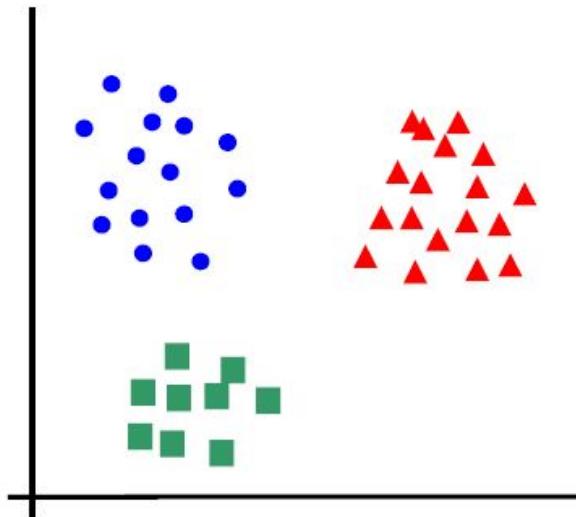
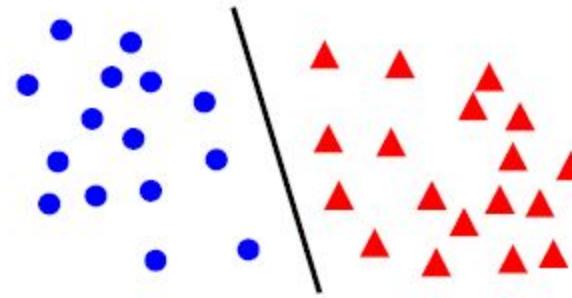
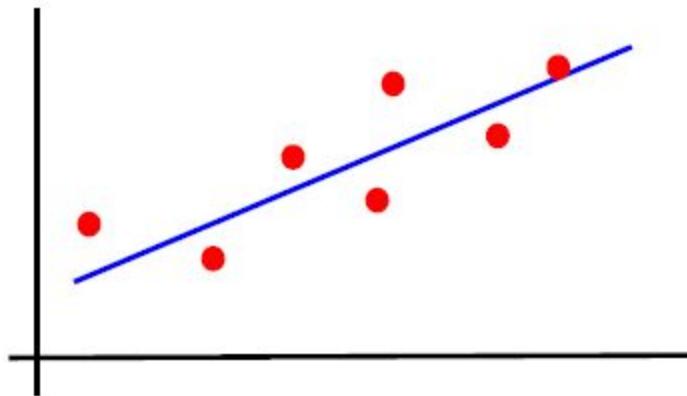
# **Machine Learning**

## **Introduction**

**Indian Institute of Information Technology  
Sri City, Chittoor**



# Welcome to Machine Learning Class



# Today's Agenda

- Course plan
  - Pre-requisite
  - Topics
  - Textbooks/References
  - Evaluation components
  - Honor code
- Introduction to machine learning
  - What is ML?
  - When do we use ML?
  - Applications
  - Relation with AI and DL
  - Relation with other fields
  - Different machine learning paradigms

# Pre-requisite

- Probability
  - Distribution, random variable, expectation, conditional probability, variance, density
- Linear algebra
  - Linear transformation, Rank, Positive definite matrix
  - Eigenvalue & Eigenvector, Diagonalization
- Calculus (Optimization – Minima, maxima)
- Basic programming
  - Python (First Priority)
  - Matlab/C/C++ (Second Priority)

# Topics

- Supervised Learning
  - Classification
  - Regression
- Unsupervised Learning
  - Clustering
- Semi-supervised Learning
- Reinforcement Learning

# Textbooks/References

1. “Pattern Classification” by R. O. Duda, P. E. Hart and D. G. Stork.
2. “An Introduction to Statistical Learning” by Gareth James, Daniela Witten, Trevor Hastie and Robert Tibshirani.
3. “Pattern Recognition and Machine Learning” by Christopher M. Bishop.
4. “Introduction to Machine Learning” by Ethem Alpaydin.
5. “Pattern Recognition: An Algorithmic Approach” by M. Narasimha Murty, V. Susheela Devi.
6. “Machine learning” by Tom Mitchell.

# Evaluation Components

- Mid-Exam(s): 20%
- End-Exam: 30%
- Assignments: 20%
- Scheduled Quiz: 10%
- Project: 20%

This is tentative and will be finalized in Class committee meeting.

# Honor Code

## Do's

- Write down the code independently
- Submit the assignment within the deadline
- Read the books/references for detail description of the topics

## Don'ts

- copy, refer to, or look at any **official or unofficial** previous years' solutions in **preparing** the answers

# Introduction to ML

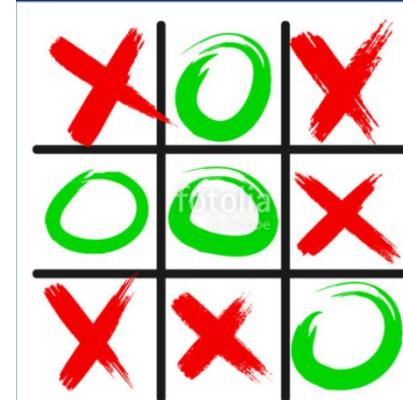
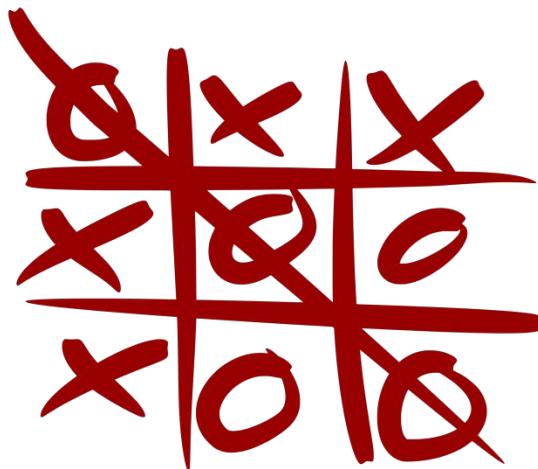
- What is ML?
- Terminologies used in ML
- When do we use ML?
- Applications
- Relation with AI and DL
- Relation with other fields
- Different machine learning paradigms

# What is ML?

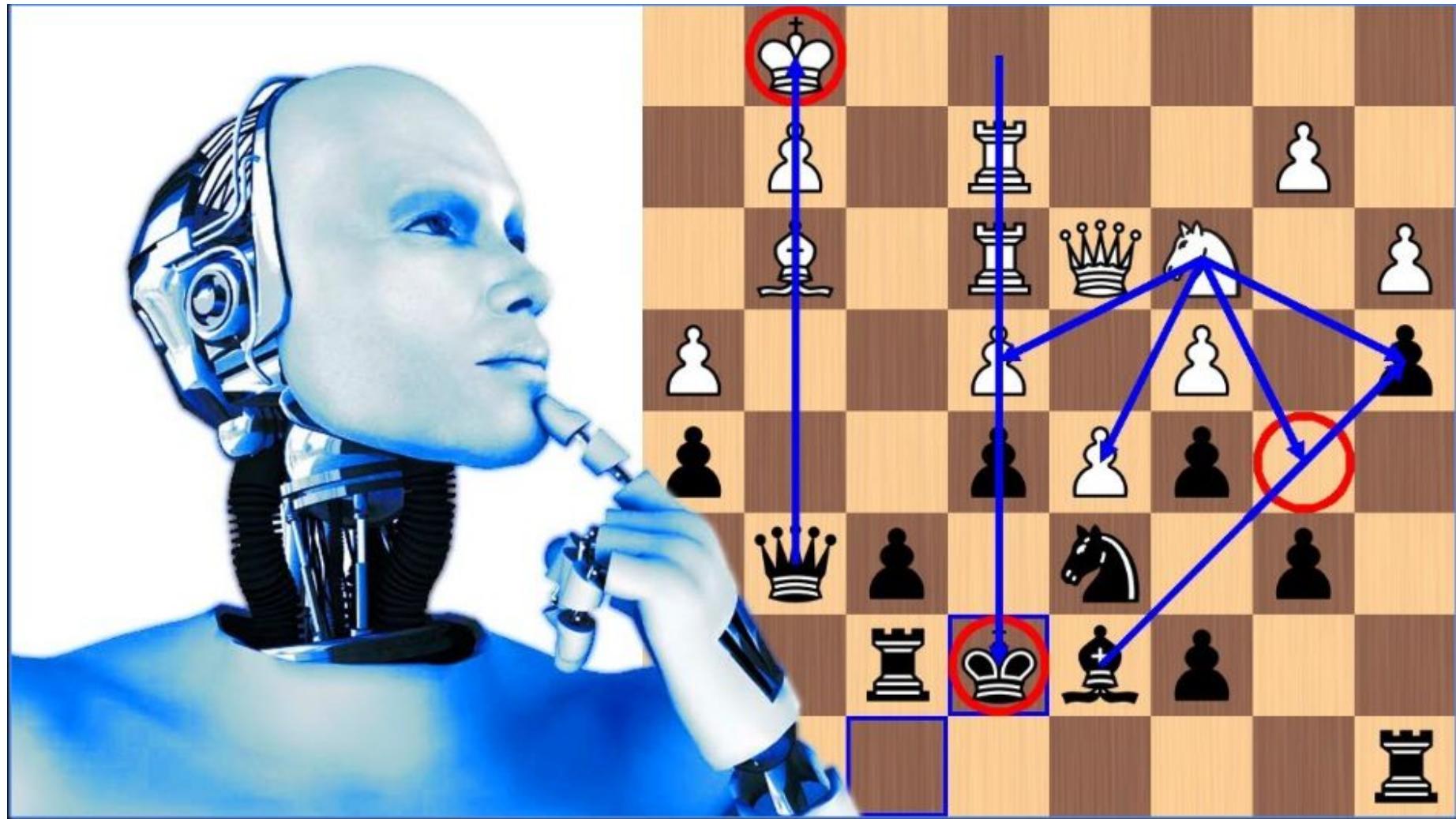
- Machine learning (ML) is the study of computer algorithms that improve automatically through experience.
- Machine-learning algorithms use statistics to find patterns in massive amounts of data.
- Traditionally, software engineering combined human created rules with data to create answers to a problem. Instead, machine learning uses data and answers to discover the rules behind a problem – **F. Chollet, Deep Learning with Python**

# Early AI – Rule based

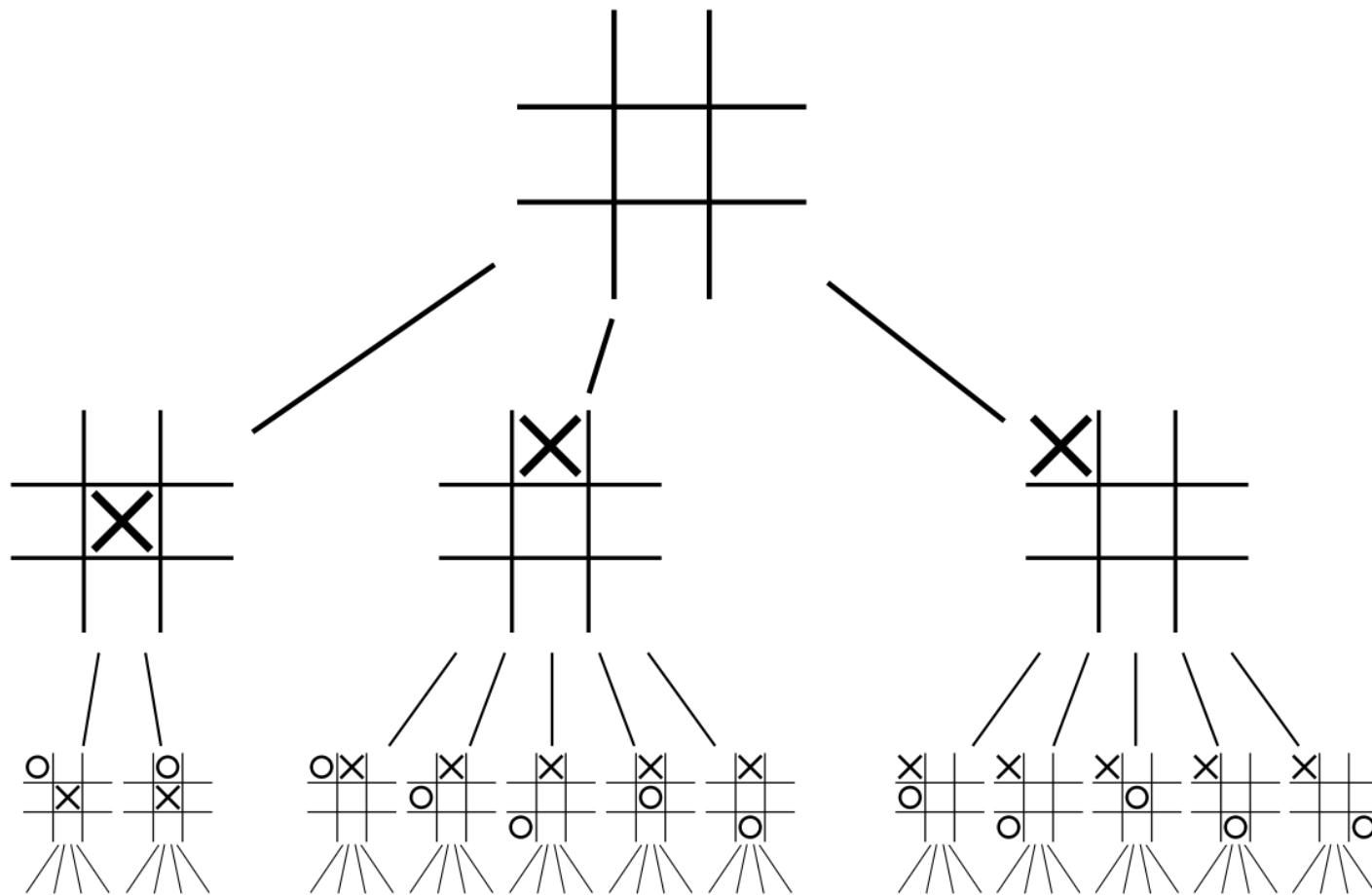
- Tic-tac-toe



# Chess --



# Tic-tac-toe – how it is done?

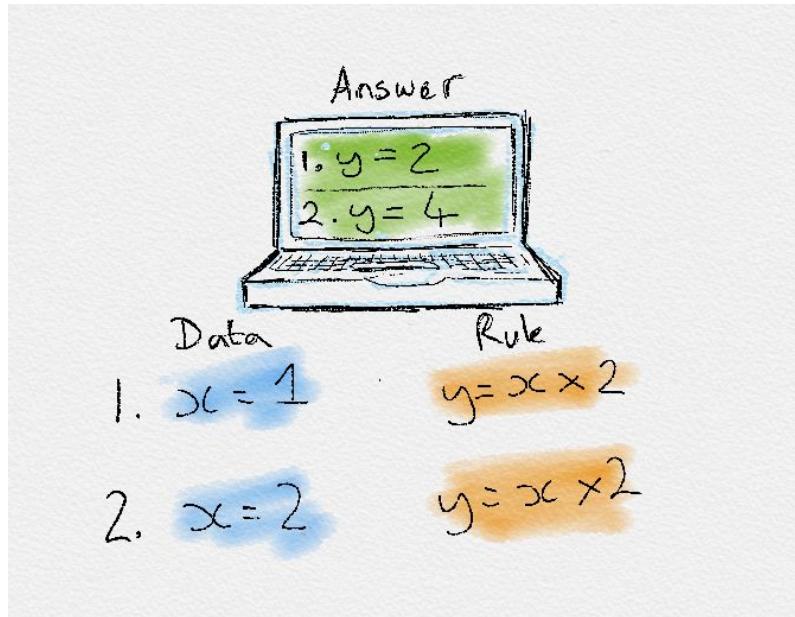


# Deductive Vs Inductive

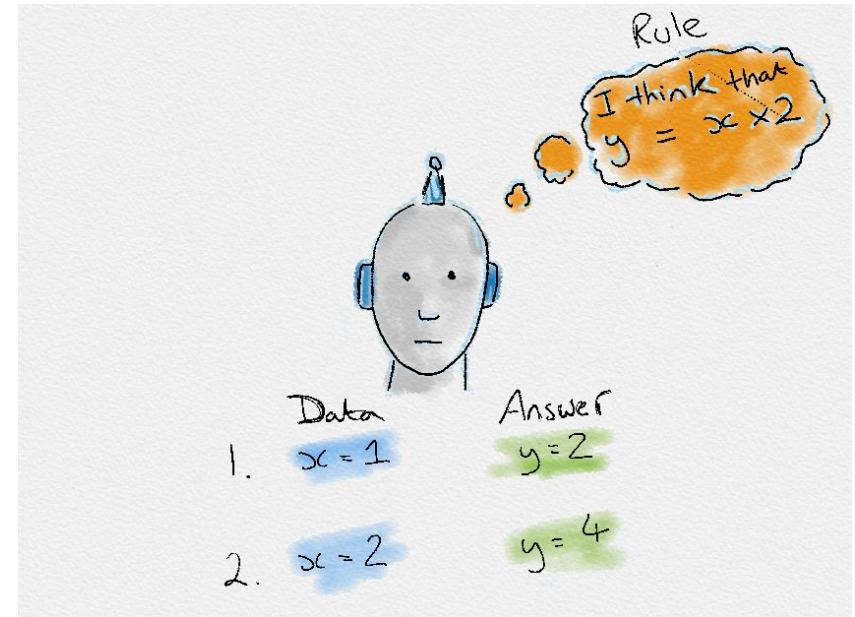
- Giving rules can work for game playing like problems.
- But, you cannot give all rules for a self-driving car.
  - This attempt failed.
- So, inductive means, we will not give explicit rules, but give examples saying which is good which is bad, so on.
  - Machine will learn rules from the data.

# Deductive Vs Inductive

Deductive Reasoning



Inductive Learning



Traditional Programming

Machine Learning

# Machine Learning

- Learning from data



ALGORITHMS



DATA

# Terminologies used in ML

- ML systems learn how to make inference from the input data samples to produce useful predictions on un-seen (test) data.
- Input data:
  - labelled examples: A labelled example includes feature(s) and the label. {features, label}: (x, y)  
For e.g.:

Features:	Label
Normal RBC, Normal HgB	Healthy
Low RBC, Low HgB	Anaemic
  - unlabelled examples: An unlabelled example contains features but not the label. {features, ?}: (x, ?)
    - For e.g.: (Normal RBC, Low HgB)

Features: House type	Price
4BHK,	40,000
4BHK,	15,000
2BHK,	25,000
2BHK,	8,000

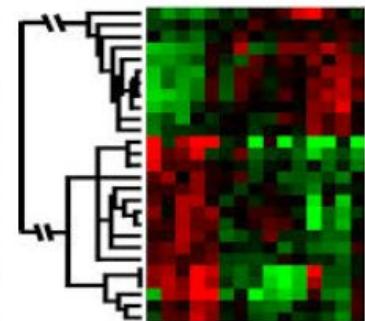
# Terminologies used in ML

- Machine Learning Model:
  - A ML model defines the relationship between the features and label.
    - For e.g.: An anaemia diagnostic model might associate certain features strongly with “anaemic” or “healthy”, and predict the labels based on the association rules it inferred.
  - Two Phases of ML model development
    - **Training** means creating or **learning** the model.
    - **Testing/Inference** means applying the trained model to unlabelled examples in order to infer the label.

# When do we use ML?

ML is used when:

- Human expertise does not exist (navigating on Mars)
- Humans can't explain their expertise (speech recognition)
- Models must be customized (personalized medicine)
- Models are based on huge amounts of data (genomics)



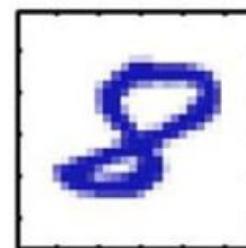
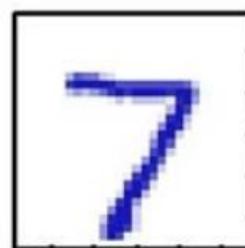
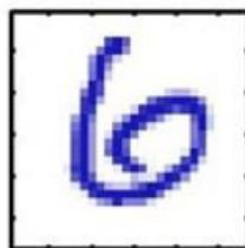
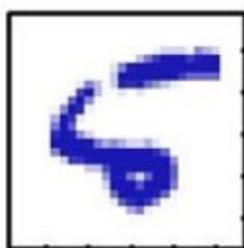
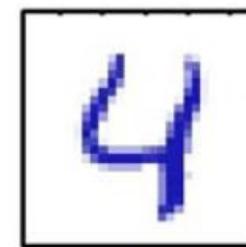
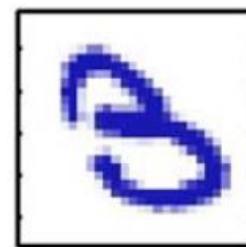
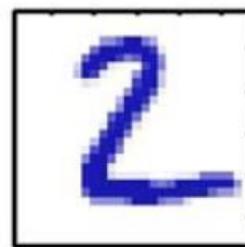
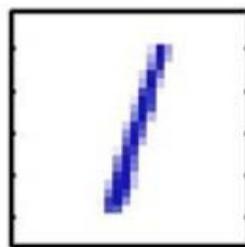
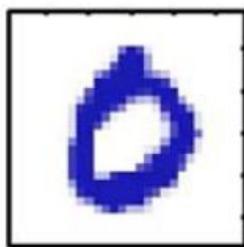
Learning isn't always useful:

- There is no need to “learn” to calculate payroll

# Applications

- Hand-written digit recognition
- Speech recognition
- Face detection
- Object classification
- Email spam detection
- Computational biology
- Autonomous cars
- Computer-aided diagnosis

# Hand-written Digit Recognition



Images are 28 x 28 pixels

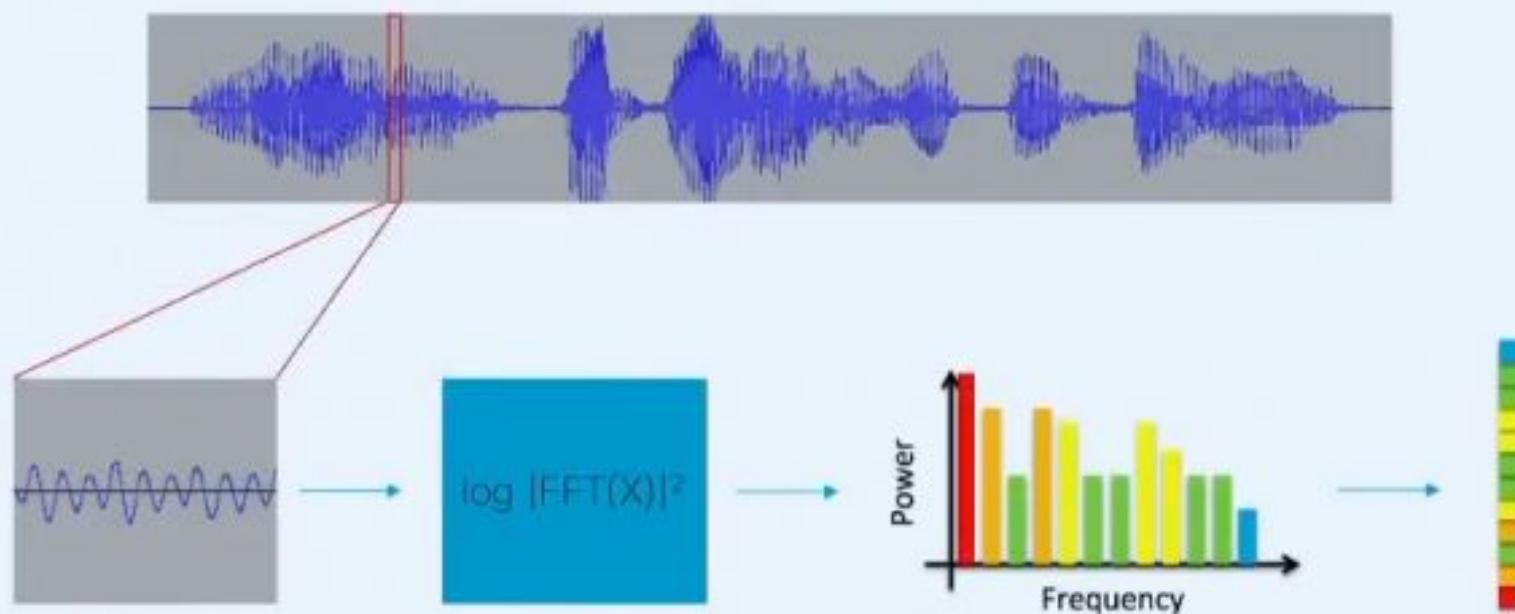
Represent input image as a vector  $\mathbf{x} \in \mathbb{R}^{784}$

Learn a classifier  $f(\mathbf{x})$  such that,

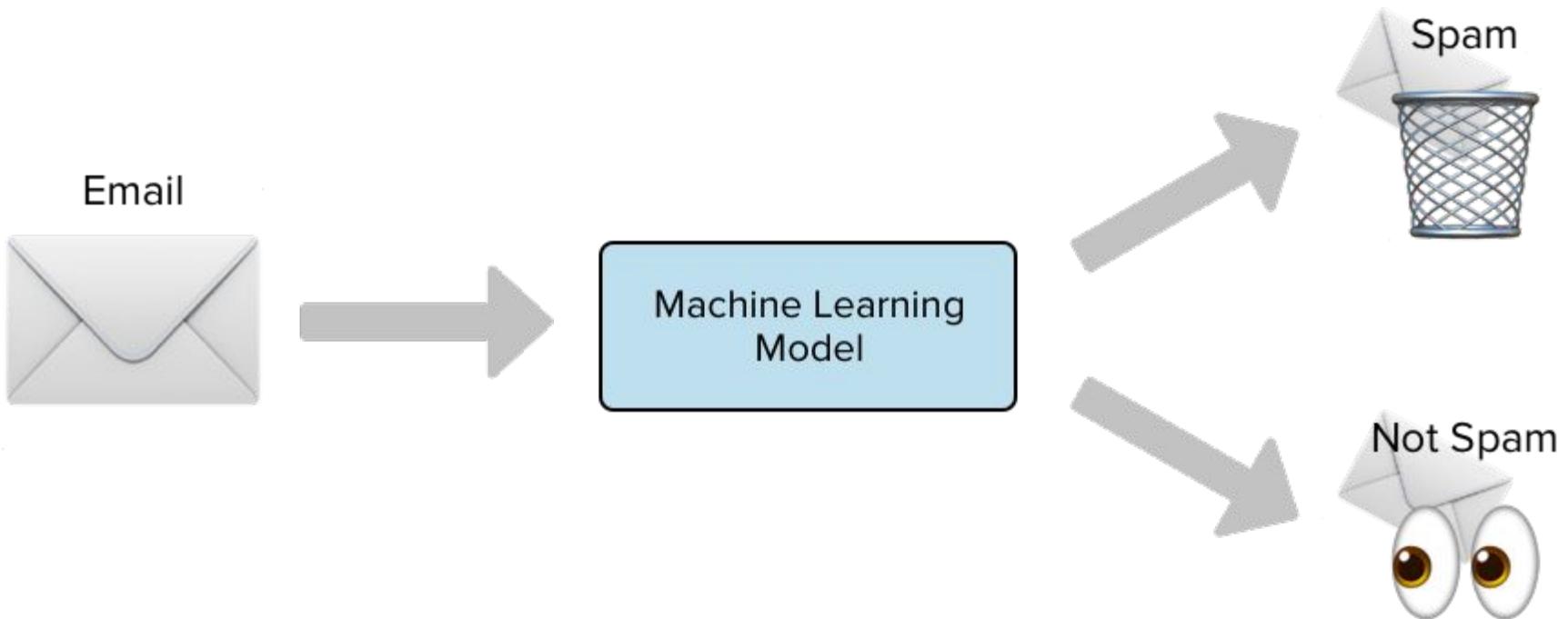
$$f : \mathbf{x} \rightarrow \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

# Speech Recognition

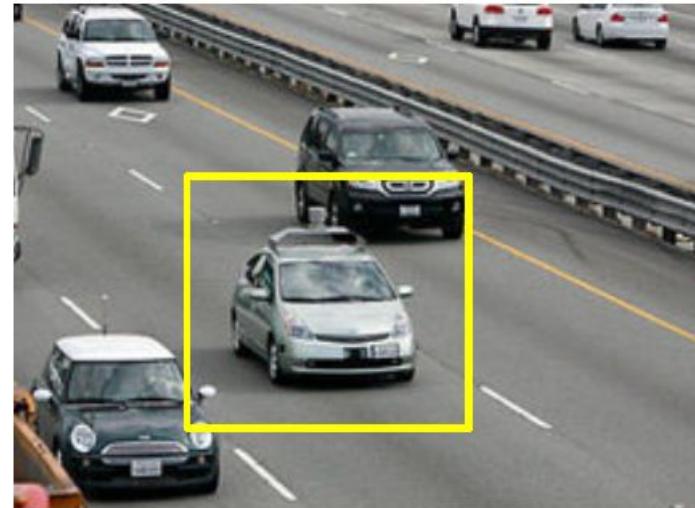
Spectrogram



# Email Spam Detection



# Autonomous Cars

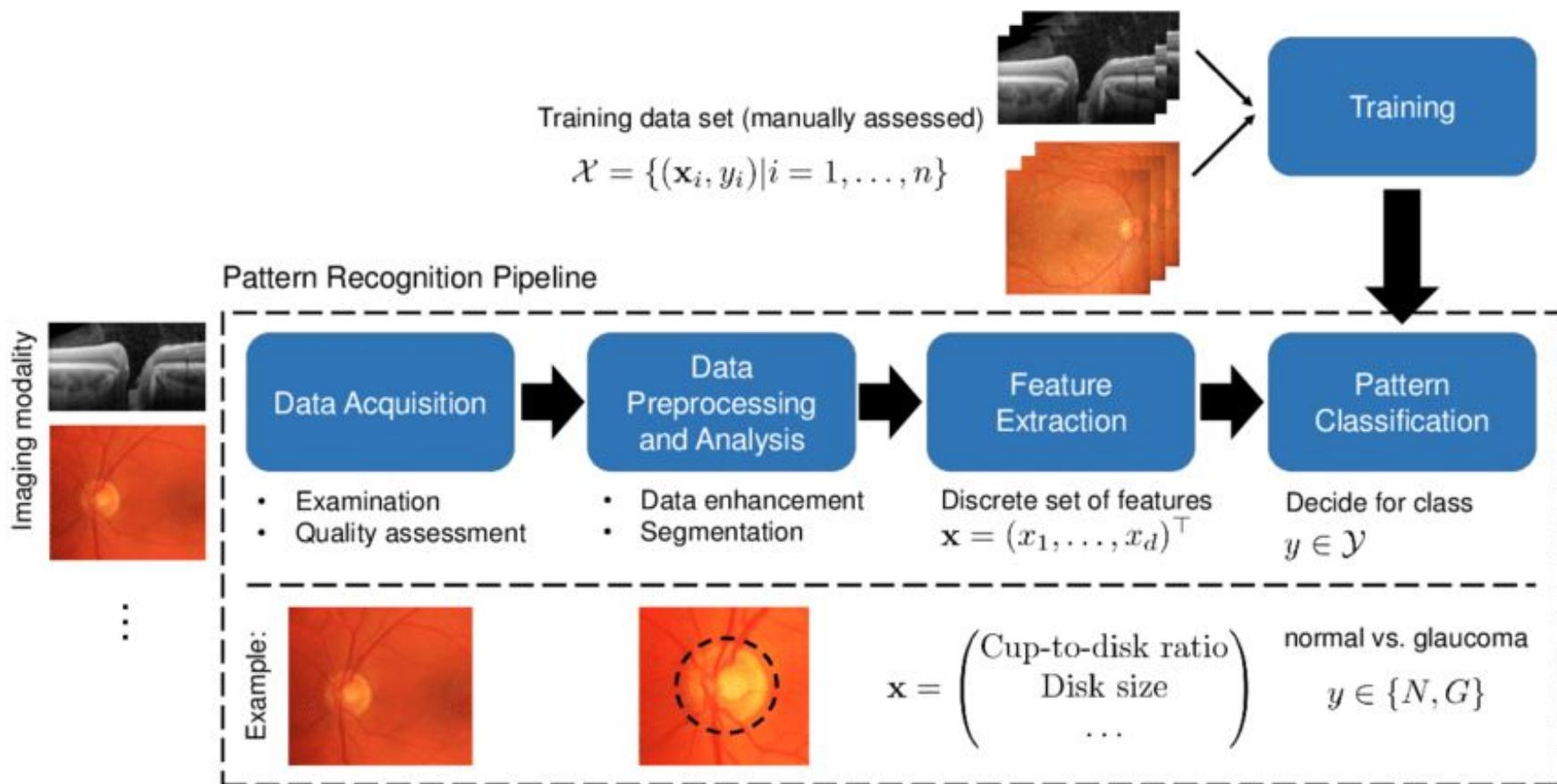


- Nevada made it legal for autonomous cars to drive on roads in June 2011
- As of 2013, four states (Nevada, Florida, California, and Michigan) have legalized autonomous cars

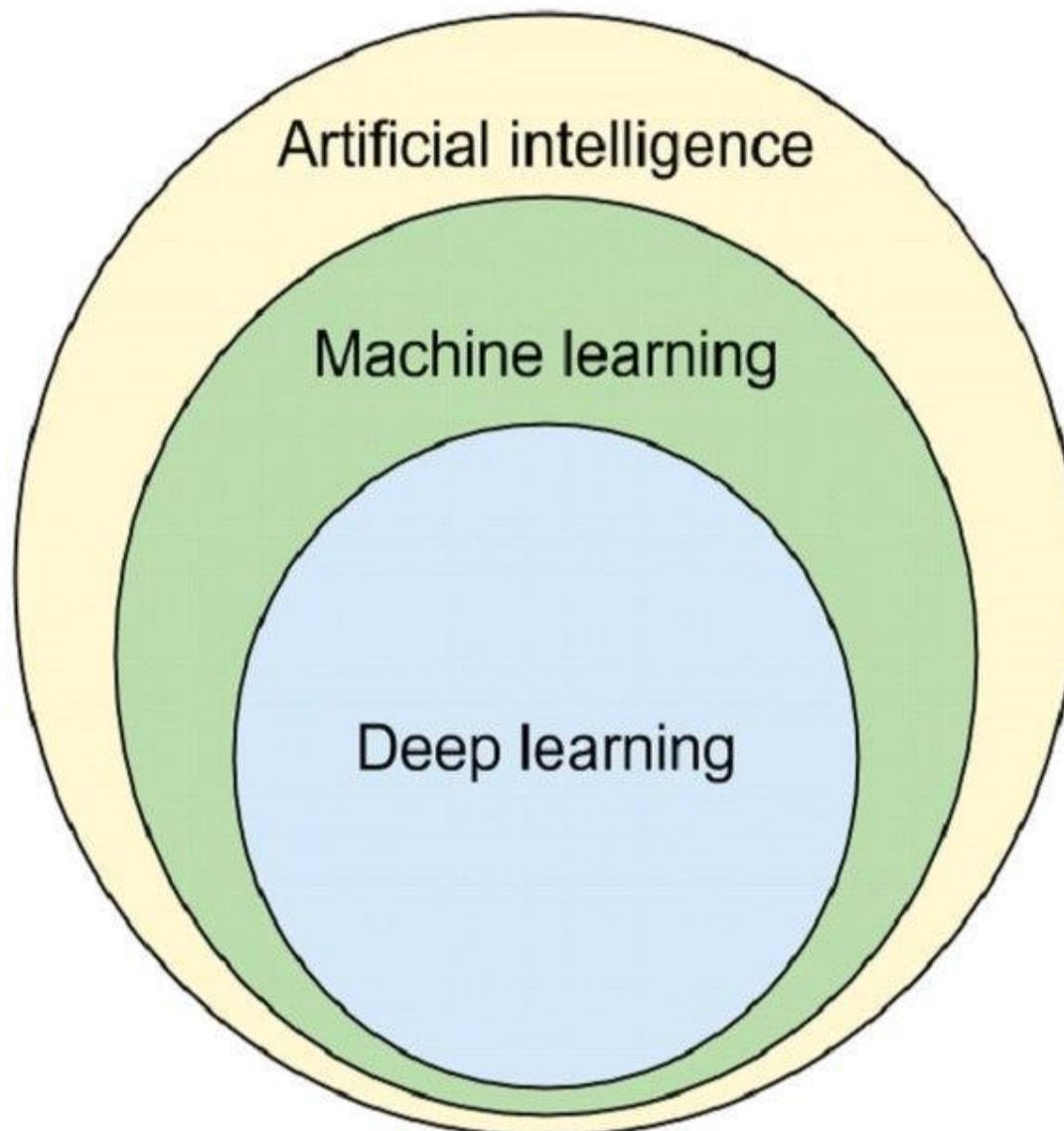
Penn's Autonomous Car →  
(Ben Franklin Racing Team)



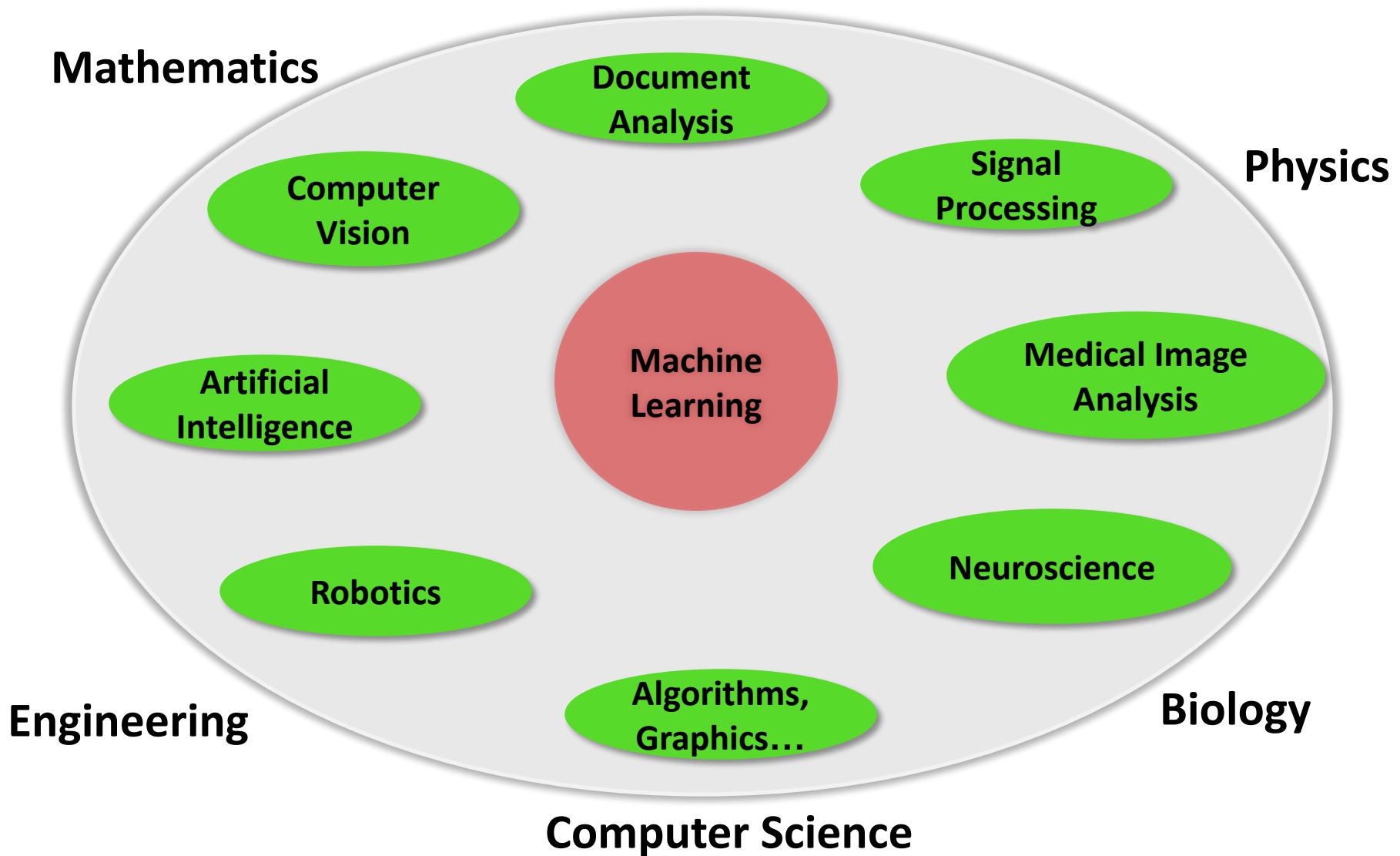
# Computer-aided Diagnosis



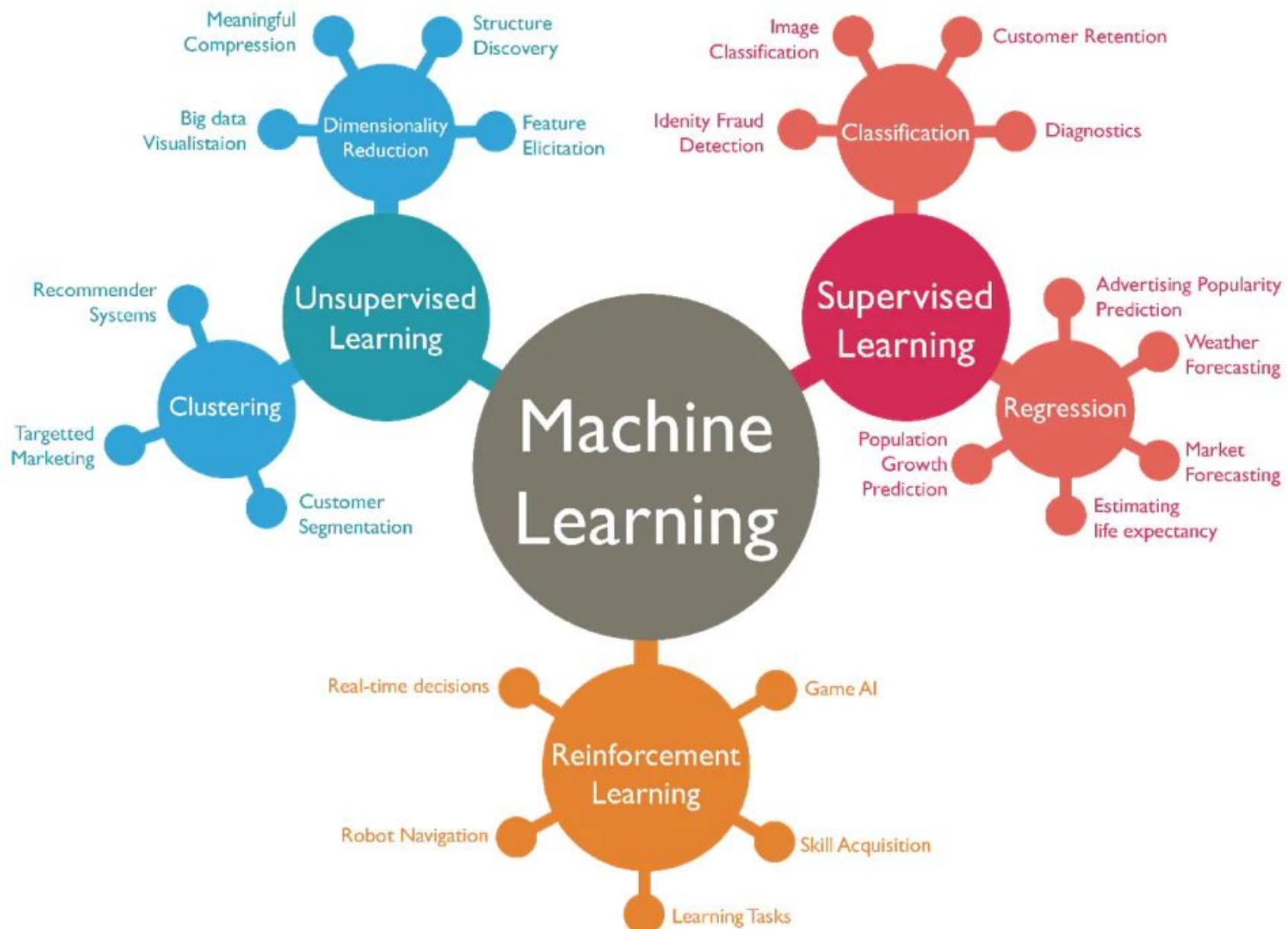
# Relation with AI and DL



# Relation with Other Fields



# Different Machine Learning Paradigms

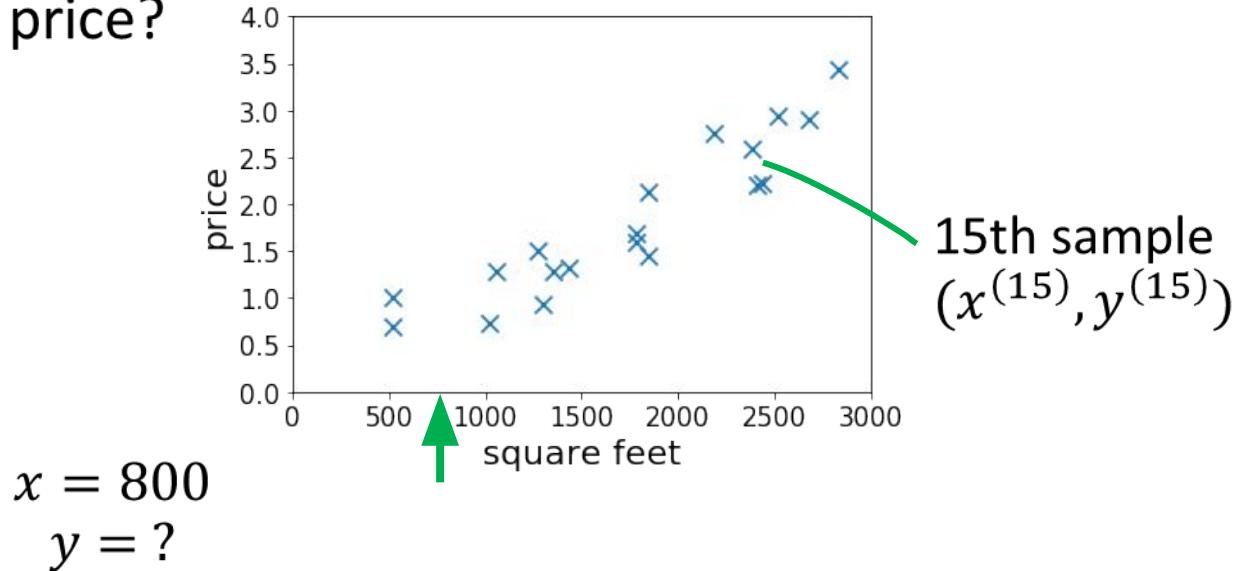


# Supervised Learning

- To learn an unknown target function  $f$
- Input: a training set of labeled examples  $(x_j, y_j)$  where  $y_j = f(x_j)$
- E.g.,  $x_j$  is an image,  $f(x_j)$  is the label “giraffe”
- Output: hypothesis  $h$  that is “close” to  $f$ , i.e., predicts well on unseen examples (“test set”)
- Many possible hypothesis families for  $h$  – Linear models, logistic regression, neural networks, decision trees, examples (nearest neighbors) etc.

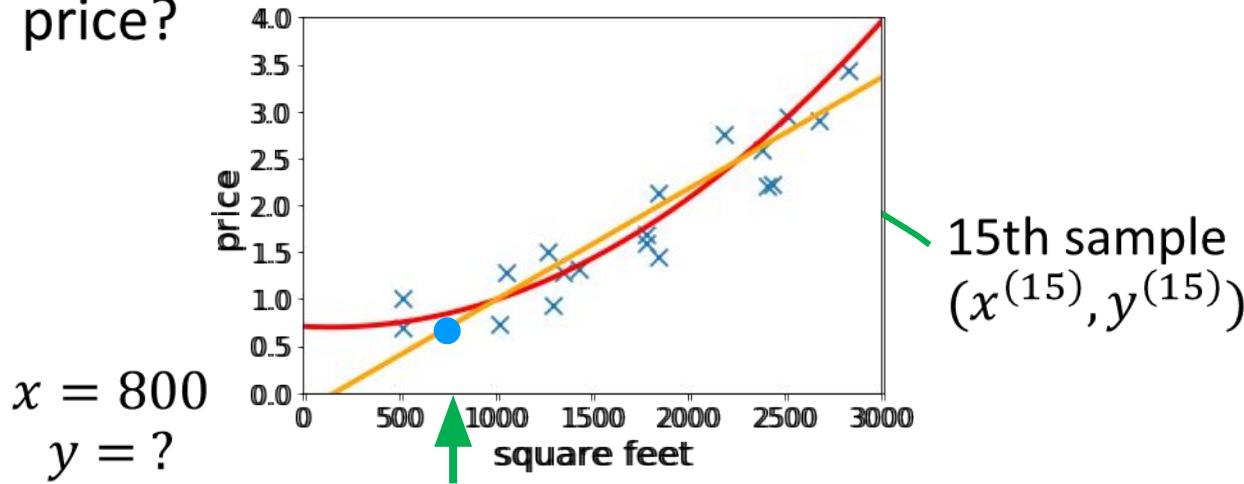
# Housing Price Prediction

- Given: a dataset that contains  $n$  samples
$$(x^{(1)}, y^{(1)}), \dots (x^{(n)}, y^{(n)})$$
- Task: If a residence has  $x$  square feet, predict its price?



# Housing Price Prediction

- Given: a dataset that contains  $n$  samples
$$(x^{(1)}, y^{(1)}), \dots (x^{(n)}, y^{(n)})$$
- Task: If a residence has  $x$  square feet, predict its price?



- Solution: fitting linear/quadratic functions to the dataset.

## High-dimensional Features

- $x \in \mathbb{R}^d$  for large  $d$
- E.g.,  
$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ \vdots \\ x_d \end{bmatrix} \rightarrow \begin{array}{l} \text{--- living size} \\ \text{--- lot size} \\ \text{--- # floors} \\ \text{--- condition} \\ \text{--- zip code} \\ \vdots \end{array} \quad y \text{ --- price}$$

# Supervised Learning in CV

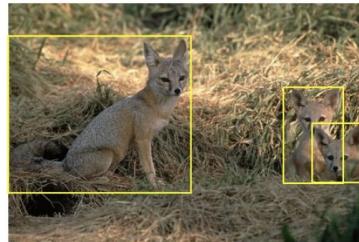
- Image Classification
  - $x$  = raw pixels of the image,  $y$  = the main object



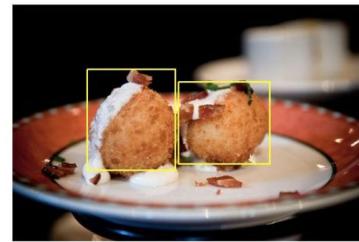
ImageNet Large Scale Visual Recognition Challenge. Russakovsky et al.'2015

# Supervised Learning in CV

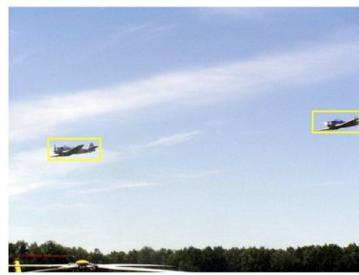
- Object localization and detection
  - $x$  = raw pixels of the image,  $y$  = the bounding boxes



kit fox



croquette



airplane



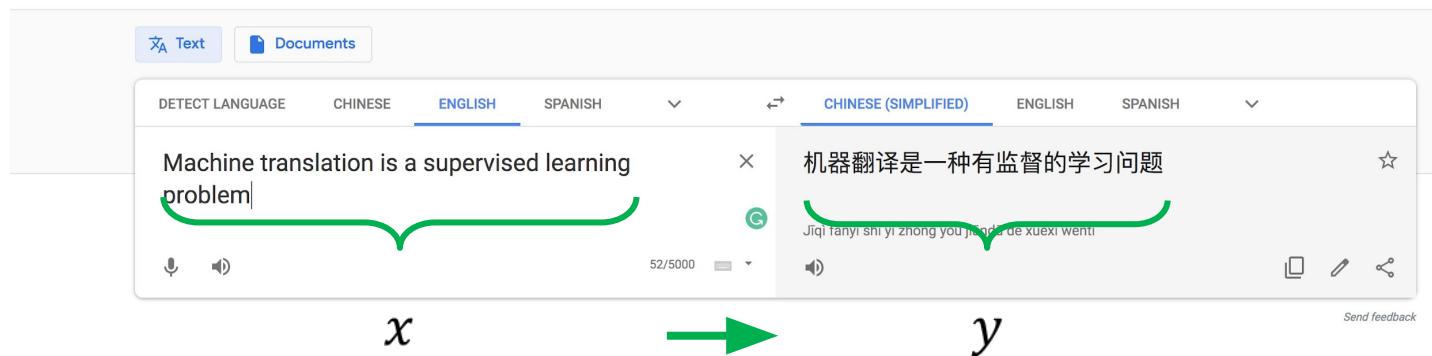
frog

ImageNet Large Scale Visual Recognition Challenge. Russakovsky et al.'2015

# Supervised Learning in Natural Language Processing

- Machine translation

Google Translate



- **Note:** This course only covers the basic and fundamental techniques of supervised learning (which are not enough for solving hard vision or NLP problems.)

# Supervised Learning

## Advantage

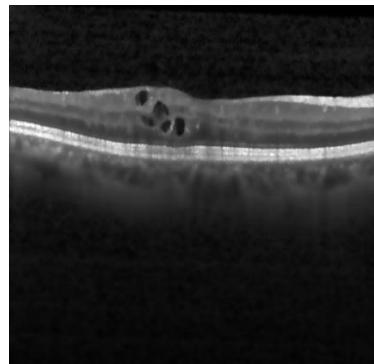
- Well established field.
  - has good quantified methods (often with theoretical bounds)
- You can easily test and debug your learning machine.
  - Since the labelled data is available you can easily inspect its output and find out what errors it's making on what type of input data.

# Supervised Learning

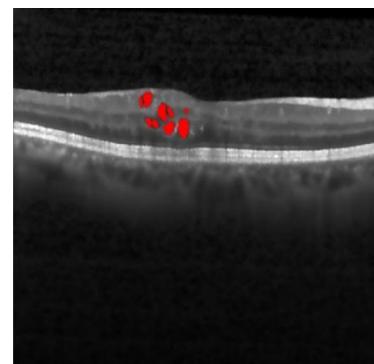
## Disadvantage

- Collecting and labelling data is expensive and time-consuming.

Example: Speech Recognition, Medical Image Analysis, etc.



Original Retinal Scan

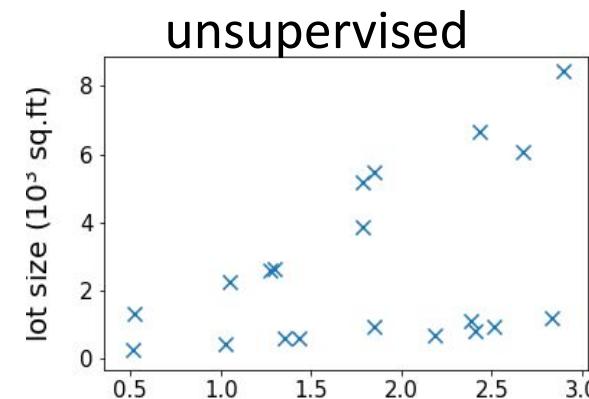
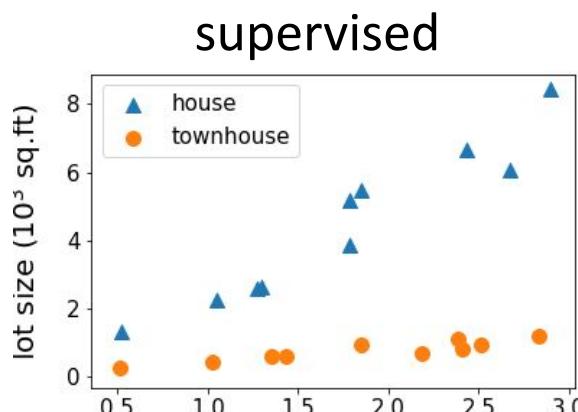


Intra-retinal cysts annotated scan

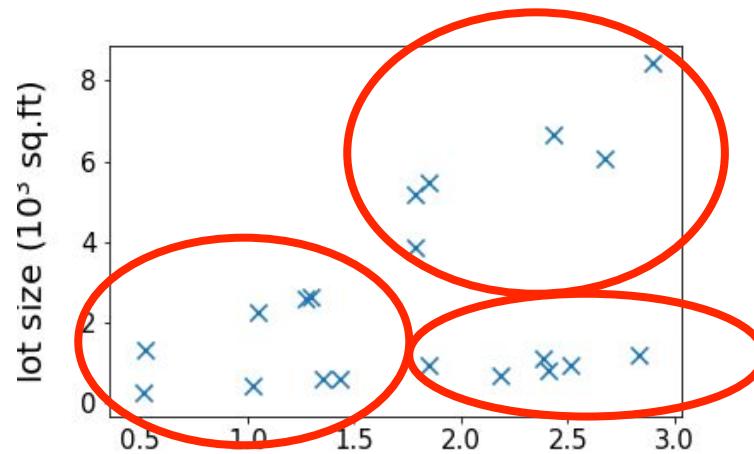
- Errors in your training data might confuse your algorithm and lower its accuracy. Garbage-in -> Garbage-out

# Unsupervised Learning

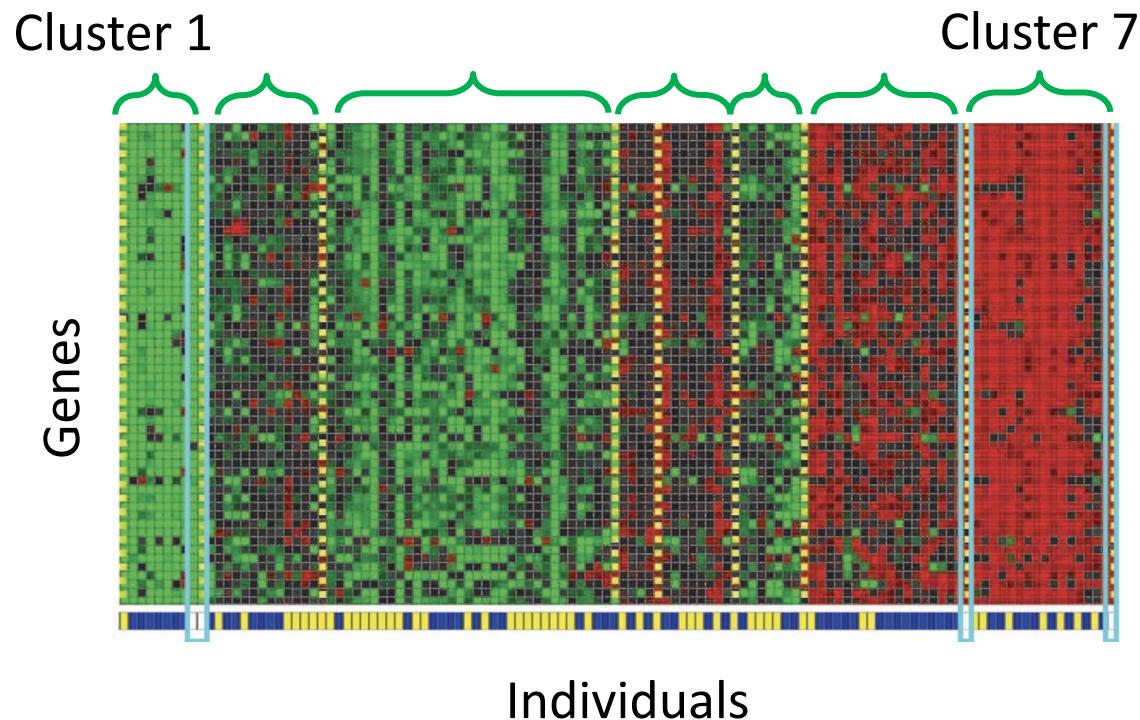
- Dataset contains **no labels**:  
 $x^{(1)}, \dots x^{(n)}$
- **Goal** (vaguely-posed): to find interesting structures in the data



# Clustering



# Clustering Genes



Identifying Regulatory Mechanisms using Individual Variation Reveals Key Role for Chromatin Modification. [Su-In Lee, Dana Pe'er, Aimee M. Dudley, George M. Church and Daphne Koller. '06]

## Need for Unsupervised Learning

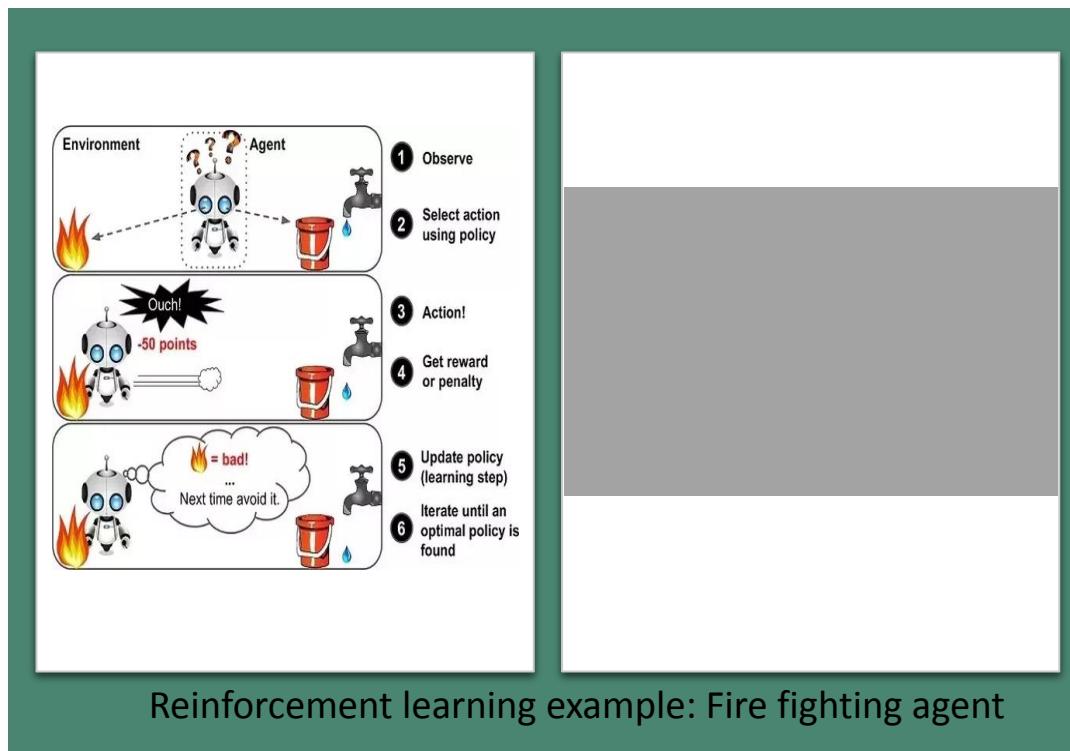
- Annotating large datasets is very costly and time consuming. Example: Speech Recognition, Medical Image Analysis, etc.
- There may be cases where we don't know how many/what classes the data is divided into.  
Example: Data Mining, Sentimental Analysis.
- We may want to use clustering to gain some insight into the structure of the data before designing a classifier.

# Disadvantages of Unsupervised Learning

- Unsupervised Learning is harder as compared to Supervised Learning. Since, quantifying the inference made is difficult, since the ground-truth is missing.
- How do we know if results are meaningful since it has unlabelled data?
  - External evaluation- Expert analysis.
  - Internal evaluation- Objective function.

# Reinforcement Learning

A reinforcement learning algorithm, or agent, learns by interacting with its environment. The agent receives rewards by performing correctly and penalties for performing incorrectly.



# Need for Reinforcement Learning

- Reinforcement learning can be used to solve very complex problems that cannot be solved by conventional techniques.
- In the absence of a training dataset, it is bound to learn from its experience.
- Reinforcement learning models can outperform humans in many tasks and learning process is similar to human learning.
- DeepMind's AlphaGo program, a reinforcement learning model, beat the world champion *Lee Sedol* at the game of *Go* in March 2016.

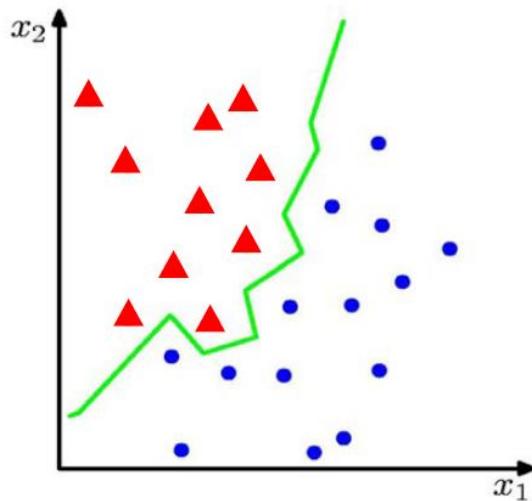
# Disadvantages of Reinforcement Learning

- Reinforcement learning needs a lot of data and a lot of computation. It is data-hungry.
  - So for solving video games and puzzles it performs well.
- Reinforcement learning assumes the world is Markovian, which it is not.
  - The Markovian model describes a sequence of possible events in which the probability of each event depends only on the state attained in the previous event.

# What is classification problem?

- Let there are two classes of objects.
  - Class 1: Set of dog pictures
  - Class 2: Set of cat pictures
- Problem is
  - Given a picture, you should say whether it is cat or dog.
  - For a human being it is easy..., but for a machine it is a non-trivial problem.

# What is classification problem?



- Suppose we are given a training set of N observations  $(x_1, \dots, x_N)$  and  $(y_1, \dots, y_N)$ ,  $x_i \in \mathbb{R}^d$ ,  $y_i \in \{-1, 1\}$
- Classification problem is to estimate  $f(x)$  from this data such that

$$f(x_i) = y_i$$

# Classification: Supervised Learning

## Training Phase

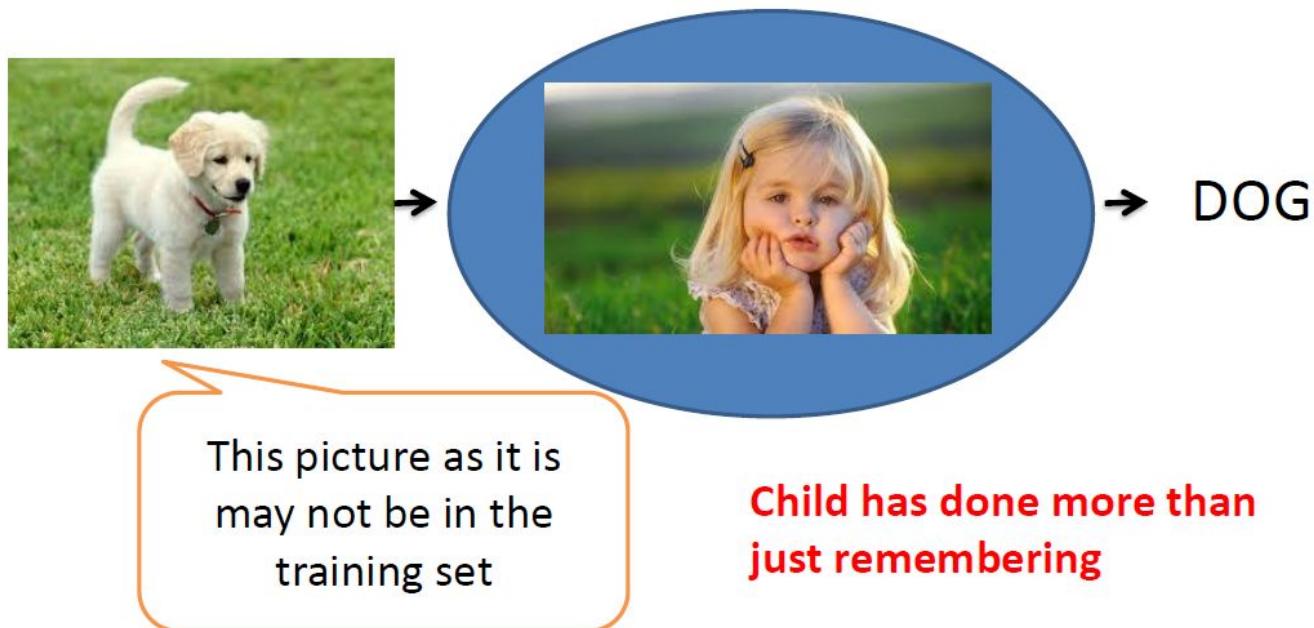


We have shown a set of dog pictures and a set of cat pictures to a child.



# Classification: Supervised Learning

## Testing Phase

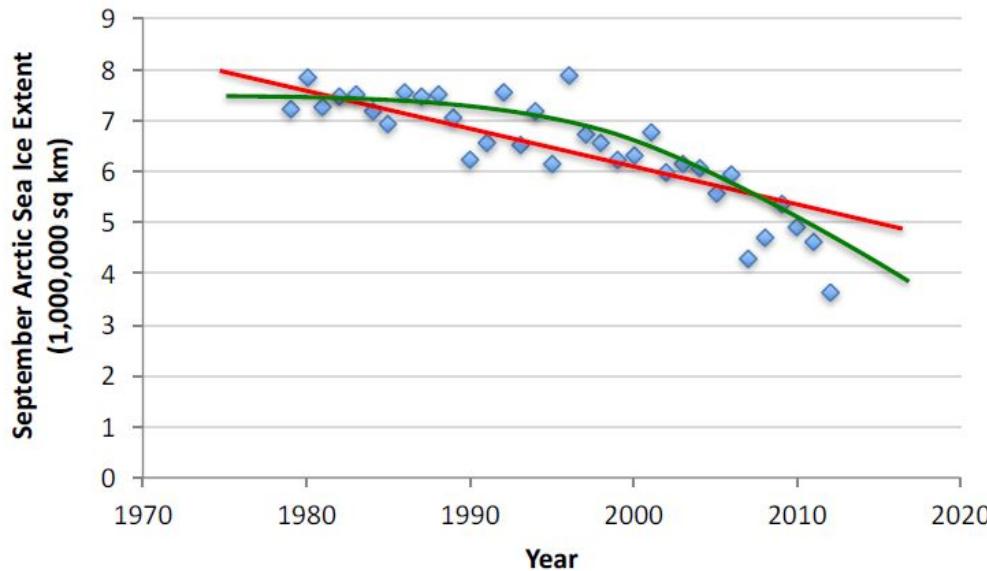


# What is Learning?

- Child has learnt what is it that is common among dogs ... and, what is it that is common among cats... also, what are the distinguishing features/attributes.
- Child has learnt the pattern (regularity) behind all dogs and the pattern behind all cats.
- Child then recognized a test image as having a particular pattern that is unique to dogs.

# What is Regression Problem?

- Given  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Learn a function  $f(x)$  to predict  $y$  given  $x$ 
  - $y$  is real-valued == regression



# Popular ML algorithms

## Classification

- Linear Classifiers
- Support Vector Machines
- Decision Trees
- K-Nearest Neighbor
- Random Forest

## Regression

- Linear Regression
- Logistic Regression
- Polynomial Regression

## Resources: Journals

- Journal of Machine Learning Research [www.jmlr.org](http://www.jmlr.org)
- Machine Learning
- IEEE Transactions on Neural Networks
- IEEE Transactions on Pattern Analysis and Machine Intelligence
- Annals of Statistics
- Journal of the American Statistical Association

## Resources: Conferences

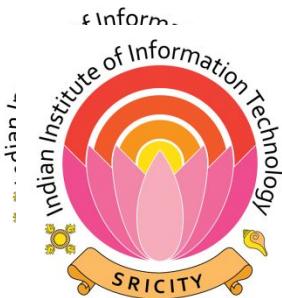
- International Conference on Machine learning (ICML)
- European Conference on Machine Learning (ECML)
- Neural Information Processing Systems (NIPS)
- Computational Learning
- International Joint Conference on Artificial Intelligence (IJCAI)
- ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)
- IEEE Int. Conf. on Data Mining (ICDM)

**Thank You:  
Question?**

# **Machine Learning**

## **Introduction to ML Model Evaluation**

IIIT Sri City



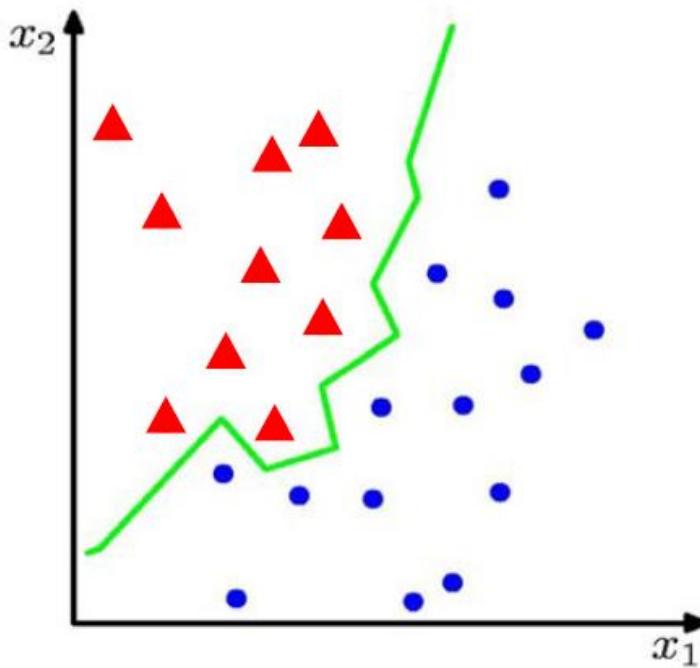
# This week's Agenda

- Recap to classification?
- How to evaluate the classification model?
- Evaluation Metrics for a Classification model
- Recap to regression problem?
- How to evaluate the regression model?
- Evaluation Metrics for a regression model
- Dataset: Train, Validation and Test sets
- Train, test and validation split
- Data Sampling Methods
- Overfit and Underfit

# Introduction to parametric and non-parametric machine learning models.

- A learning model that summarizes data with a set of parameters of fixed size (independent of the number of training examples) is called a parametric model.
  - No matter how much data you throw at a parametric model, it won't change its mind about how many parameters it needs.
- Some more examples of parametric machine learning algorithms include:
  - Logistic Regression
  - Linear Discriminant Analysis
  - Perceptron
- Algorithms that do not make strong assumptions about the form of the mapping function are called nonparametric machine learning algorithms.
  - Nonparametric methods are good when you have a lot of data and no prior knowledge, and when you don't want to worry too much about choosing just the right features.
- Some more examples of popular nonparametric machine learning algorithms are:
  - k-Nearest Neighbors
  - Decision Trees like CART and C4.5
  - Support Vector Machines

# Recap: What is classification problem?



- Suppose we are given a training set of  $N$  observations  $(x_1, \dots, x_N)$  and  $(y_1, \dots, y_N)$ ,  $x_i \in \mathbb{R}^d$ ,  $y_i \in \{-1, 1\}$
- Classification problem is to estimate  $f(x)$  from this data such that

$$f(x_i) = y_i$$

# Classification: Supervised Learning

## Training Phase



We have shown a set of dog pictures and a set of cat pictures to a child.



# Classification: Supervised Learning

## Testing Phase



→ DOG

A black arrow pointing from the oval to the word "DOG".

This picture as it is  
may not be in the  
training set

**Child has done more than  
just remembering**

# Evaluating a classifier

# Classification Accuracy and Classification Error

Let the test set size be  $n$ .

Classification Accuracy (CA):

Number of correct predictions

$$CA = \frac{\text{Number of correct predictions}}{N}$$

Number of wrong predictions

$$Error = \frac{\text{Number of wrong predictions}}{N}$$

# CA and Error

In the previous formula we assumed that the test examples are all equally likely.

In other words, the Probability of drawing any of the test example is the same.

However, this need not be the case always.

# Recap: Probability

Probability theory is built on 3 axioms.

The terms one should be clear:

1. Experiment
2. Sample space
3. Event
4. Probability

# Kolmogorov Axioms of Probability

To the probabilities of outcomes/events of an experiment must obey the axioms:

- **Axiom 1:** For any event  $A$ ,  $\Pr(A) \geq 0$
- **Axiom 2:**  $\Pr(\Omega) = 1$
- **Axiom 3:** For a collection of mutually exclusive events,  $A_1, A_2, \dots, A_n$

$$\Pr(A_1 \cup A_2 \dots \cup A_n) = \sum_{i=1}^n \Pr(A_i)$$

- Everything else in probability theory can be deduced starting with these axioms

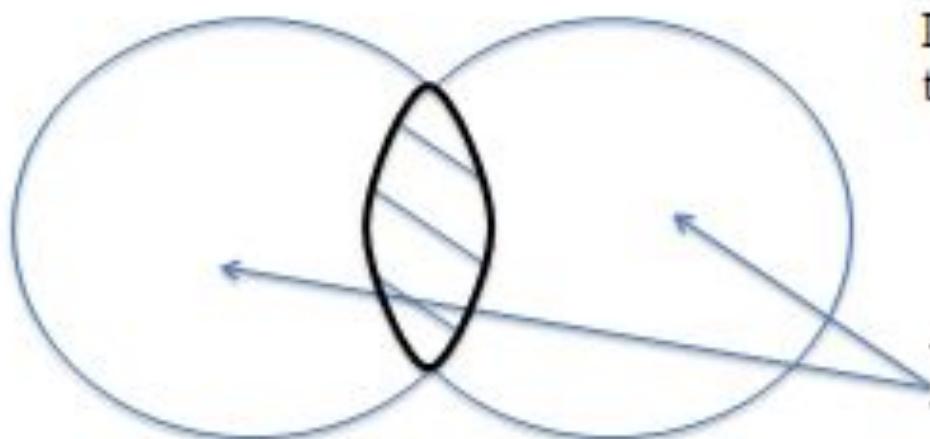
## Handy Consequences of Kolmogorov Axioms

- Important consequences:
  - Probability of a union of non-disjoint events

$$\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B)$$

In words: The probability of A *or* B is the probability of A  
plus the probability of B minus the probability of A *and* B

Don't count the probabilities of A and B  
twice if there is overlap between the events



$$\Pr(A \cup B)$$

# Loss function

Let  $X$  be the given example with  $y$  being the target label (ground-truth).

Let the classifier predicted  $y'$

Loss of this prediction,  $P(\text{error}|X)$  and  $P(\text{error})$  :

$$\mathcal{L}(x) = \begin{cases} 0, & \text{if } y = y' \\ 1, & \text{otherwise} \end{cases}$$

$$p(\text{error}|x) = \mathcal{L}(x)$$

{This assumes that the feature space  
is a discrete one}

$$\begin{aligned} p(\text{error}) &= E[p(\text{error}|x)] \\ &= \sum_x \mathcal{L}(x) p(x) \end{aligned}$$

If the space is continuous,

$$p(\text{error}|x) = L(x)$$

$$\begin{aligned} p(\text{error}) &= E[p(\text{error}|x)] \\ &= \int L(x) p(x) dx \end{aligned}$$

{The integration is over the feature space}

# Example (discrete case)

Data point	Ground truth	Predicted label	Probability
X1	Cat	Cat	1/8
X2	Cat	Dog	1/4
X3	Dog	Dog	1/4
X4	Dog	Cat	3/8

$$P(\text{error}) = \text{Average Error rate} = 0(\frac{1}{8}) + 1(\frac{1}{4}) + 0(\frac{1}{4}) + 1(\frac{3}{8}) = \frac{5}{8}.$$

Note: The entire space is covered by these four examples.

## Eg 2

Assume we are working with 1D, two class problem. Let the class labels are 1 and 2.

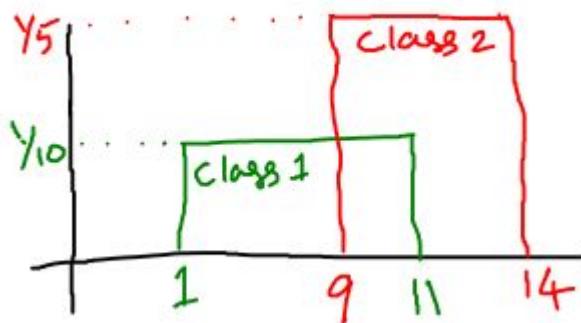
Let the classifier is a simple rule: if the feature value ( $x$ ) is less-than 10 the class assigned is 1, else the class assigned is 2.

$$x < 10 \Rightarrow \text{class} = 1$$

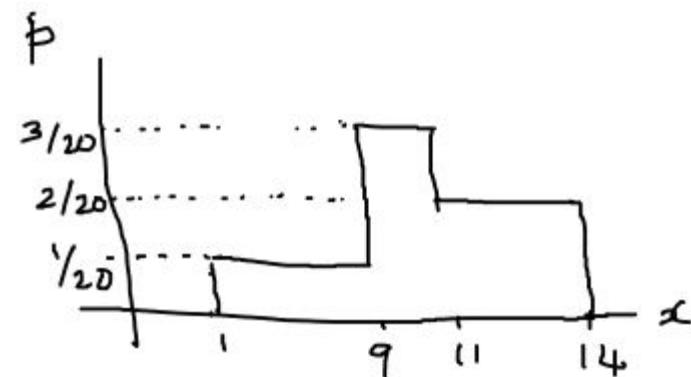
$$x \geq 10 \Rightarrow \text{class} = 2$$

## Eg 2 Continued

Assume that the data is drawn from the mixture of the two class-conditionals (shown below) where each class is equally likely.



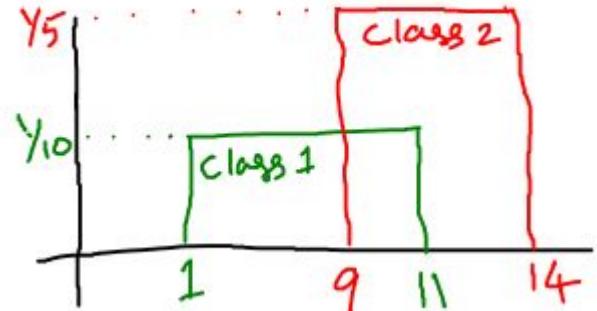
Class-conditional distributions



Distribution from which  $x$  is drawn

## Eg 2

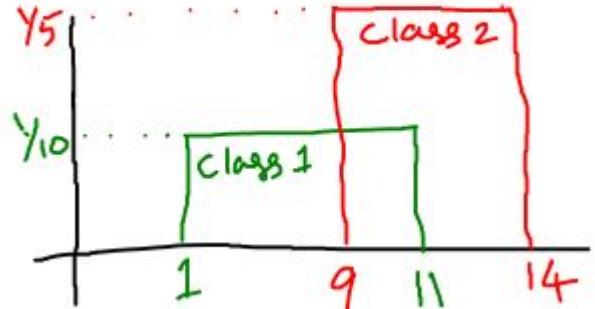
The classifier is not making any mistake for  $x < 9$  and for  $x > 11$ . Only when  $x$  is in  $[9, 11]$  one has to worry about the mistakes made.



## Eg 2

[9, 10] error: Let X is an arbitrary element drawn in [9,10]

$$P(\text{err} \mid X) = P(y=\text{class2} \mid X)$$

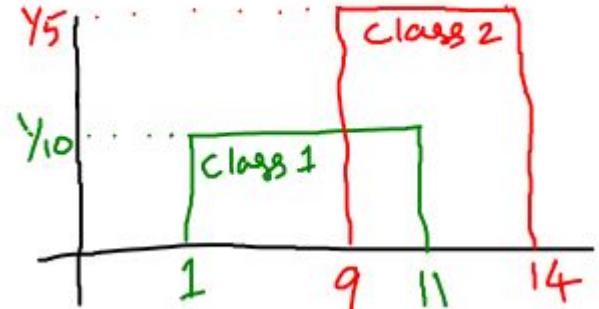


We use the notation  $P(\text{class2} \mid X)$  instead of  $P(y=\text{class2} \mid X)$

## Eg 2

[9, 10] error: Let X is an arbitrary element drawn in [9,10]

$$P(\text{err} \mid X) = P(\text{class2} \mid X)$$



$$P(\text{class2} \mid X) = \frac{p(X, \text{class2})}{p(X)}$$

$$= \frac{p(X \mid \text{class2}) P(\text{class2})}{p(X)}$$

$$p(X) = p(X|class1)P(class1) + p(X|class2)P(class2)$$

So, for  $X$  in  $[9,10]$ ,

$$P(class2|X) = \frac{\frac{1}{5} \cdot \frac{1}{2}}{\frac{1}{5} \cdot \frac{1}{2} + (1/10) \cdot \frac{1}{2}} = \frac{2}{3}$$

[10,11] error: Let X is an arbitrary element drawn in [10,11]

$$\begin{aligned} P(\text{err} \mid X) &= P(\text{class1} \mid X) \\ &= \frac{1}{3} \end{aligned}$$

# P(error)

$$P(\text{error}) = \int P(\text{error}(x)) p(x) dx$$

$$= \int_9^9 P(\text{class 2}|x) p(x) dx$$

$$+ \int_1^{10} P(\text{class 2}|x) p(x) dx$$

$$+ \int_{10}^{11} P(\text{class 1}|x) p(x) dx$$

$$+ \int_{11}^{14} P(\text{class 1}|x) p(x) dx$$

$$P(\text{error}) = 0 + \frac{1}{5} \cdot \frac{1}{2} + \frac{1}{10} \cdot \left(\frac{1}{2}\right) + 0 = \frac{3}{20}$$

## Eg2 Contd

By shifting the threshold of the classifier can we reduce the error rate of the classifier?

If so, what would be the best possible threshold?

This actually points to the Bayes Classifier.

# Actual error-rate

The actual error-rate calculation is theoretical and is possible when the distributions (Probability structure) of the given problem is known.

In reality we will not be having this.

# Test set

Test set is an independently drawn set from the same distribution as the training set.

Since we are given with an example set, normal practice is to divide this into **the training set** and **the test set**.

Note that this division has to be random and in practice often disjoint division is followed.

# Multiple test sets are possible

Since multiple training and test set divisions are possible, often we measure the accuracy or error as the average one over multiple divisions.

This can give us variance or standard deviation of the computed accuracy or error.

Later we will see about **bias and variance** of the learning.

**Other evaluation metrics (especially  
for binary classification)**

Total number of Dog Images (10), Cat (10), N=20

**Confusion Matrix:** Normally this matrix is of size  $K*K$  where  $K$  is number of classes.

N=20		Predicted	
		Dog (+)	Cat (-)
Actual	Dog (+)	7	3
	Cat (-)	4	6

N=20		Predicted	
		Dog (+)	Cat (-)
Actual	Dog (+)	TP	FN
	Cat (-)	FP	TN

True Positives: 7 (Dogs images were classified as Dog)

False Positives (*Type 1 Error*): 4 (Cats images were classified as Dog)

True Negatives: 6 (Cats images were classified as cat)

False Negatives (*Type 2 Error*): 3 (Dogs images were classified as cat)

**Total (N) =TP+FP+FN+TN=20**

# Evaluation Metrics for a classification model

- **Accuracy:** Accuracy is number of correct predictions out of total records.

$$\text{Accuracy} = (TP + TN) / \text{Total} = 13/20 = 65\%.$$

- Misclassification rate or error rate:

$$\text{Error rate} = (FP + FN) / \text{Total} = 7/20 = 35\%.$$

## Accuracy Paradox:

Consider, Total number of Dog Images (19), Cat (1), N=20

		Predicted	
		Dog (+)	Cat (-)
Actual	Dog (+)	19	0
	Cat (-)	1	0

		Predicted	
		Dog (+)	Cat (-)
Actual	Dog (+)	TP	FN
	Cat (-)	FP	TN

$$\text{Accuracy} = (TP + TN) / \text{Total} = (19 + 0) / 20 = 99\%$$

# Evaluation Metrics for a Classification model

- Precision (positive predicted value): Is the fraction of documents retrieved that are relevant.
  - It is the number of positive predictions divided by the total number of positive class values predicted.  
 $Precision=TP/(TP+FP)=19/(19+1)=19/20=99\%.$
  - Your classifier said something is positive, how much precise this decision is?
- Recall (sensitivity or true positive rate): It is the fraction of relevant documents that are retrieved.
  - It is the number of positive predictions divided by the number of positive class values in the test data.  
 $Recall=TP/(TP+FN)=19/(19+0)=100\%.$
  - Fractions of positives that are actually said to be positive by your classifier.

**PRECISION, RECALL ARE USED IN INFORMATION RETRIEVAL SYSTEMS. These can be used to evaluate a binary classification.**

# F-1 Measure

- F-1 Measure: A balanced measure between precision and recall.

$$F-1 \text{ Measure} = 2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$$

$$F-1 = 2 * (0.99 * 1.0) / (0.99 + 1.0) = 2 * (0.99) / (1.99) = 1.98 / 1.99 = 0.99$$

# Evaluation Metrics for a Classification model

- Specificity (True Negative Rate): How often the ML model predicts negative samples correctly out of total negative samples.

$$\text{Specificity} = \text{TN}/(\text{TN}+\text{FP}) = 0/(0+1) = 0\%$$

*-Model has 0% effective in predicting negative samples as negative.*

- False positive rate (FPR): How often the ML model predict the negative samples (cat) as positive (dog)?

$$\text{FPR} = \text{FP}/(\text{TN}+\text{FP}) = 1/(0+1) = 100\%$$

**Note: It can also be calculated as FPR=1-Specificity=1-0=100%**

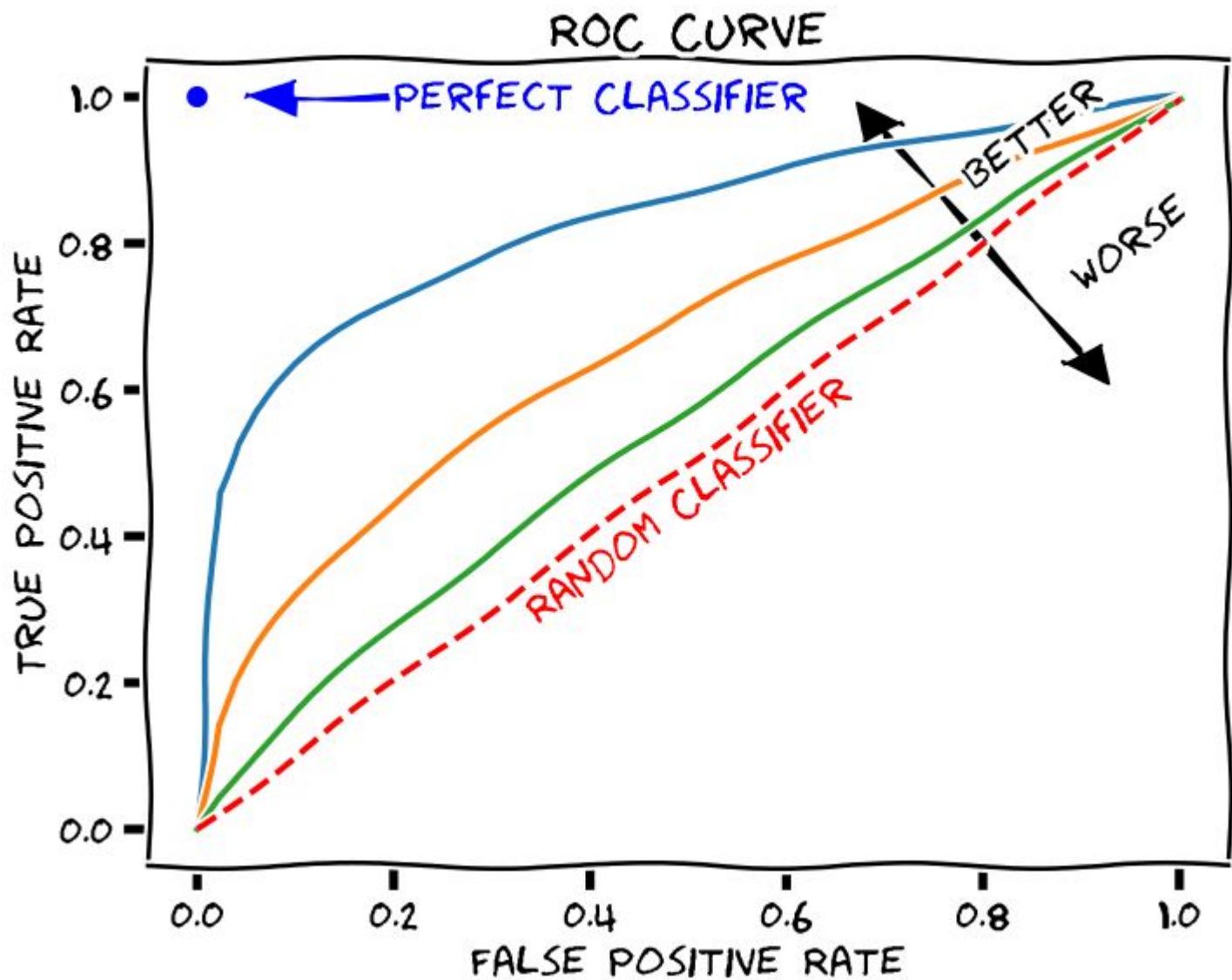
*ML model has 100% FPR, which means that every time model will classifies every Negative (Cat) sample as Positive (Dog).*

- *Sensitivity and Specificity are most important in medical diagnosis.*

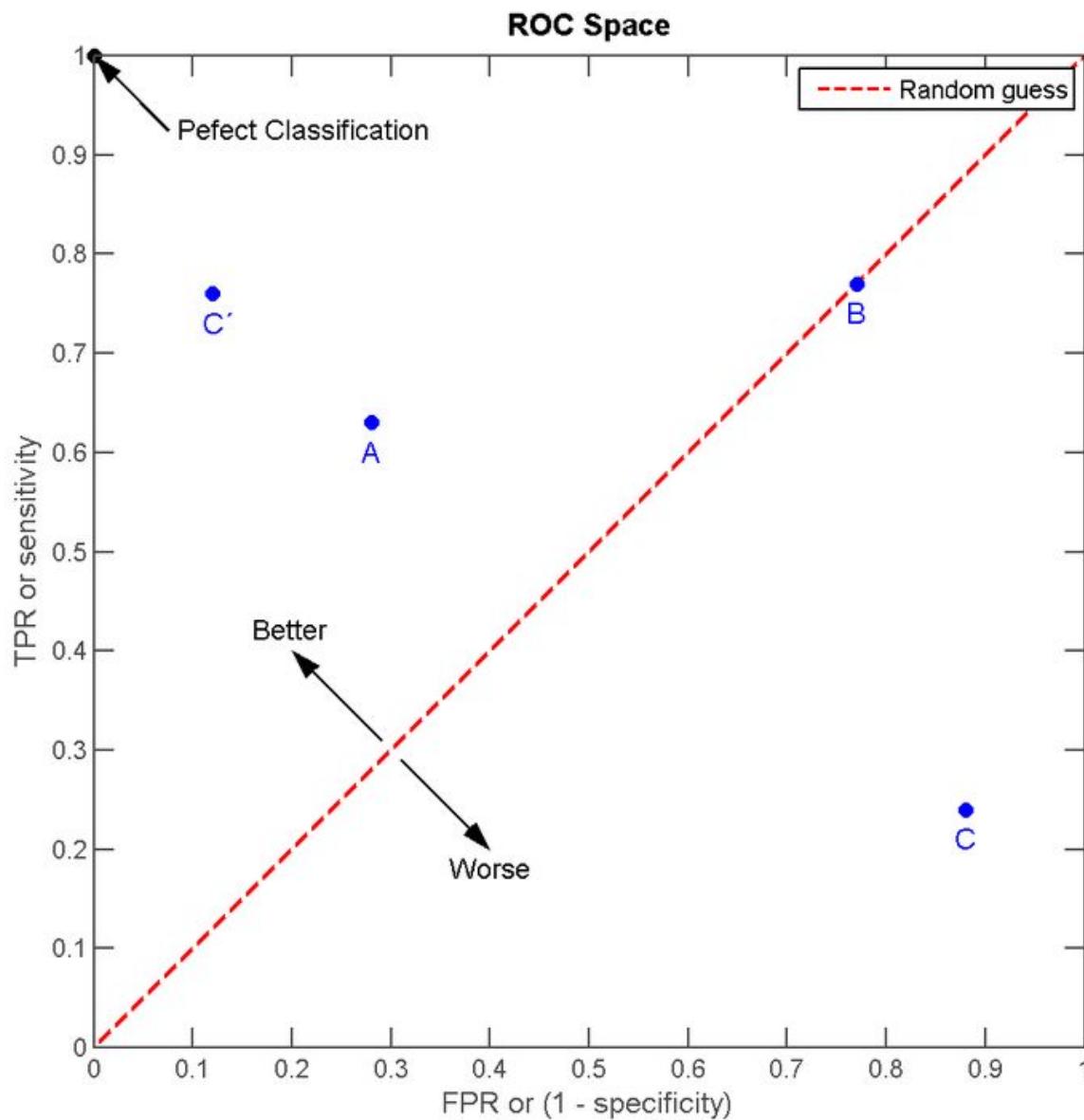
# ROC Curve

- A receiver operating characteristic curve, or ROC curve, is a graphical plot that illustrates the diagnostic ability of a binary classifier system as its discrimination threshold is varied.
- The method was originally developed for operators of military radar receivers, which is why it is so named.
- The ROC curve is created by plotting the true positive rate (TPR) against the false positive rate (FPR) at various threshold settings.

# ROC Curve



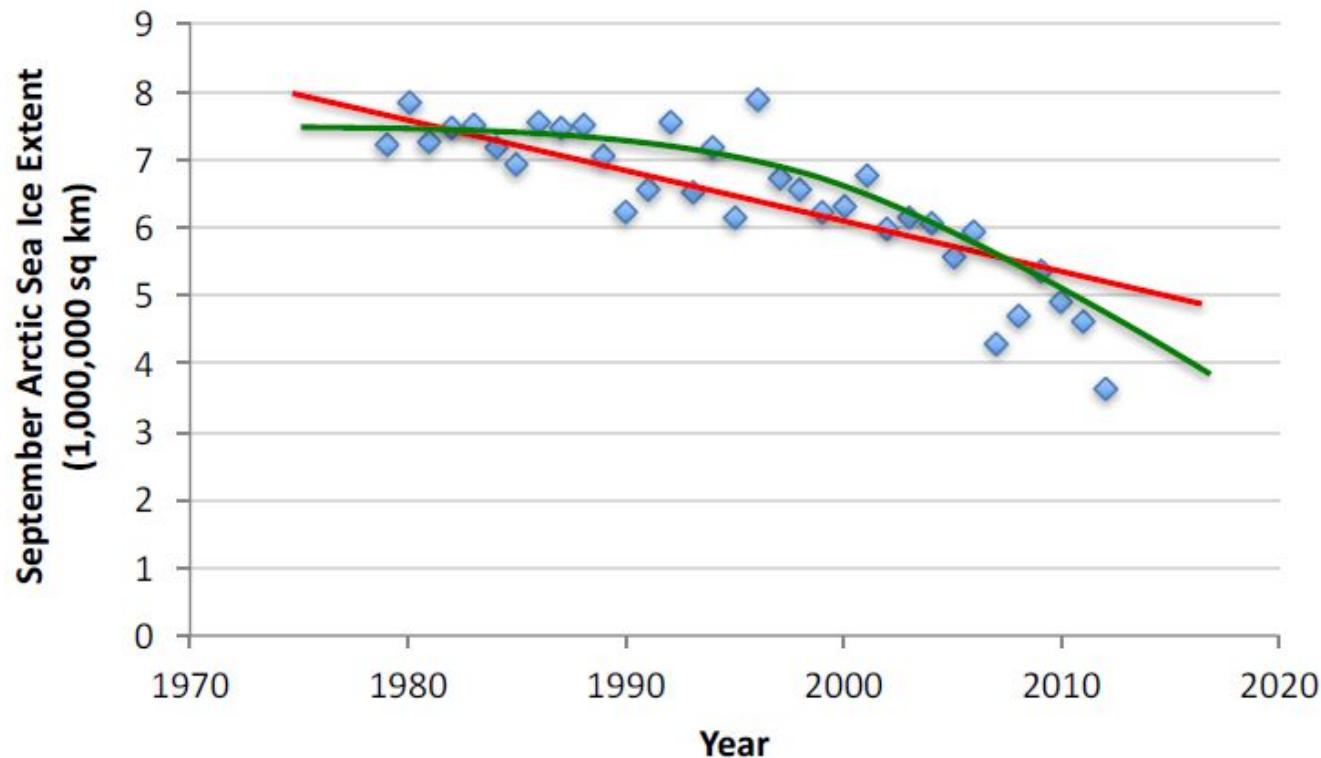
# ROC Curve



# **Regression**

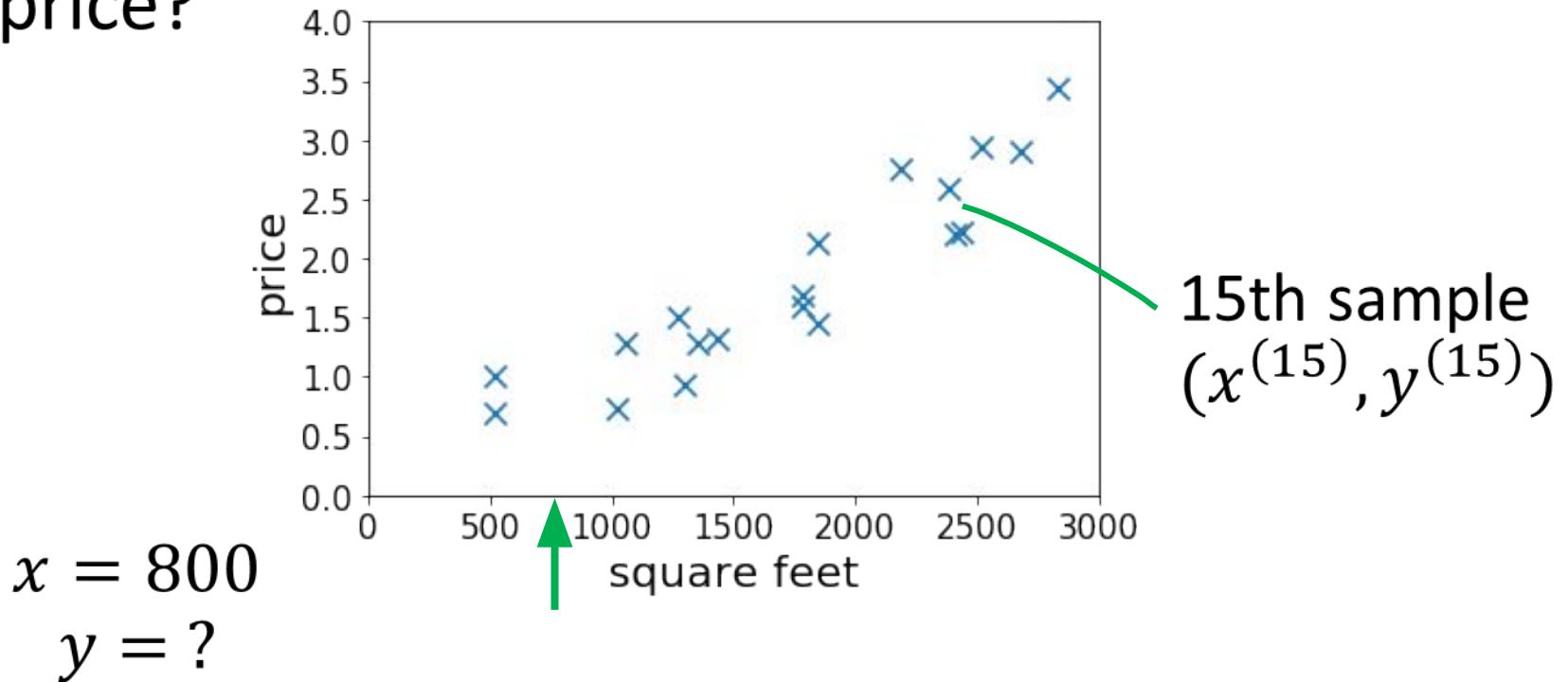
# What is the Regression Problem?

- Given  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Learn a function  $f(x)$  to predict  $y$  given  $x$ 
  - $y$  is real-valued == regression



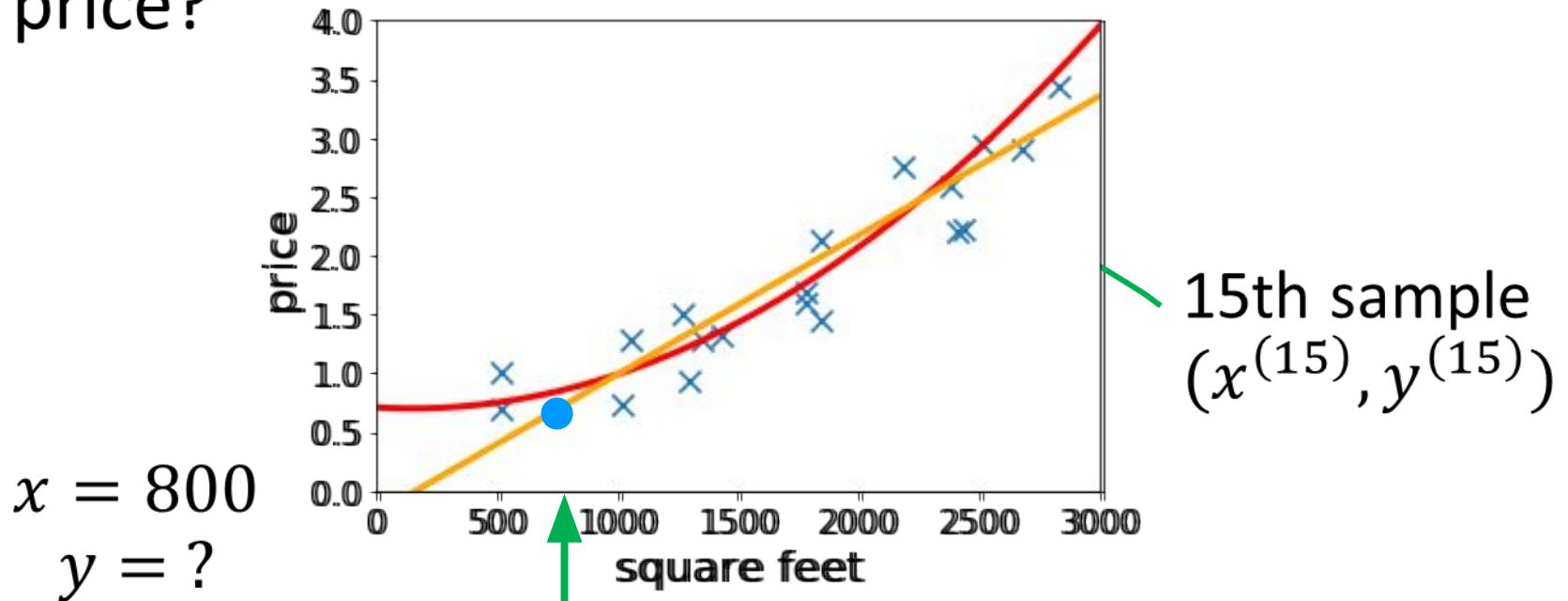
# Housing Price Prediction

- Given: a dataset that contains  $n$  samples  $(x^{(1)}, y^{(1)}), \dots (x^{(n)}, y^{(n)})$
- Task: If a residence has  $x$  square feet, predict its price?



# Housing Price Prediction

- Given: a dataset that contains  $n$  samples  $(x^{(1)}, y^{(1)}), \dots (x^{(n)}, y^{(n)})$
- Task: If a residence has  $x$  square feet, predict its price?



- Solution: fitting linear/quadratic functions to the dataset.

# How to evaluate the regression model?

- What is the accuracy of your model prediction?
- Compared to classification (where you have discrete set of labels as prediction) predicting the accuracy in regression (where you have continuous real value number) is slightly difficult.!
- It might be impossible for your ML model to predict the exact value.
- So to calculate the accuracy, you can compare the predicted value as how close it is against the real value.

# Evaluation Metrics for a Regression model

- *There are three main metrics used for evaluating the regression models:*
  - *R Square/Square of the Correlation Coefficient*
  - *Mean Square Error(MSE)/Root Mean Square Error(RMSE)*
  - *Mean Absolute Error(MAE)*

## Evaluating a regression model: R square value

- *R Square/coefficient of determination:* It measures the proportion of the variance in the dependent variable that is predictable from the independent variable(s).
- *R Square* value is between 0 to 1 and bigger value indicates a better fit between prediction and actual value

$$R^2 = 1 - \frac{SS_{Regression}}{SS_{Total}} = 1 - \frac{\sum_i(y_i - \hat{y}_i)^2}{\sum_i(y_i - \bar{y})^2}$$

Where  $y_i$  is the original value,  $\hat{y}_i$  is the predicted value and  $\bar{y}$  is the mean of original values.

# Evaluating a regression model: Mean Squared Error (MSE)

- MSE is calculated by the sum of square of prediction error which is real output minus predicted output and then divide by the number of data points.
- It gives you an absolute number on how much your predicted results deviate from the actual number.
- Root Mean Square Error(RMSE) is the square root of MSE.

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

Where  $y_i$  is the original value and  $\hat{y}_i$  is the predicted value.

# Evaluating a regression model: Mean Absolute Error (MAE)

- Mean Absolute Error(MAE) is similar to Mean Square Error(MSE).
- Unlike MSE where we take the sum of square of errors, MAE computes the sum of absolute value of error.

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

Where  $y_i$  is the original value and  $\hat{y}_i$  is the predicted value.

- Note: MSE gives larger penalisation to big prediction error by square it while MAE treats all errors the same.

# Dataset: Train, Validation and Test sets

- **Train set:** The model is initially fit on a training dataset,[3] which is a set of examples used to fit the parameters (e.g. weights of connections between neurons in artificial neural networks) of the model.
- In practice, the training dataset often consists of pairs of an input vector (or scalar) and the corresponding output vector (or scalar), where the answer key is commonly denoted as the target (or label).

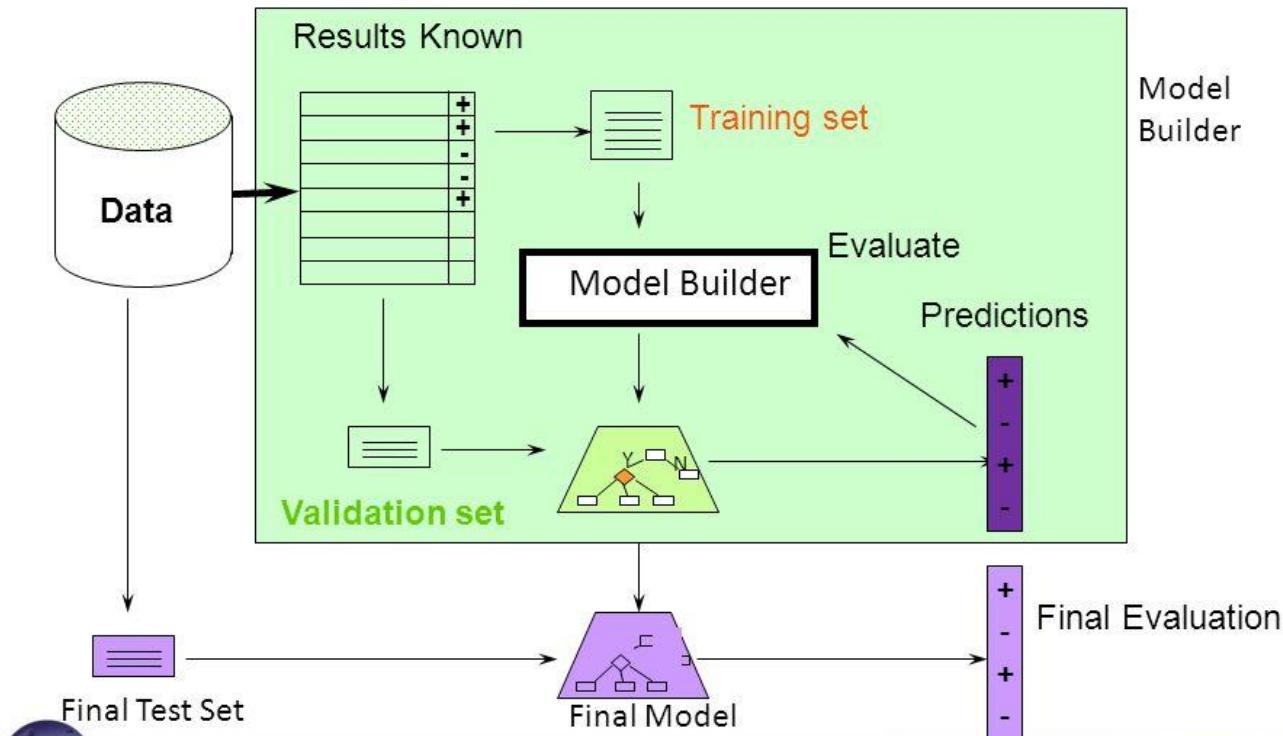
# Dataset: Train, Validation and Test sets

- **Validation set:** Successively, the fitted model is used to predict the responses for the observations in a second dataset called the validation dataset.
- The validation dataset provides an unbiased evaluation of a model fit on the training dataset while tuning the model's hyperparameters e.g. the number of hidden units (layers and layer widths) in a neural network.
- The validation dataset functions as a hybrid: it is training data used for testing, but neither as part of the low-level training nor as part of the final testing.

# Dataset: Train, Validation and Test sets

- **Test set:** Finally, the test dataset is a dataset used to provide an unbiased evaluation of a final model fit on the training dataset.
- The test dataset is typically used to assess the final model that is selected during the validation process.
- It is only used once a model is completely trained(using the train and validation sets).
  - Many a times the validation set is used as the test set, but it is not good practice.
- The test set generally contains carefully sampled data that spans the various classes that the model would face, when used in the real world.

# Train, test and validation split



# Data Sampling Methods

- **Hold out method:** In the holdout method, we randomly assign data points to two sets  $d_0$  and  $d_1$ , usually called the training set and the test set, respectively.
- The size of each of the sets is arbitrary although typically the test set is smaller than the training set.



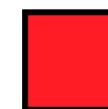
# Data Sampling Methods

- **Cross Validation Method:** Cross-validation involves partitioning a sample of data into complementary subsets, performing the analysis on one subset (called the training set), and validating the analysis on the other subset (Validation set).
- To reduce variability, in most methods multiple rounds of cross-validation are performed using different partitions, and the validation results are combined (e.g. averaged) over the rounds to give an estimate of the model's predictive performance.

$n = 8$



Test



Train

Model 1



# Data Sampling Methods

**k-fold Cross validation:** This procedure has a single parameter called “k” that refers to the number of groups that a given data sample is to be split into.

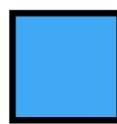
The general procedure is as follows:

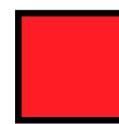
1. Shuffle the dataset randomly.
2. Split the dataset into k groups
3. For each unique group:
  1. Take the group as a hold out or test data set
  2. Take the remaining groups as a training data set
  3. Fit a model on the training set and evaluate it on the test set
  4. Retain the evaluation score and discard the model
4. Summarize the skill of the model using the sample of model evaluation scores

## k-fold Cross validation:

$n = 12$

$k = 3$

 Test

 Train

Data



# Overfit vs Underfit

- **Overfitting:** refers to a model that models the training data too well.
- Overfitting happens when a model learns the detail and noise in the training data to the extent that it negatively impacts the performance of the model on new data.
- In overfitting, ML model learns the inference rules based on the noise or random fluctuations in the training data.

# Overfit vs Underfit

- **Underfitting:** refers to a model that can neither model the training data nor generalize to new data.
- An underfit machine learning model is not a suitable model and will be obvious as it will have poor performance on the training data.
- Underfitting is often not discussed as it is easy to detect given a good performance metric. The remedy is to move on and try alternate machine learning algorithms.

**Thank You:  
Question?**

# **Machine Learning**

## **Non-parametric Algorithms: k-NN Classifier and Parzen Window**

**Indian Institute of Information Technology  
Sri City, Chittoor**



# Agenda

- Parzen Window
- Defining  $R_n$  (region in the feature-space)
- Two different approaches - fixed volume vs. fixed number of samples in a variable volume
- Example 3D hypercube
- The window function and estimation
- Critical parameters of the Parzen-window technique: window width and kernel
- Selecting Window
- Selecting Kernel

# Introduction

Probability density estimation is one of the oldest problems in statistics and machine learning.

There are two approaches, viz.,

1. Parametric, and
2. Non-parametric.

# Parametric Density Estimation

- We assume a parametric class (the form) from which the data is drawn. For eg., Gaussian distribution.
- Then we try to estimate the parameters of the Gaussian distribution ie., mean and covariance matrix.
- For doing the parameter estimation we use the training examples.

We study about this Later.

# Non-parametric Methods

No assumption is made about the form of the distribution.

Depends totally on the data set.

# Non-parametric Methods

1. Parzen Window based
2. Nearest Neighbors based

# Parzen window

- The Parzen-window method (also known as Parzen-Rosenblatt window method) is a widely used non-parametric approach to estimate probability density  $p(x)$ , for a specific point  $x$ .
- Notation: The estimate of  $p(x)$  when we use dataset of size  $n$  is denoted by  $p_n(x)$ .
- It doesn't require any knowledge or assumption about the underlying distribution.
- A popular application of the Parzen-window technique is to estimate the class-conditional densities (or also often called 'likelihoods').
- Likelihoods,  $p(x | \omega_i)$  in a supervised pattern classification problem from the training dataset (where  $p(x)$  refers to the probability density that the sample  $x$  belongs to the particular class  $\omega_i$ ).

# Where would this method be useful?

- Imagine that we are about to design a Bayes classifier for solving a statistical pattern classification task using Bayes' rule:

$$P(\omega_i \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid \omega_i) \cdot P(\omega_i)}{p(\mathbf{x})}$$

$$\Rightarrow \text{posterior probability} = \frac{\text{likelihood} \cdot \text{prior probability}}{\text{evidence}}$$

- If the parameters of the class-conditional densities (also called likelihoods) are known, it is pretty easy to design the classifier.
- Imagine we are about to design a classifier for a pattern classification task where the parameters of the underlying sample distribution are not known.
- Therefore, we wouldn't need the knowledge about the whole range of the distribution; it would be sufficient to know the probability of the particular point, which we want to classify, in order to make the decision.

# Parzen Window

- In parzen window we are going to see how we can estimate this probability from the training sample.
- However, the only problem of this approach would be that we would seldom have exact values - if we consider the histogram of the frequencies for a arbitrary training dataset.
- Therefore, we define a certain region (i.e., the Parzen-window) around the particular value to make the estimate.

[1] *Parzen, Emanuel*. On Estimation of a Probability Density Function and Mode. The Annals of Mathematical Statistics 33 (1962), no. 3, 1065–1076.

[2] *Rosenblatt, Murray*. Remarks on Some Nonparametric Estimates of a Density Function. The Annals of Mathematical Statistics 27 (1956), no. 3, 832–837.

# The most fundamental technique

- The probability  $P$  that a vector  $x$  will fall in a region  $R$  is given by

$$P = \int_{\mathcal{R}} p(\mathbf{x}') d\mathbf{x}'.$$

- $p(x) \approx \frac{\left( \int_R p(x') dx' \right)}{V}$  where  $V$  is the volume of the region  $R$

# Defining the Region $R_n$

- The basis of this approach is to count how many data-points fall within a specified region  $R_n$  (or “window”). Our intuition tells us, that (based on the observation), the probability that one sample falls into this region is:

$$P = \frac{\# \text{ of samples in } R}{\text{Total samples}}$$

- The probability of observing  $k$  points out of  $n$  in a Region  $R_n$  : we consider a **binomial distribution**:

$$P_k = \binom{n}{k} \cdot P^k \cdot (1 - P)^{n-k}$$

- In a binomial distribution, the probability peaks sharply at the mean

# Defining the Region R<sub>n</sub>

- Let  $p(x)$  be the probability density at  $x$ . Let over the small region  $R$  it is uniformly distributed.

$$P = \int_R p(x') dx' = p(x) \cdot V$$

where  $V$  is the volume of the region  $R$ , and if we rearrange those terms, so that we arrive at the following equation:

$$\begin{aligned}\frac{k}{n} &= p(x) \cdot V \\ \Rightarrow p(x) &= \frac{k/n}{V}\end{aligned}$$

- This simple equation above (i.e, the “probability estimate”) lets us calculate the probability density of a point  $x$  by counting how many points  $k$  fall in a defined region (or volume).

To estimate the density at  $\mathbf{x}$ , we form a sequence of regions  $\mathcal{R}_1, \mathcal{R}_2, \dots$ , containing  $\mathbf{x}$  — the first region to be used with one sample, the second with two, and so on. Let  $V_n$  be the volume of  $\mathcal{R}_n$ ,  $k_n$  be the number of samples falling in  $\mathcal{R}_n$ , and  $p_n(\mathbf{x})$  be the  $n$ th estimate for  $p(\mathbf{x})$ :

$$p_n(\mathbf{x}) = \frac{k_n/n}{V_n}.$$

If  $p_n(\mathbf{x})$  is to converge to  $p(\mathbf{x})$ , three conditions appear to be required:

- $\lim_{n \rightarrow \infty} V_n = 0$
- $\lim_{n \rightarrow \infty} k_n = \infty$
- $\lim_{n \rightarrow \infty} k_n/n = 0$ .

# Theoretically it can be shown that

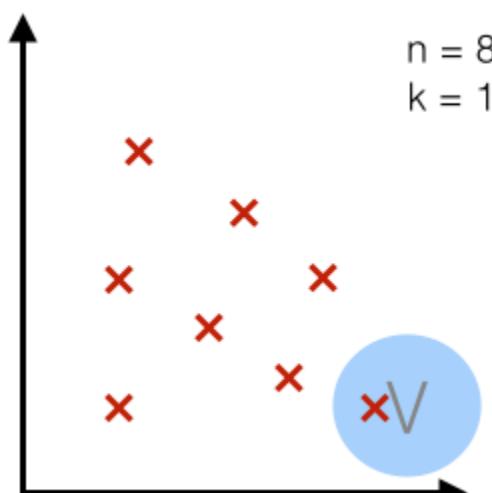
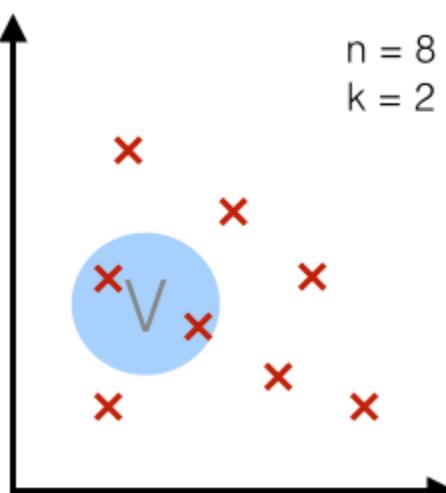
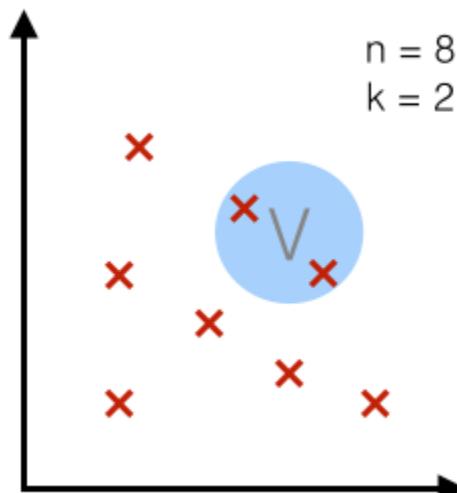
- $V_n$  can be reduced as n increases: Like  $V_n = 1/\sqrt{n}$ . Starting with  $V_1 = 1$ .
- Or,  $k_n$  can be increased as n increases: Like  $k_n = \sqrt{n}$ . Here the volume of the region is increased to fit the  $k_n$  points exactly.
- Both these approaches are shown to converge to the true density *asymptotically*.

**In Practice: Two approaches  
followed are-**

# Two different approaches - fixed volume vs. fixed number of samples in a variable volume

Case 1 - fixed volume:

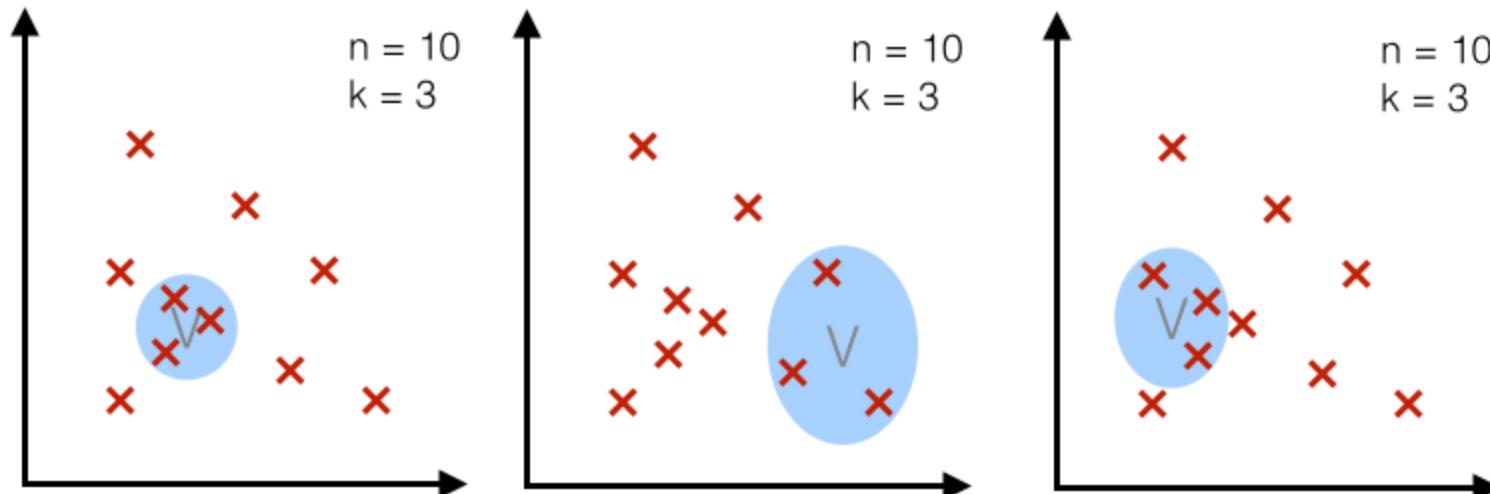
- For a particular number  $n$  (= number of total points), we use volume  $V$  of a fixed size and observe how many points  $k$  fall into the region.



# Two different approaches - fixed volume vs. fixed number of samples in a variable volume

Case 2 - fixed  $k$ :

- For a particular number  $n$  (= number of total points), we use a fixed number  $k$  (number of points that fall inside the region or volume) and adjust the volume accordingly..

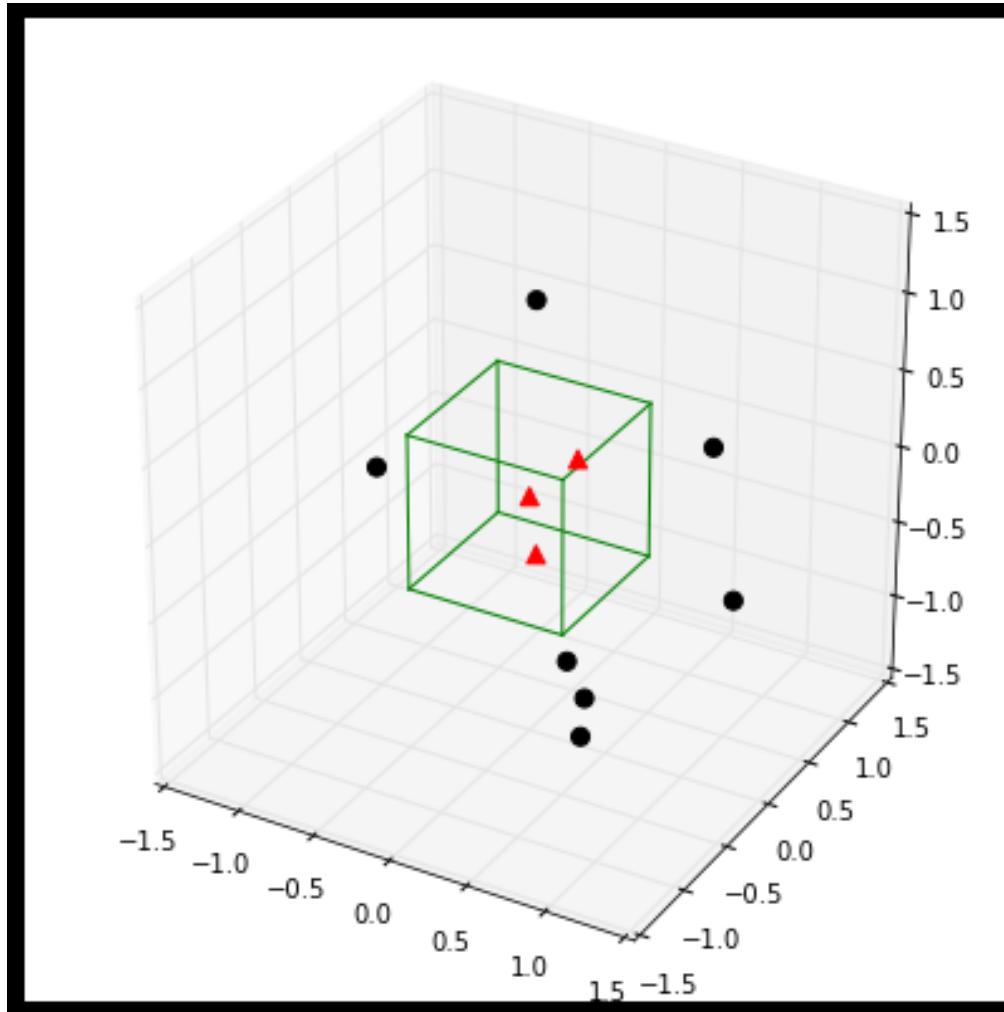


- The second approach, fixed  $k$ , is nothing but the  $k$ -Nearest Neighbor Classifier.
- We will see about this, in detail later.

## Example 3D-hypercubes

- To illustrate this with an example and a set of equations, let us assume this region  $R_n$  is a hypercube.
- The volume of this hypercube is defined by  $V_n = h_n^d$ , where  $h_n$  is the length of the hypercube, and  $d$  is the number of dimensions.
- For an 2D-hypercube with length 1, for example, this would be  $V_1 = 1^2$  and for a 3D hypercube  $V_1 = 1^3$ , respectively.

Example: A typical 3-dimensional unit hypercube ( $h_1 = 1$ ) representing the region  $R_1$ , and 10 sample points, where 3 of them lie within the hypercube (red triangles), and the other 7 outside (black dots).



Credits: "Kernel density estimation via the Parzen-Rosenblatt window method"

By Sebastian Raschka

- Each point falling within the window (hyper-cube) contributes to the density.
- Points falling outside will not.
- We can formalize this in to a window function.

# The window function

- Once we visualized the region  $R_1$  like above, it is easy and intuitive to count how many samples fall within this region, and how many lie outside.
- To approach this problem more mathematically, we would use the following equation to count the samples  $k_n$  within this hypercube, where  $\varphi$  is our so-called window function.

$$\varphi(\mathbf{u}) = \begin{cases} 1 & |u_j| \leq 1/2 ; \quad j = 1, \dots, d \\ 0 & otherwise \end{cases}$$

for a hypercube of unit length 1 centered at the coordinate system's origin.

- If we extend on this concept, we can define a more general equation that applies to hypercubes of any length  $h_n$  that are centered at  $x$ :

$$k_n = \sum_{i=1}^n \varphi\left(\frac{x - x_i}{h_n}\right)$$

$$\text{where } \mathbf{u} = \left( \frac{x - x_i}{h_n} \right)$$

# Parzen-window estimation

- we can now formulate the Parzen-window estimation with a hypercube kernel as follows:

$$p_n(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h^d} \phi\left[\frac{\mathbf{x}-\mathbf{x}_i}{h_n}\right]$$

where

$$h^d = V_n \quad \text{and} \quad \phi\left[\frac{\mathbf{x}-\mathbf{x}_i}{h_n}\right] = k$$

- And applying this to our unit-hypercube example above, for which 3 out of 10 samples fall inside the hypercube (into region  $R$ ), we can calculate the probability  $p(\mathbf{x})$  that  $\mathbf{x}$  samples fall within region  $R$  as follows:

$$\mathbf{x} = \frac{k/n}{h^d} = \frac{3/10}{1^3} = \frac{3}{10} = 0.3$$

# An important observation

- A point falling slightly outside the window will not contribute to the density.
- This is incorrect.
- Also, intuitively, very near point to  $x$  should contribute more to the density than far point (even though both are within the window).

# An important observation

- Each point in the training set should contribute to density.
  - Near contributes more than far.
- This gives smooth estimates
- This avoids selecting the window width problem.
- Empty window problem can be avoided.

# An important Observation

- For  $p(x)$ , let near-by points contribute more, and far-away points less.
- For example one can use a Gaussian function to do this.

# An important Observation

- Assume that a Gaussian  $N(0,1)$  is kept at each of the data points. For example, let  $x_i$  be the data point.
- Let the dataset size is  $n$ .
- Then contribution of  $x_i$  to  $p(x)$  will be
- 

$$\varphi(x - x_i) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-x_i)^2}$$

- If  $x$  is a multivariate data point, say with  $d$  dimensions, then

$$\phi(x - x_i) = \frac{1}{(2\pi)^{d/2}} \exp\left[-\frac{1}{2} \|x - x_i\|^2\right]$$

# The Gaussian Kernel

- The estimation of  $p(x)$  when we have  $n$  data points is:

$$p_n(x) = \frac{1}{n} \sum_{i=1}^n \phi(x - x_i) = \frac{1}{n \cdot (2\pi)^{d/2}} \sum_{i=1}^n \exp\left[-\frac{1}{2} \|x - x_i\|^2\right]$$

- This approach is called “Parzen-window density estimation using the Gaussian Kernel”

Note, division by the volume (of the hypercube !)  $V$  is not appearing here. Since the contribution of each data point to density is itself density.

- Classification example

In classifiers based on Parzen-window estimation:

- We estimate the densities for each category and classify a test point by the label corresponding to the maximum posterior
- The decision region for a Parzen-window classifier depends upon the choice of window function as illustrated in the following figure.

In the next Lecture we will see...

**In Practice, K-NNC is the most used classifier (which is nothing but a non-parametric density estimation based method only!!)**

**Thank You:  
Question?**

# K-Nearest Neighbor Classifier

A non-parametric density estimation  
based classifier

# Introduction

- K-Nearest Neighbor Classifier (k-NNC) is a simple classifier.
- It does not have a design phase, except for the value of k and the distance measure.
- Many people think that this is a weak classifier
  - That is, not a good one.
- In contrary, this is a strong classifier with well established asymptotic bounds.

- Given  $x$  to be classified –
- Let the set of classes be  $\Omega = \{\omega_1, \dots, \omega_c\}$ .
- We find  $P(\omega_i|x)$  for all classes, i.e.,  
 $i = 1, 2, \dots, c$ .
- Note that  $P(\omega_i|x)$  is the Posterior Probability  
that  $x$  belong to the class  $\omega_i$
- The classifier's decision is the class-label  
 $\omega_{max}$  where  $P(\omega_{max}|x) \geq P(\omega_i|x)$  for all  $i$ .

# Posterior Probability

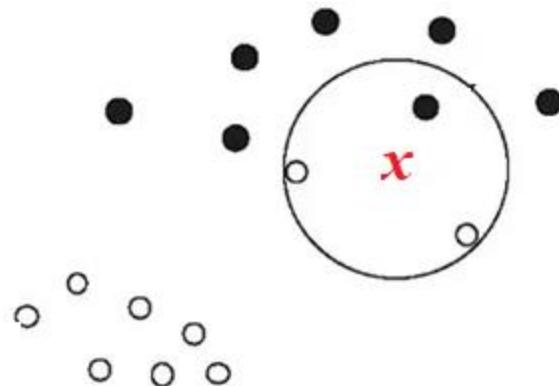
- We first find class-conditional densities, in order to find the Posterior Probability.
- To estimate the class-conditional  $p(x|\omega_i)$  we do:

# Class-conditional

- To estimate the class-conditional  $p(x|\omega_i)$  we do:
  - For the given positive integer  $k$ , draw a hyper-sphere at  $x$  such that exactly  $k$  points are within the sphere.
  - count the number of points from class  $\omega_i$  falling in the sphere. Let this be  $k_i$
  - Let the number of training examples in class  $\omega_i$  be  $n_i$
  - Let the total number of training examples be  $n$

- Then  $p(x|\omega_i) = \frac{k_i}{n_i \cdot V}$
- Here,  $V$  is the volume of the hyper-sphere.

### 3-Nearest Neighbor



In the example,

$$p(x|Black) = \frac{1}{7 \cdot V}$$

$$p(x|White) = \frac{2}{9 \cdot V}$$

# Apriori Probabilities

- We do:  $P(\omega_i) = \frac{n_i}{n}$
- So in the example  $P(Black) = \frac{7}{16}$ , and  
 $P(White) = \frac{9}{16}$ .

# Posterior Probability

$$\begin{aligned} \bullet P(\omega_i|x) &= \frac{p(x|\omega_i)P(\omega_i)}{p(x)} = \frac{\frac{k_i}{n_iV} \cdot \frac{n_i}{n}}{p(x)} \\ &= \frac{k_i}{Vnp(x)} \end{aligned}$$

- Decision :  $\operatorname{argmax}_{\omega} \{P(\omega_1|x), \dots, P(\omega_c|x)\}$
- Decision :  $\operatorname{argmax}_{\omega} \left\{ \frac{k_1}{Vnp(x)}, \dots, \frac{k_c}{Vnp(x)} \right\}$

- Decision :  $\operatorname{argmax}_{\omega} \left\{ \frac{k_1}{Vnp(x)}, \dots, \frac{k_c}{Vnp(x)} \right\}$
- Decision :  $\operatorname{argmax}_{\omega} \{k_1, \dots, k_c\}$

*k-NNC is an approximation of Bayes classifier!*

Mathematically it can be shown that when  $n \rightarrow \infty$ ,  $k \rightarrow \infty$  and  $k/n \rightarrow 0$ , k-NNC is exactly the Bayes classifier.

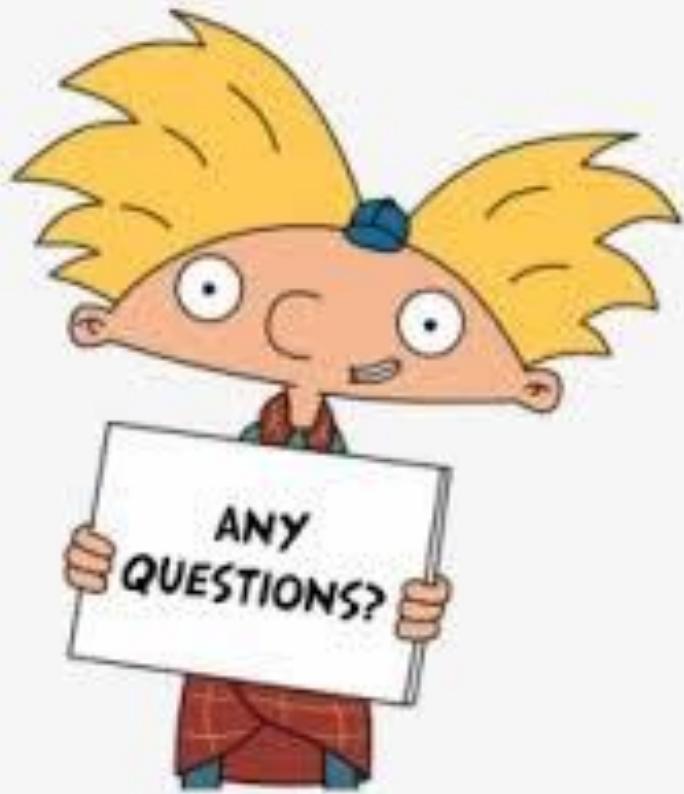
When  $k = 1$ , k-NNC is simply called the *nearest neighbor classifier (NNC)*.

- Asymptotically, it can be shown that the error of the NNC is less than twice the error of the Bayes Classifier.

- Cross-validation can be used to find the value of k.

# Problems with k-NNC

- Test time.
  - Scanning of the entire training set is required.
- Test space
  - Entire training set has to be stored.
- Both time and space complexities of k-NNC (assuming that the k value is a small constant):  $O(n)$ .
- This is quiet heavy. For Neural networks this is  $O(1)$ .



?????



# Bayesian Decision Theory

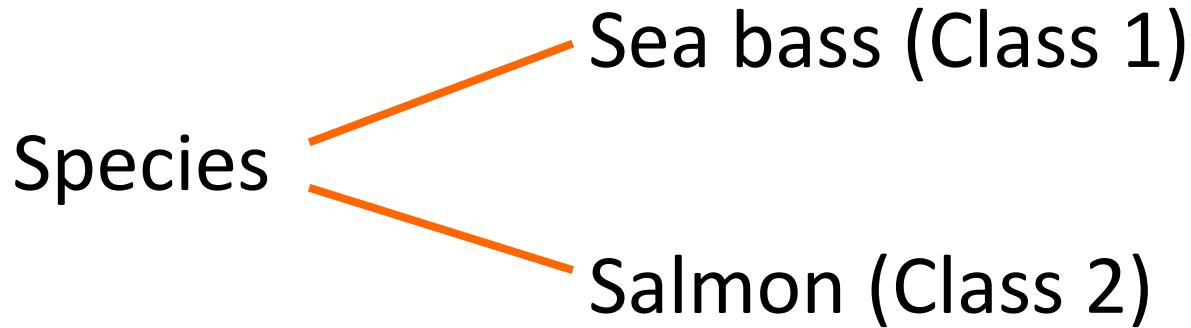
Primary source of reference: *Pattern Classification* – Duda  
and Hart

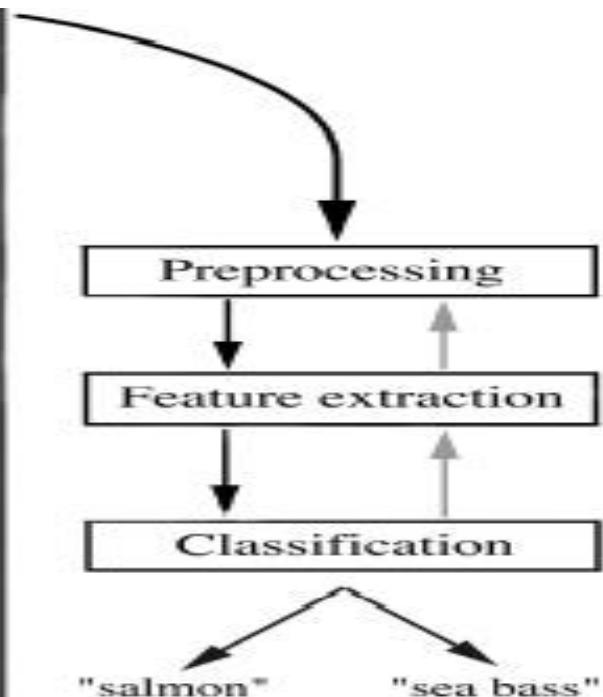
## Introduction

Bayesian Decision Theory–Continuous Features

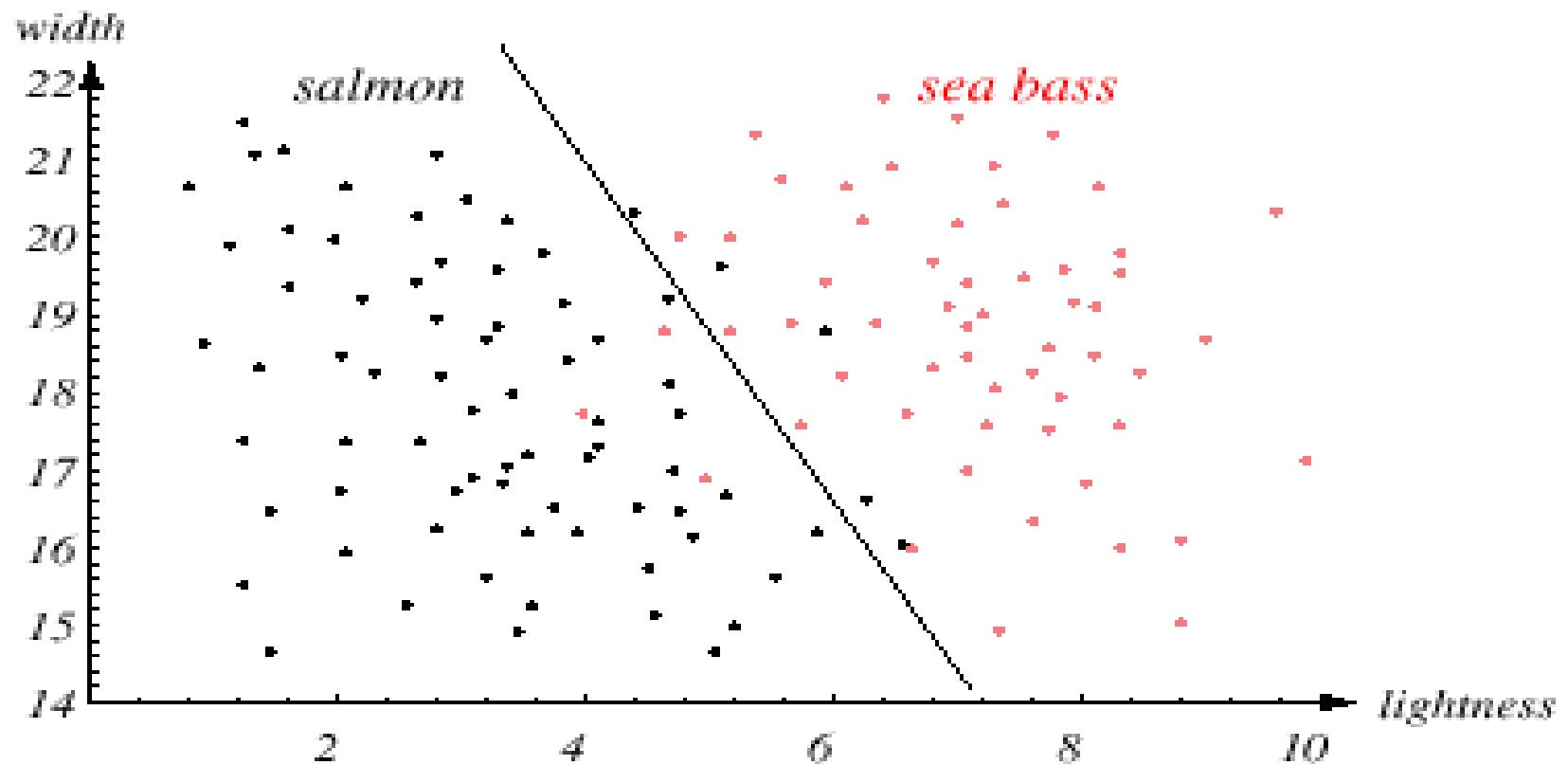
# An Example

- “Sorting incoming Fish on a conveyor according to species using optical sensing”





- Problem Analysis
  - Set up a camera and take some sample images to extract features like
    - Length of the fish
    - Lightness (based on the gray level)
    - Width of the fish



# Introduction

- The sea bass/salmon example  
(a two class problem)
- For example if we randomly catch 100 fishes and out of this if 75 are *sea bass* and 25 are *salmon*.
- Let the rule, in this case is: For any fish say its class is *sea bass*.
- What is the error rate of this rule?
- This information which is independent of feature values is called **apriori** knowledge.

- Let the two classes are  $\omega_1$  and  $\omega_2$ 
  - $P(\omega_1) + P(\omega_2) = 1$
  - State of nature (class) is a random variable
  - If  $P(\omega_1) = P(\omega_2)$ , we say it is of uniform priors
    - The catch of salmon and sea bass is equi-probable

- Decision rule with only the prior information
  - Decide  $\omega_1$  if  $P(\omega_1) > P(\omega_2)$ , otherwise decide  $\omega_2$
- *This is not a good classifier.*
- *We should take feature values into account !*
- *If  $x$  is the pattern we want to classify, then use the rule:*

*If  $P(\omega_1 | x) > P(\omega_2 | x)$  then assign class  $\omega_1$*

*Else assign class  $\omega_2$*

- *$P(\omega_1 | x)$  is called posterior probability of class  $\omega_1$  given that the pattern is  $x$ .*

# Bayes rule

- From data it might be possible for us to estimate  $p(x | \omega_i)$ , where  $i = 1$  or  $2$ . These are called class-conditional distributions.
- Also it is easy to find apriori probabilities  $P(\omega_1)$  and  $P(\omega_2)$ . How this can be done?
- Bayes rule combines apriori probability with class conditional distributions to find posteriori probabilities.

# Bayes Rule

$$P(B|A) = \frac{P(A, B)}{P(A)} = \frac{P(A|B) * P(B)}{P(A)}$$

This is Bayes Rule



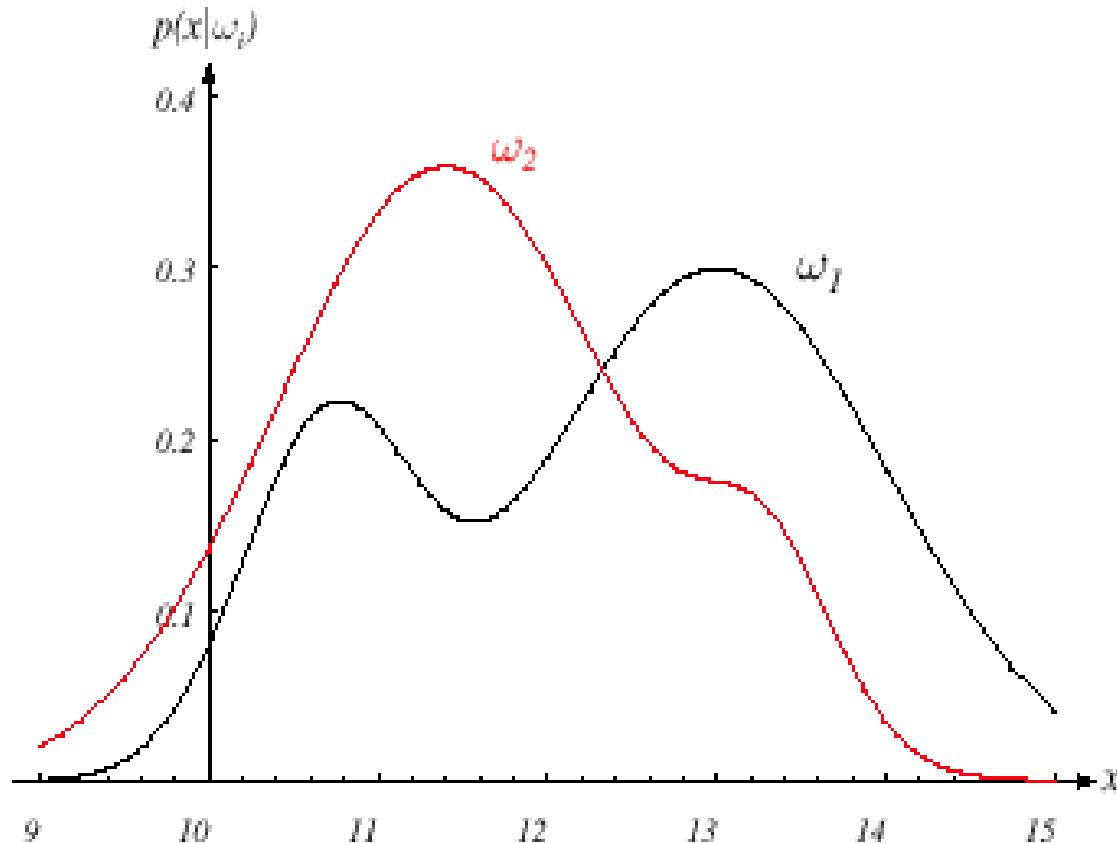
**Bayes, Thomas (1763)** An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London*, 53:370-418

$$P(\omega_j | x) = \frac{p(x | \omega_j) \cdot P(\omega_j)}{p(x)}$$

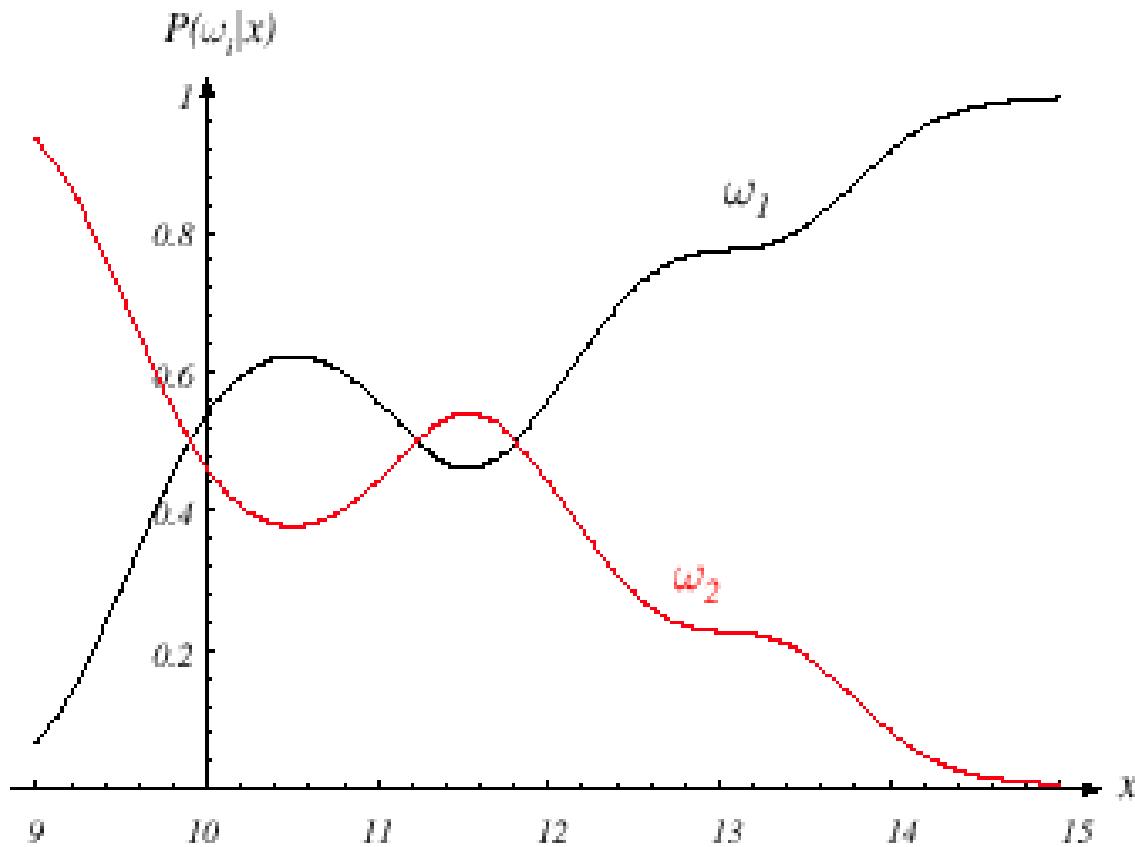
– Where in case of two categories

$$p(x) = \sum_{j=1}^{j=2} p(x | \omega_j) P(\omega_j)$$

$$\text{Posterior} = \frac{\text{Likelihood} \cdot \text{Prior}}{\text{Evidence}}$$



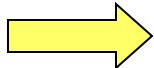
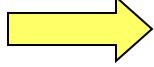
**FIGURE 2.1.** Hypothetical class-conditional probability density functions show the probability density of measuring a particular feature value  $x$  given the pattern is in category  $\omega_i$ . If  $x$  represents the lightness of a fish, the two curves might describe the difference in lightness of populations of two types of fish. Density functions are normalized, and thus the area under each curve is 1.0. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



**FIGURE 2.2.** Posterior probabilities for the particular priors  $P(\omega_1) = 2/3$  and  $P(\omega_2) = 1/3$  for the class-conditional probability densities shown in Fig. 2.1. Thus in this case, given that a pattern is measured to have feature value  $x = 14$ , the probability it is in category  $\omega_2$  is roughly 0.08, and that it is in  $\omega_1$  is 0.92. At every  $x$ , the posteriors sum to 1.0. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

- Decision given the posterior probabilities

X is an observation for which:

if  $P(\omega_1 | x) > P(\omega_2 | x)$   True state of nature =  $\omega_1$   
if  $P(\omega_1 | x) < P(\omega_2 | x)$   True state of nature =  $\omega_2$

Therefore:

whenever we observe a particular x, the probability of error is :

$$P(\text{error} | x) = P(\omega_1 | x) \text{ if we decide } \omega_2$$
$$P(\text{error} | x) = P(\omega_2 | x) \text{ if we decide } \omega_1$$

- Minimizing the probability of error
- Decide  $\omega_1$  if  $P(\omega_1 | x) > P(\omega_2 | x)$ ;  
otherwise decide  $\omega_2$

Therefore:

$$P(\text{error} | x) = \min [P(\omega_1 | x), P(\omega_2 | x)]$$

*(error of Bayes decision)*

# Average error rate

Average probability of error,  $P(\text{error})$  is :

$$\int P(\text{error} \mid x) p(x) dx$$

This is the expected value of  $P(\text{error}/x)$  w.r.t.  $x$ ,

i.e.,  $E_x[P(\text{error} / x)]$

1. Consider a one dimensional two class problem. The feature used is color of fish. Color can be either white or dark  $P(\omega_1) = 0.75$ ,  $P(\omega_2) = 0.25$ ,  $P(\text{white} | \omega_1) = 0.2$ ,  $P(\text{white} | \omega_2) = 0.6$ ,  $P(\text{dark} | \omega_1) = 0.8$ ,  $P(\text{dark} | \omega_2) = 0.4$  Find  $P(\text{error})$  of the Bayes Classifier.

1. Consider a one dimensional two class problem. The feature used is color of fish. Color can be either white or dark  $P(\omega_1) = 0.75$ ,  $P(\omega_2) = 0.25$ ,  $P(\text{white} | \omega_1) = 0.2$ ,  $P(\text{white} | \omega_2) = 0.6$ ,  $P(\text{dark} | \omega_1) = 0.8$ ,  $P(\text{dark} | \omega_2) = 0.4$  Find  $P(\text{error})$  of the Bayes Classifier.

$$P(\text{white}) = P(\text{white} | \omega_1)P(\omega_1) + P(\text{white} | \omega_2)P(\omega_2)$$

$$P(\text{white}) = 0.2 * 0.75 + 0.6 * 0.25 = 0.3$$

1. Consider a one dimensional two class problem. The feature used is color of fish. Color can be either white or dark  $P(\omega_1) = 0.75$ ,  $P(\omega_2) = 0.25$ ,  $P(\text{white} | \omega_1) = 0.2$ ,  $P(\text{white} | \omega_2) = 0.6$ ,  $P(\text{dark} | \omega_1) = 0.8$ ,  $P(\text{dark} | \omega_2) = 0.4$  Find  $P(\text{error})$  of the Bayes Classifier.

$$P(\text{white}) = P(\text{white} | \omega_1)P(\omega_1) + P(\text{white} | \omega_2)P(\omega_2)$$

$$P(\text{white}) = 0.2 * 0.75 + 0.6 * 0.25 = 0.3$$

$$P(\text{dark}) = P(\text{dark} | \omega_1)P(\omega_1) + P(\text{dark} | \omega_2)P(\omega_2)$$

$$P(\text{dark}) = 0.8 * 0.75 + 0.4 * 0.25 = 0.7$$

1. Consider a one dimensional two class problem. The feature used is color of fish. Color can be either white or dark  $P(\omega_1) = 0.75$ ,  $P(\omega_2) = 0.25$ ,  $P(\text{white} | \omega_1) = 0.2$ ,  $P(\text{white} | \omega_2) = 0.6$ ,  $P(\text{dark} | \omega_1) = 0.8$ ,  $P(\text{dark} | \omega_2) = 0.4$  Find  $P(\text{error})$  of the Bayes Classifier.

$$P(\text{white}) = P(\text{white} | \omega_1)P(\omega_1) + P(\text{white} | \omega_2)P(\omega_2)$$

$$P(\text{white}) = 0.2 * 0.75 + 0.6 * 0.25 = 0.3$$

$$P(\text{dark}) = P(\text{dark} | \omega_1)P(\omega_1) + P(\text{dark} | \omega_2)P(\omega_2)$$

$$P(\text{dark}) = 0.8 * 0.75 + 0.4 * 0.25 = 0.7$$

$$P(\omega_1 | \text{white}) = \frac{P(\text{white} | \omega_1)P(\omega_1)}{P(\text{white})} = \frac{0.2 * 0.75}{0.3} = 0.5$$

$$P(\omega_2 | \text{white}) = \frac{P(\text{white} | \omega_2)P(\omega_2)}{P(\text{white})} = \frac{0.6 * 0.25}{0.3} = 0.5$$

1. Consider a one dimensional two class problem. The feature used is color of fish. Color can be either white or dark  $P(\omega_1) = 0.75$ ,  $P(\omega_2) = 0.25$ ,  $P(\text{white} | \omega_1) = 0.2$ ,  $P(\text{white} | \omega_2) = 0.6$ ,  $P(\text{dark} | \omega_1) = 0.8$ ,  $P(\text{dark} | \omega_2) = 0.4$  Find  $P(\text{error})$  of the Bayes Classifier.

$$P(\text{white}) = P(\text{white} | \omega_1)P(\omega_1) + P(\text{white} | \omega_2)P(\omega_2)$$

$$P(\text{white}) = 0.2 * 0.75 + 0.6 * 0.25 = 0.3$$

$$P(\text{dark}) = P(\text{dark} | \omega_1)P(\omega_1) + P(\text{dark} | \omega_2)P(\omega_2)$$

$$P(\text{dark}) = 0.8 * 0.75 + 0.4 * 0.25 = 0.7$$

$$P(\omega_1 | \text{white}) = \frac{P(\text{white} | \omega_1)P(\omega_1)}{P(\text{white})} = \frac{0.2 * 0.75}{0.3} = 0.5$$

$$P(\omega_2 | \text{white}) = \frac{P(\text{white} | \omega_2)P(\omega_2)}{P(\text{white})} = \frac{0.6 * 0.25}{0.3} = 0.5$$

$$P(\omega_1 | \text{dark}) = \frac{P(\text{dark} | \omega_1)P(\omega_1)}{P(\text{dark})} = \frac{0.8 * 0.75}{0.7} = \frac{6}{7}$$

$$P(\omega_2 | \text{dark}) = \frac{P(\text{dark} | \omega_2)P(\omega_2)}{P(\text{dark})} = \frac{0.4 * 0.25}{0.7} = \frac{1}{7}$$

$$P(error) = P(error|white)P(white) + P(error|dark)P(dark)$$

$$P(error) = 0.5 * 0.3 + \frac{1}{7} * 0.7 = 0.25$$

$$P(error) = P(error|white)P(white) + P(error|dark)P(dark)$$

$$P(error) = 0.5 * 0.3 + \frac{1}{7} * 0.7 = 0.25$$

- But, what is the error, if we use only apriori probabilities?

Since,  $P(\omega_1) = 0.75$ ,  $P(\omega_2) = 0.25$ , every pattern is assigned to  $\omega_1$ , So the error,

$$P(error) = P(\omega_2 | white)P(white) + P(\omega_2 | dark)P(dark)$$

Since,  $P(\omega_1) = 0.75$ ,  $P(\omega_2) = 0.25$ , every pattern is assigned to  $\omega_1$ , So the error,

$$P(error) = P(\omega_2 | white)P(white) + P(\omega_2 | dark)P(dark)$$

$$P(error) = \frac{P(white|\omega_2)P(\omega_2)}{P(white)}P(white) + \frac{P(dark|\omega_2)P(\omega_2)}{P(dark)}P(dark)$$

Since,  $P(\omega_1) = 0.75$ ,  $P(\omega_2) = 0.25$ , every pattern is assigned to  $\omega_1$ , So the error,

$$P(error) = P(\omega_2 | white)P(white) + P(\omega_2 | dark)P(dark)$$

$$P(error) = \frac{P(white|\omega_2)P(\omega_2)}{P(white)}P(white) + \frac{P(dark|\omega_2)P(\omega_2)}{P(dark)}P(dark)$$

$$P(error) = (P(white|\omega_2) + P(dark|\omega_2))P(\omega_2)$$

$$P(error) = P(\omega_2) = 0.25$$

Since,  $P(\omega_1) = 0.75$ ,  $P(\omega_2) = 0.25$ , every pattern is assigned to  $\omega_1$ , So the error,

$$P(error) = P(\omega_2 | white)P(white) + P(\omega_2 | dark)P(dark)$$

$$P(error) = \frac{P(white|\omega_2)P(\omega_2)}{P(white)}P(white) + \frac{P(dark|\omega_2)P(\omega_2)}{P(dark)}P(dark)$$

$$P(error) = (P(white|\omega_2) + P(dark|\omega_2))P(\omega_2)$$

$$P(error) = P(\omega_2) = 0.25$$

- Same error? Where is the advantage?!

Consider  $P(\omega_1) = 0.5$ ,  $P(\omega_2) = 0.5$

$$P(\text{white}) = P(\text{white}|\omega_1)P(\omega_1) + P(\text{white}|\omega_2)P(\omega_2)$$

$$P(\text{white}) = 0.2 * 0.5 + 0.6 * 0.5 = 0.4$$

$$P(\text{dark}) = P(\text{dark}|\omega_1)P(\omega_1) + P(\text{dark}|\omega_2)P(\omega_2)$$

$$P(\text{dark}) = 0.8 * 0.5 + 0.4 * 0.5 = 0.6$$

$$P(\omega_1|\text{white}) = \frac{P(\text{white}|\omega_1)P(\omega_1)}{P(\text{white})} = \frac{0.2 * 0.5}{0.4} = 0.25$$

$$P(\omega_2|\text{white}) = \frac{P(\text{white}|\omega_2)P(\omega_2)}{P(\text{white})} = \frac{0.6 * 0.5}{0.4} = 0.75$$

$$P(\omega_1|\text{dark}) = \frac{P(\text{dark}|\omega_1)P(\omega_1)}{P(\text{dark})} = \frac{0.8 * 0.5}{0.6} = \frac{2}{3}$$

$$P(\omega_2|\text{dark}) = \frac{P(\text{dark}|\omega_2)P(\omega_2)}{P(\text{dark})} = \frac{0.4 * 0.5}{0.6} = \frac{1}{3}$$

$$P(\text{error}) = P(\text{error}|\text{white})P(\text{white}) + P(\text{error}|\text{dark})P(\text{dark})$$

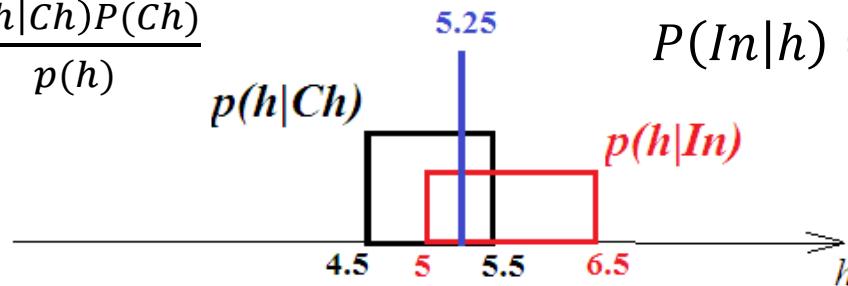
$$P(\text{error}) = 0.25 * 0.4 + \frac{1}{3} * 0.6 = 0.3$$

- But,  $P(\text{error})$  based on apriori probabilities only is 0.5.
- Error based on the Bayes classifier is the lower bound.
  - Any classifier's error is greater than or equal to this.
- One can prove this!

2. We know that 30% of people in Sri City are from China. Remaining are Indians. We know that Chinese height is a uniformly distributed random variable with parameters 4.5 and 5.5. We also know that Indians height is also uniformly distributed with parameters 5 and 6.5. By measuring his/her height we want to classify the person as "Chinese" or "Indian" (We know that the person is living in Sri City). We follow the rule "if height is 5.25 or below classify the person as Chinese, otherwise classify the person as Indian". There are two types of mistakes in this classification, (1) Chinese being classified as Indians, (2) Indians being classified as Chinese. Find the probability of making each of these two types of mistakes.

$$P(Ch) = 0.3$$

$$P(Ch|h) = \frac{p(h|Ch)P(Ch)}{p(h)}$$



$$P(In) = 0.7$$

$$P(In|h) = \frac{p(h|In)P(In)}{p(h)}$$

$$p(h) = p(h|Ch)P(Ch) + p(h|In)P(In)$$

$$= \begin{cases} 0.3, & \text{for } h \text{ in } [4.5, 5] \\ 1 * 0.3 + \frac{2}{3} * 0.7 = 0.767, & \text{for } h \text{ in } [5, 5.5] \\ \frac{2}{3} * 0.7 = 0.467, & \text{for } h \text{ in } [5.5, 6.5] \end{cases}$$

For  $h \in [4.5, 5]$

$$P(Ch | h) = (1 * 0.3) / 0.3 = 1$$

$$P(In|h) = 0$$

For  $h \in [5, 5.5]$

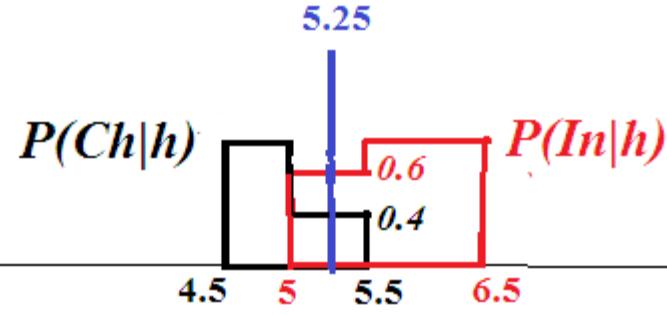
$$P(Ch | h) = (1 * 0.3) / 0.767 = 0.39$$

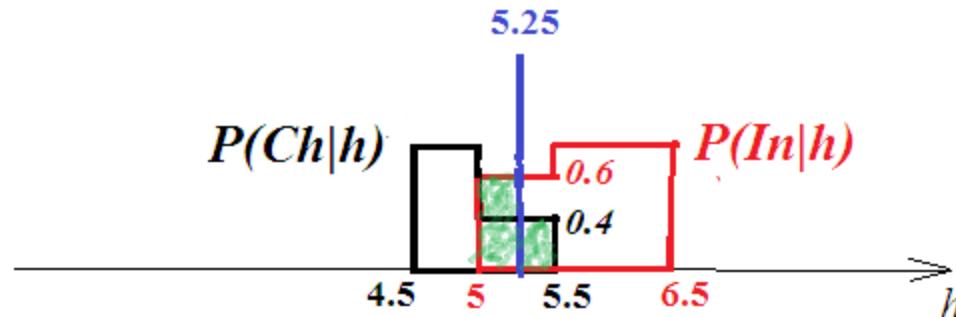
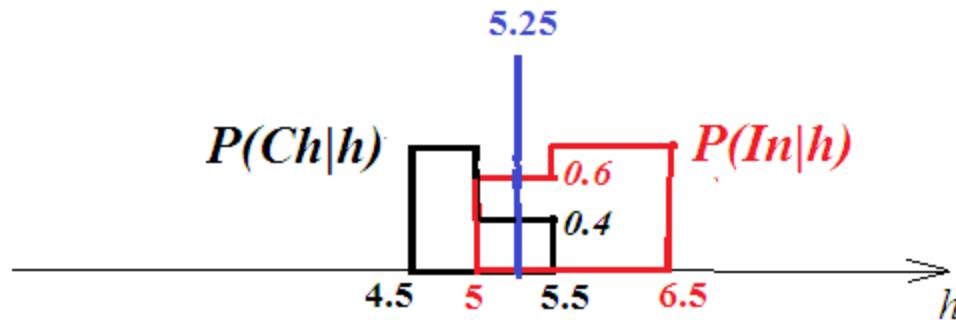
$$P(In | h) = 0.67 * 0.7 / 0.767 = 0.61$$

For  $h \in [5, 6.5]$

$$P(Ch | h) = 0$$

$$P(In|h) = 1$$





*Error is due to this region*

$$\begin{aligned}
 P(\text{error}) &= \int_5^{5.25} 0.609 p(h) dh + \int_{5.25}^{5.5} 0.391 p(h) dh \\
 &= 0.117 + 0.075 = 0.192
 \end{aligned}$$

- But, this is not the Bayes classifier.
- What does the Bayes classifier do?
- What is the error of the Bayes classifier?

# Bayesian Decision Theory – Continuous Features

- Generalization of the preceding ideas
  - Use of more than one feature
  - Use more than two states (classes) of nature
  - Allowing actions (decisions) other than just classification.
  - Introduce a *loss function* which is more general than the probability of error.

- Allowing actions other than classification primarily allows the possibility of rejection
- Refusing to make a decision in close or bad cases!
- The loss function states how costly each action taken is

# Problem setting

What is given to us:

Loss function :  $\lambda(\alpha_i / \omega_j)$  is the loss of taking action  $\alpha_i$  when the state of nature is  $\omega_j$

Let  $\Omega = \{\omega_1, \omega_2, \dots, \omega_c\}$  be the set of c states of nature  
(or “categories” or “classes”)

Let  $A = \{\alpha_1, \alpha_2, \dots, \alpha_a\}$  be the set of possible actions  
For example :  $\alpha_1$  is the action *ring the alarm*  
 $\alpha_2$  is the action *shutdown the system*

Let  $\Omega = \{\omega_1, \omega_2, \dots, \omega_c\}$  be the set of c states of nature  
(or “categories” or “classes”)

Let  $A = \{\alpha_1, \alpha_2, \dots, \alpha_a\}$  be the set of possible actions  
For example :  $\alpha_1$  is the action *ring the alarm*  
 $\alpha_2$  is the action *shutdown the system*

**Objective:** Given a pattern  $x$ , find the action to take.  
That is, to find a function  $\alpha(x)$  which maps  $x$  to action.

# How to find the best action?

Let  $R(\alpha_i | x)$  is the risk of taking action  $\alpha_i$  when the given pattern is  $x$  (this is called conditional risk).

We can take the action  $\alpha_k$  provided  $R(\alpha_k | x)$  is minimum in  $\{ R(\alpha_1 | x), R(\alpha_2 | x), \dots, R(\alpha_a | x) \}$

How to relate  $R(\alpha_i | x)$  with  $\lambda(\alpha_i | \omega_j)$  values.

$$R(\alpha_i | x) = \sum_{j=1}^{j=c} \lambda(\alpha_i | \omega_j) P(\omega_j | x)$$

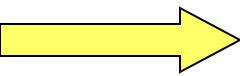
# $\alpha(x)$ is given, how good is this?

- Overall risk of this rule can be found

$$R = \int R(\alpha(x)|x)p(x) dx,$$

- The smaller this quantity, better the decision rule.

Select the action  $\alpha_i$  for which  $R(\alpha_i | x)$  is minimum



R is minimum and R in this case is called the  
*Bayes risk = best performance that can be achieved!*

- Two-category classification

$\alpha_1$  : deciding  $\omega_1$

$\alpha_2$  : deciding  $\omega_2$

$\lambda_{ij} = \lambda(\alpha_i \mid \omega_j)$

loss incurred for deciding  $\omega_i$  when the true state of nature is  $\omega_j$

Conditional risk:

$$R(\alpha_1 \mid x) = \lambda_{11}P(\omega_1 \mid x) + \lambda_{12}P(\omega_2 \mid x)$$

$$R(\alpha_2 \mid x) = \lambda_{21}P(\omega_1 \mid x) + \lambda_{22}P(\omega_2 \mid x)$$

Our rule is the following:

if  $R(\alpha_1 | x) < R(\alpha_2 | x)$

action  $\alpha_1$ : “decide  $\omega_1$ ” is taken

This results in the equivalent rule :

decide  $\omega_1$  if:

$$(\lambda_{21} - \lambda_{11}) p(x | \omega_1) P(\omega_1) > (\lambda_{12} - \lambda_{22}) p(x | \omega_2) P(\omega_2)$$

and decide  $\omega_2$  otherwise

# An Example

- Let the two actions are:

$$\begin{array}{ll} \alpha_1 & \xrightarrow{\text{yellow arrow}} x \text{ is criminal} \\ \alpha_2 & \xrightarrow{\text{yellow arrow}} x \text{ is innocent} \end{array}$$

Let  $\lambda(\alpha_1 | \omega_1) = 0$ ;  $\lambda(\alpha_1 | \omega_2) = 10$ ;  
 $\lambda(\alpha_2 | \omega_1) = 1$ ;  $\lambda(\alpha_2 | \omega_2) = 0$ ;

Assume equal priors, and

Let  $p(x | \omega_1) = 0.8$ ,  $p(x | \omega_2) = 0.6$

What action you will take?

$\omega_1$  is class of criminals,  
 $\omega_2$  is class of innocents.

# Example: Contd...

$$p(x | \omega_j) \cdot P(\omega_j)$$

$$\bullet \quad P(\omega_j | x) = \frac{p(x | \omega_j) \cdot P(\omega_j)}{p(x)}$$

$$0.8(0.5)$$

$$\bullet \quad P(\omega_1 | x) = \frac{0.8(0.5)}{0.8(0.5) + 0.6(0.5)} = 8/14 = 0.57$$

$$\bullet \quad P(\omega_2 | x) = 0.43$$

*• As per plain Bayes rule we declare the person to be a criminal.*

# Example: Contd ...

- $$\begin{aligned} R(\alpha_1 | x) &= \lambda(\alpha_1 | \omega_1) P(\omega_1 | x) + \lambda(\alpha_1 | \omega_2) P(\omega_2 | x) \\ &= 0 (0.57) + 10 (0.43) \\ &= 4.3 \end{aligned}$$
- $$\begin{aligned} R(\alpha_2 | x) &= \lambda(\alpha_2 | \omega_1) P(\omega_1 | x) + \lambda(\alpha_2 | \omega_2) P(\omega_2 | x) \\ &= 1 (0.57) + 0 (0.43) \\ &= 0.57 \end{aligned}$$

Action taken:  $x$  is innocent

## Likelihood ratio:

The preceding rule is equivalent to the following rule:

$$\text{if } \frac{p(x | \omega_1)}{p(x | \omega_2)} > \frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \cdot \frac{P(\omega_2)}{P(\omega_1)}$$

Then take action  $\alpha_1$  (decide  $\omega_1$ )

Otherwise take action  $\alpha_2$  (decide  $\omega_2$ )

# Example 1

- Let  $\Omega = \{\omega_1, \omega_2, \omega_3\}$ ,  $A = \{\alpha_1, \alpha_2, \alpha_3\}$   
 $P(\omega_1 | x) = 0.1$ ,  $P(\omega_2 | x) = 0.4$ ,  $P(\omega_3 | x) = 0.5$

Loss function is

Loss	$\omega_1$	$\omega_2$	$\omega_3$
$\alpha_1$	0	1	2
$\alpha_2$	1	0	2
$\alpha_3$	3	10	0

Find  $R(\alpha_k | x)$  for  $k = 1, 2, 3$ .

Find what is the best action?

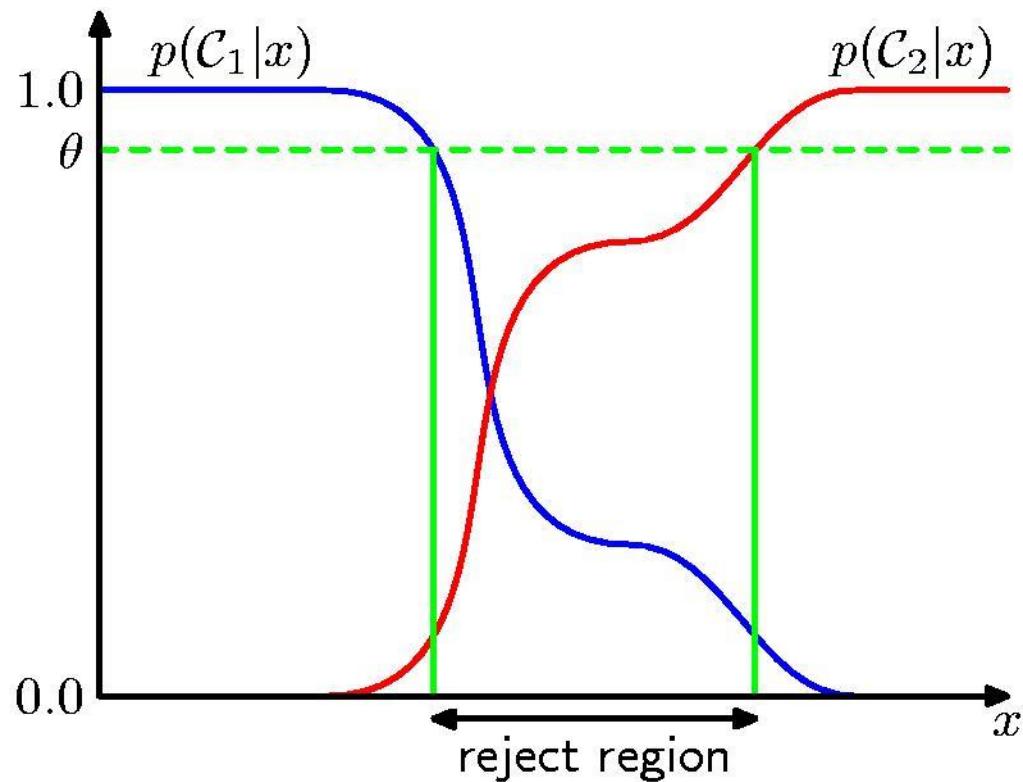
## Reject option (I do not want to classify)

Actions are : Assign a class label or reject

Sometimes, when misclassification is costly, we can reject to classify it.

May be some other expert can look into it and can take appropriate action.

# Reject Option



- Read Duda and Hart book and try to solve some problems related to this.



# Naïve Bayes Classification

# Introduction

- The Bayes Classifier requires probability structure of the problem to be known.
- Density estimation (using non-parametric or parametric methods) is one way to handle the problem.
- There are several problems ....

# Problems with density estimation

- Large datasets are needed.
- Numeric valued features are required.
- In practice these two may not be satisfied.

# How to overcome the problem

- One has to work with the given data set.
- So, Probability estimations needs to be done using the given data only.
- Often marginal probabilities can be better estimated than the joint probabilities.
- Also, marginal probabilities are easy to compute.

## Play-tennis data

Outlook	Temperature	Humidity	Windy	Class
sunny	hot	high	false	N
sunny	hot	high	true	N
overcast	hot	high	false	P
rain	mild	high	false	P
rain	cool	normal	false	P
rain	cool	normal	true	N
overcast	cool	normal	true	P
sunny	mild	high	false	N
sunny	cool	normal	false	P
rain	mild	normal	false	P
sunny	mild	normal	true	P
overcast	mild	high	true	P
overcast	hot	normal	false	P
rain	mild	high	true	N

- $P(\langle \text{sunny}, \text{cool}, \text{high}, \text{false} \rangle | N) = 0$
- But,  $P(\text{sunny} | N) = 3/5$ ,  $P(\text{cool} | N) = 1/5$ ,  
 $P(\text{high} | N) = 4/5$ ,  $P(\text{false} | N) = 2/5$ .

- $P(\langle \text{sunny}, \text{cool}, \text{high}, \text{false} \rangle | N) = 0$
- This may be because of the smaller dataset.
- If we increase the dataset size, this may become a positive number.
- This problem is often referred to as “the curse of dimensionality”.

# Assumption

- Make the assumption that for a given class, features are independent of each other.
- In practice, this assumption holds very often.
- Then  $P(<\text{sunny}, \text{cool}, \text{high}, \text{false}> | N) = P(\text{sunny} | N) \cdot P(\text{cool} | N) \cdot P(\text{high} | N) \cdot P(\text{false} | N) = 3/5 \cdot 1/5 \cdot 4/5 \cdot 2/5 = 24/625.$

# Naïve Bayesian Classification

- Naïve assumption: for a given class, features are independent of each other

$$P(x_1, \dots, x_k | C) = P(x_1 | C) \cdot \dots \cdot P(x_k | C)$$

- $P(x_i | C)$  is estimated as the relative freq of samples having value  $x_i$  as i-th attribute in class C
- It often makes the problem a feasible and easy one to solve.

# Play-tennis example: estimating $P(x_i | C)$

Outlook	Temperature	Humidity	Windy	Class
sunny	hot	high	false	N
sunny	hot	high	true	N
overcast	hot	high	false	P
rain	mild	high	false	P
rain	cool	normal	false	P
rain	cool	normal	true	N
overcast	cool	normal	true	P
sunny	mild	high	false	N
sunny	cool	normal	false	P
rain	mild	normal	false	P
sunny	mild	normal	true	P
overcast	mild	high	true	P
overcast	hot	normal	false	P
rain	mild	high	true	N

$$P(p) = 9/14$$

$$P(n) = 5/14$$

outlook	
$P(\text{sunny}   p) = 2/9$	$P(\text{sunny}   n) = 3/5$
$P(\text{overcast}   p) = 4/9$	$P(\text{overcast}   n) = 0$
$P(\text{rain}   p) = 3/9$	$P(\text{rain}   n) = 2/5$
temperature	
$P(\text{hot}   p) = 2/9$	$P(\text{hot}   n) = 2/5$
$P(\text{mild}   p) = 4/9$	$P(\text{mild}   n) = 2/5$
$P(\text{cool}   p) = 3/9$	$P(\text{cool}   n) = 1/5$
humidity	
$P(\text{high}   p) = 3/9$	$P(\text{high}   n) = 4/5$
$P(\text{normal}   p) = 6/9$	$P(\text{normal}   n) = 2/5$
windy	
$P(\text{true}   p) = 3/9$	$P(\text{true}   n) = 3/5$
$P(\text{false}   p) = 6/9$	$P(\text{false}   n) = 2/5$

# Play-tennis example: classifying X

- An unseen sample  $X = \langle \text{rain}, \text{hot}, \text{high}, \text{false} \rangle$
- $P(X|p) \cdot P(p) =$   
 $P(\text{rain}|p) \cdot P(\text{hot}|p) \cdot P(\text{high}|p) \cdot P(\text{false}|p) \cdot P(p) = 3/9 \cdot 2/9 \cdot 3/9 \cdot 6/9 \cdot 9/14 = 0.010582$
- $P(X|n) \cdot P(n) =$   
 $P(\text{rain}|n) \cdot P(\text{hot}|n) \cdot P(\text{high}|n) \cdot P(\text{false}|n) \cdot P(n) = 2/5 \cdot 2/5 \cdot 4/5 \cdot 2/5 \cdot 5/14 = 0.018286$
- Sample **X is** classified in class **n** (don't play)

# With Continuous features

- In order to use the Naïve Bayes classifier, the features has to be discretized appropriately (**otherwise what happens?**)

# With Continuous features

- In order to use the Naïve Bayes classifier, the features has to be discretized appropriately (**otherwise what happens?**)
- Height = 4.234 will not occur anywhere in that column; but 4.213, 4.285 may be occurring. If you discretize (eg., rounding) then frequency ratio's are meaningful.
- Clustering of feature values of a feature may be done to achieve a better discretization.