

Artificial Intelligence for Knowledge Discovery

Hatim Chahdi

Lecture 3 : Machine Learning

Course objectives

- The main objective of this course is to give you an overview of the process **of building a end-to-end Knowledge Discovery System.**
- Through the lectures, the labs and the project, you will gain an understanding and a hands-on experience in the different parts and the necessary components to set up such systems
- **Ontologies – Reasoning - Text Mining - Knowledge extraction et representation**

Course overview

- Introduction
- Ontologies for knowledge management
- Deductive Reasoning for Ontologies
- Applied Machine Learning to Text
- Knowledge Extraction for Ontology Construction

Project : Text mining for Knowledge Extraction

Course logistics

Course planning :

- Lectures : 10,5h
- Tutorials : 14h
- Project : 10,5h

Prerequisites :

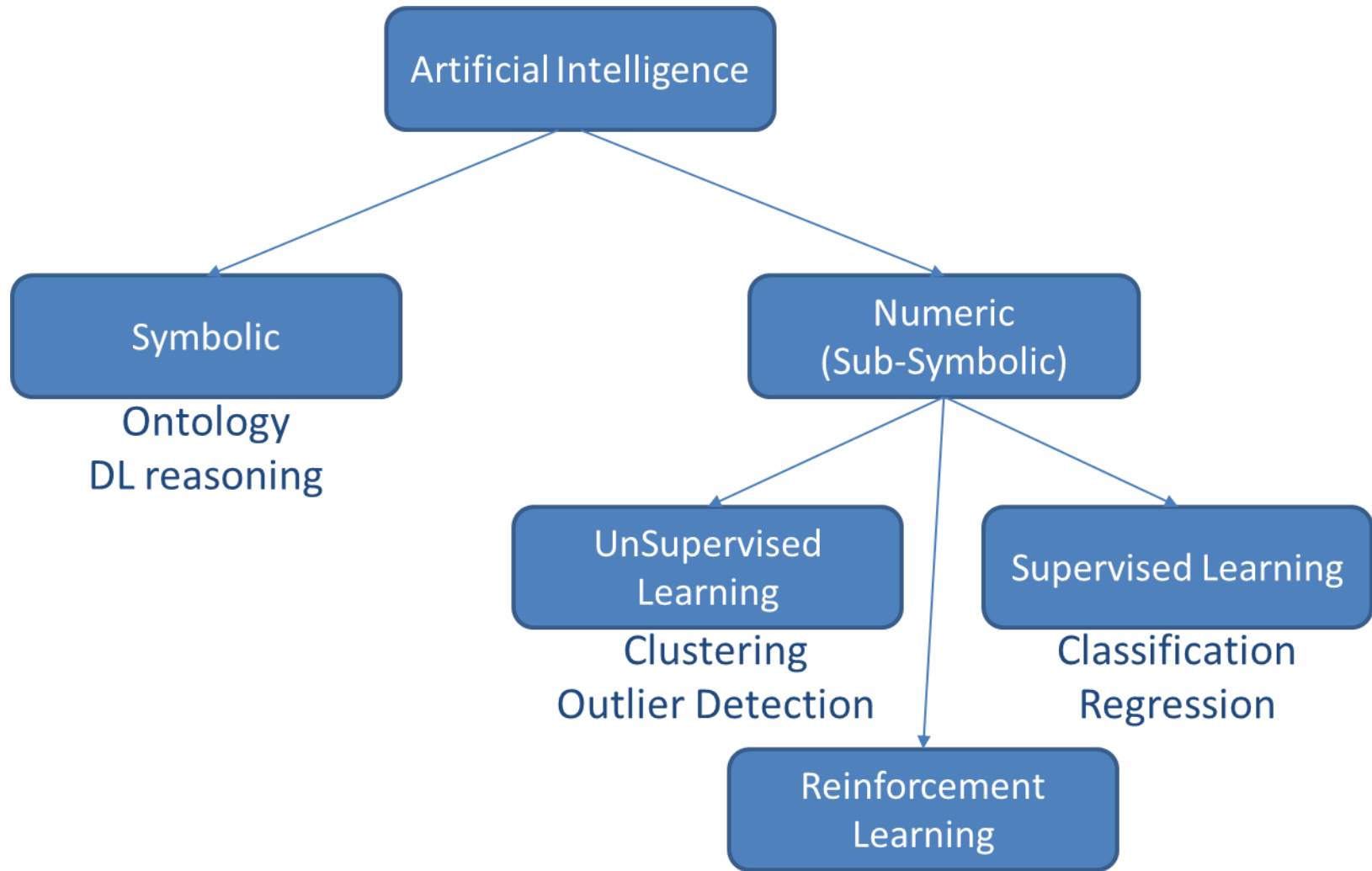
- Intermediate Java & Python
- Basic Statistics

Lecture Outline:

- Overview of the knowledge extraction process
- Machine Learning techniques and concepts
- Overview of the Knowledge discovery process
- Conclusion and discussions about the project

- Recall of the course objective : Artificial Intelligence for Knowledge Discovery
 - Ontology => Knowledge Representation and Management
 - Machine Learning => Knowledge Extraction
 - Decision Support System => Using both to deliver insights

An Artificial Intelligence Vision



Supervised Learning :

- Classification
- Regression
- A lot of other variants :
 - Sequence generation
 - Syntax tree prediction
 - Object detection
 - Image segmentation

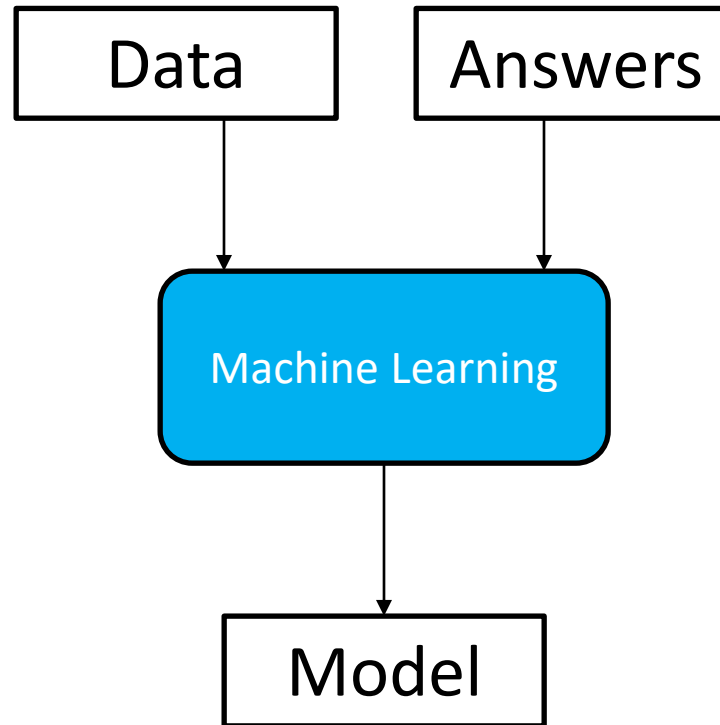
UnSupervised Learning :

- Clustering
- Dimensionality reduction

What will we see in the rest of this lecture :

- Classification algorithms
- Clustering algorithms
- The general framework of the machine learning process

Classification :



- The principle

Given a collection of a training set containing a set of attributes (variables) and the corresponding output (class label), learn a representation (rules, model) that maps between the attributes and there labels

- Formal 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 ”

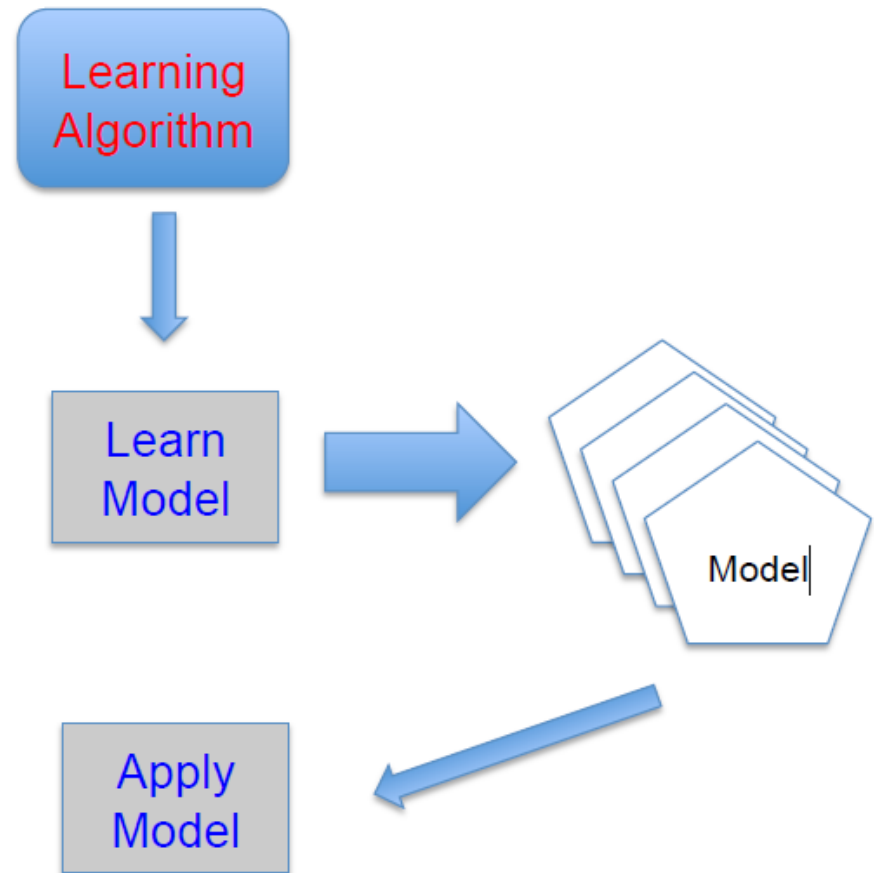
Tom M. Mitchell

Supervised Learning :

- Classification : The framework

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Tid	Refund	Marital Status	Taxable Income	Cheat
1	No	Married	80K	?
2	No	Single	100K	?
3	Yes	Single	90K	?
4	No	Married	120K	?
5	Yes	Divorced	130K	?



Supervised Learning :

- Classification examples :
 - Filtering spam from emails
 - Predicting tumor cells as benign or malignant
 - Classifying credit card transactions as legitimate or fraudulent
 - Classifying secondary structures of protein as alpha-helix, beta-sheet, or random coil
 - Categorizing news stories as finance (Topics)
 - Movies Recommendation
 - Named Entity Recognition

Supervised Learning :

- Classification techniques :
 - Decision Tree
 - Naïve Bayes
 - Instance Based Learning
 - Rule-based Methods
 - Neural Networks
 - Bayesian Belief Networks
 - Support Vector Machines

Supervised Learning :

Before we continue : Terminology (1/2)

- Sample or input : One data point that goes into the model
- Prediction or output : What comes out of your model
- Target – The truth : What your model should ideally have predicted
- Prediction error or loss value : measure of the distance between your model's prediction and the target
- Classes : A set of possible labels to choose from in a classification problem

Before we continue : Terminology (2/2)

- Label : A specific instance of a class annotation in a classification problem
- Ground-truth or annotations : All targets for a dataset, typically collected by humans
- Binary classification : A classification task where each input sample should be categorized into two exclusive categories
- Multiclass classification : A classification task where each input sample should be categorized into more than two categories
- Multilabel classification : A classification task where each input sample can be assigned multiple labels.

Supervised Learning :

Google ML Glossary :

<https://developers.google.com/machine-learning/glossary/>

Decision Tree : (Principles)

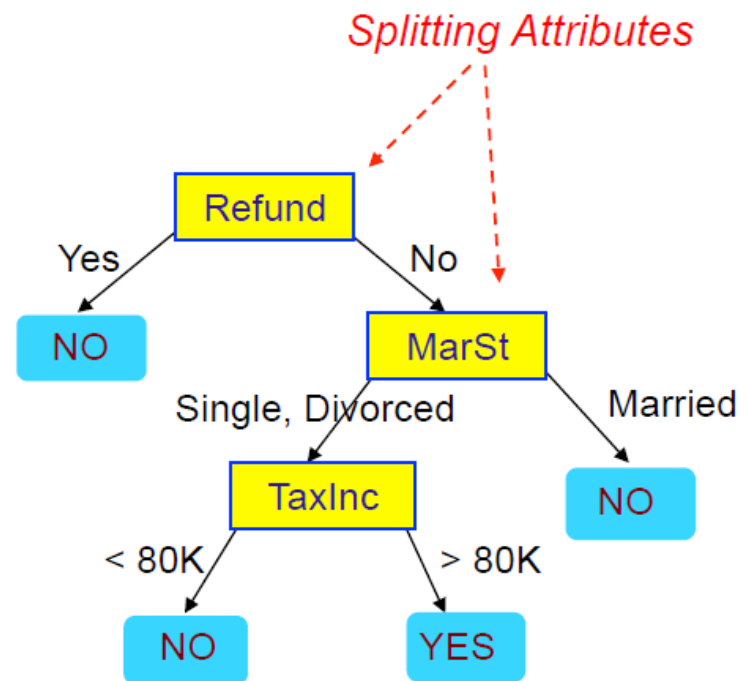
- Uses a tree structure to model the training set
- Classifies a new record following the path in the tree
- Inner nodes represent attributes and leaves nodes represent the class

Supervised Learning :

Decision Tree : Example

categorical
categorical
continuous
class

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



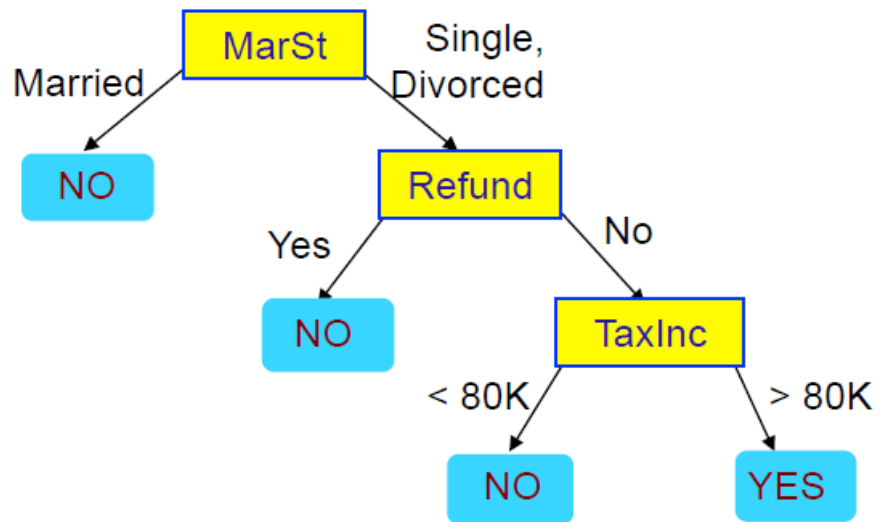
Training Data

Model: Decision Tree

Supervised Learning :

Decision Tree : Example of an other tree

<i>Tid</i>	<i>Refund</i>	<i>Marital Status</i>	<i>Taxable Income</i>	<i>Cheat</i>
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

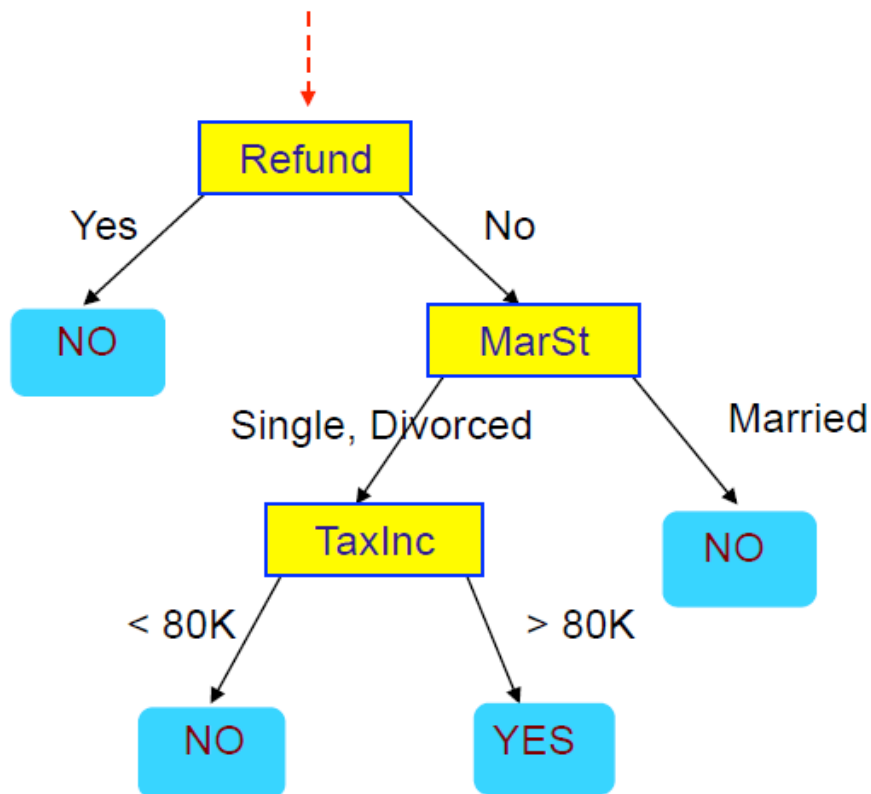


There could be more than one tree that fits the same data!

Supervised Learning :

Decision Tree : Model application

Start from the root of tree.



Test Data

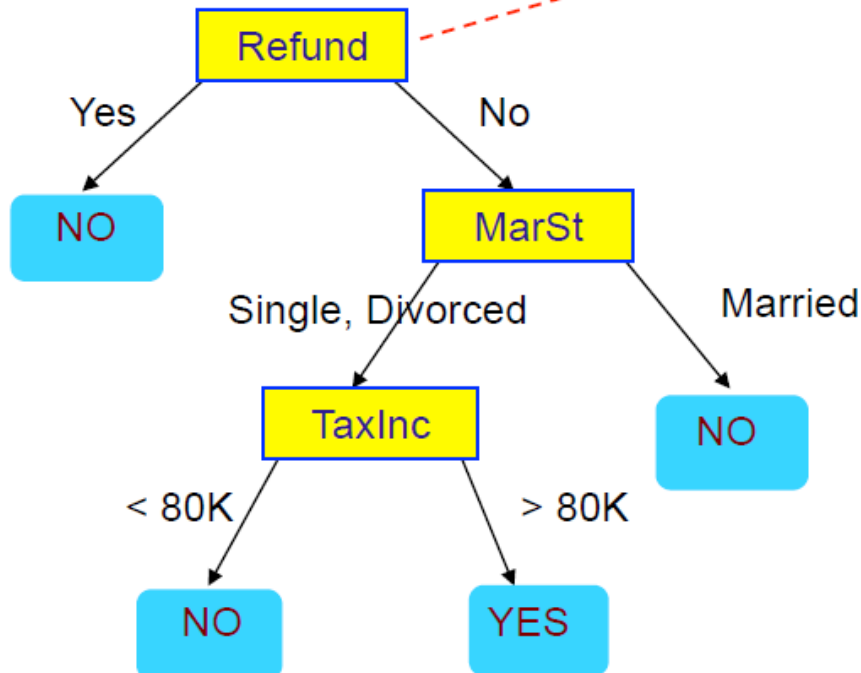
Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

Supervised Learning :

Decision Tree : Model application

Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

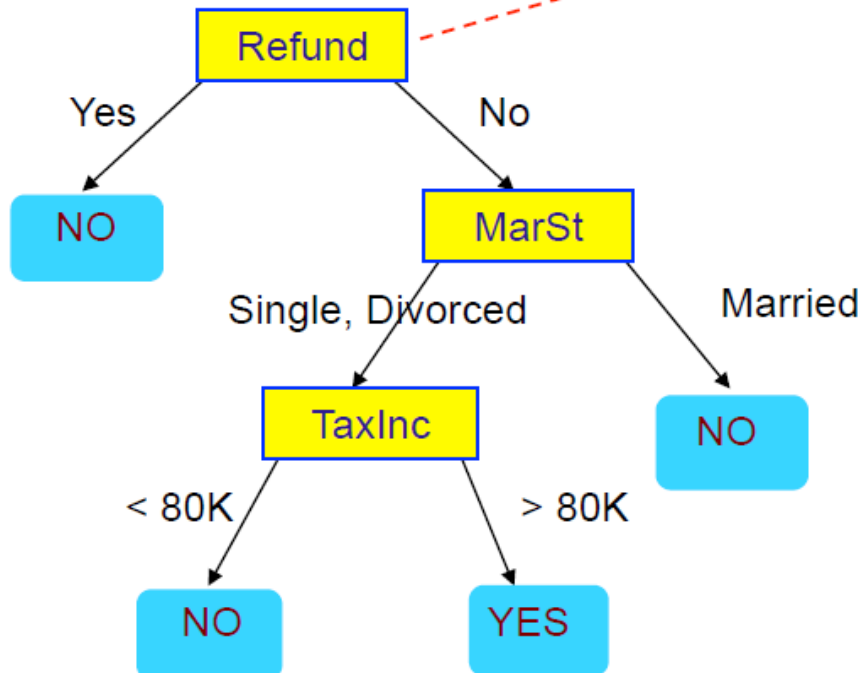


Supervised Learning :

Decision Tree : Model application

Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

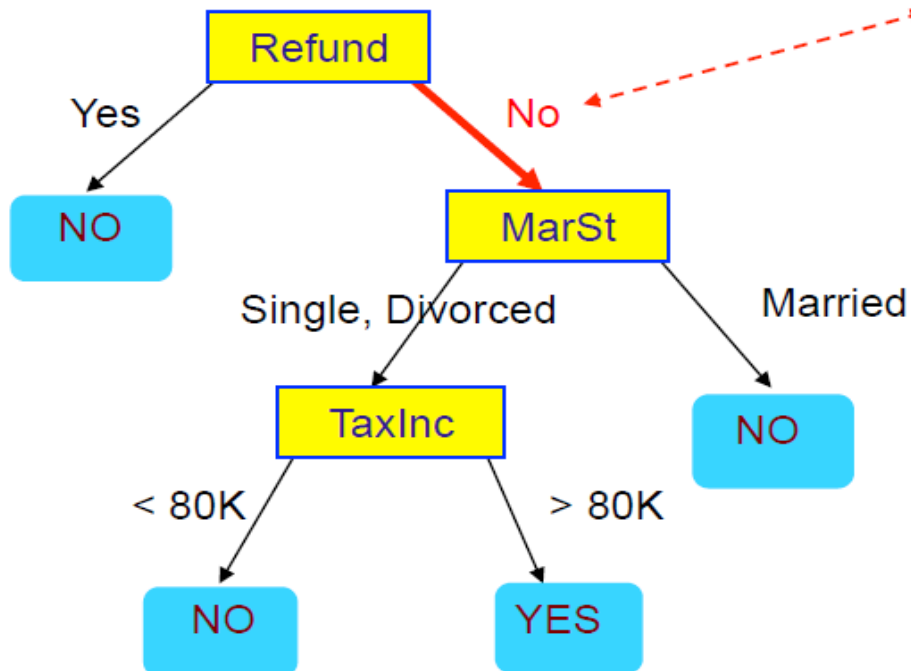


Supervised Learning :

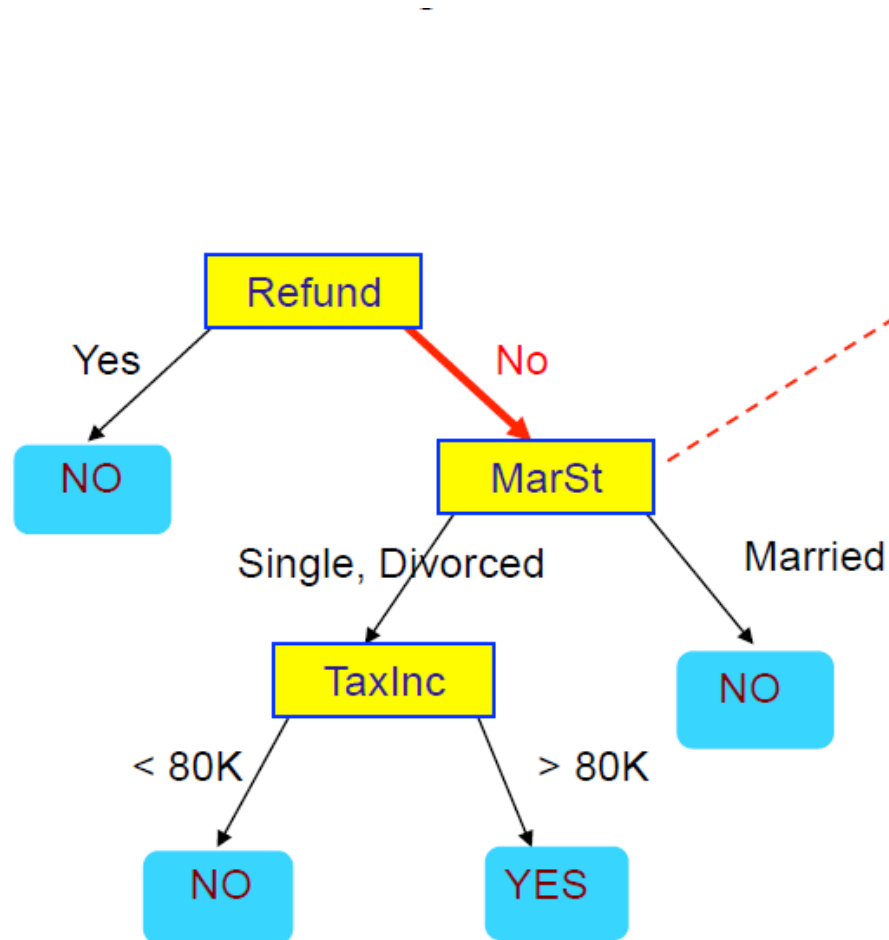
Decision Tree : Model application

Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



Decision Tree : Model application

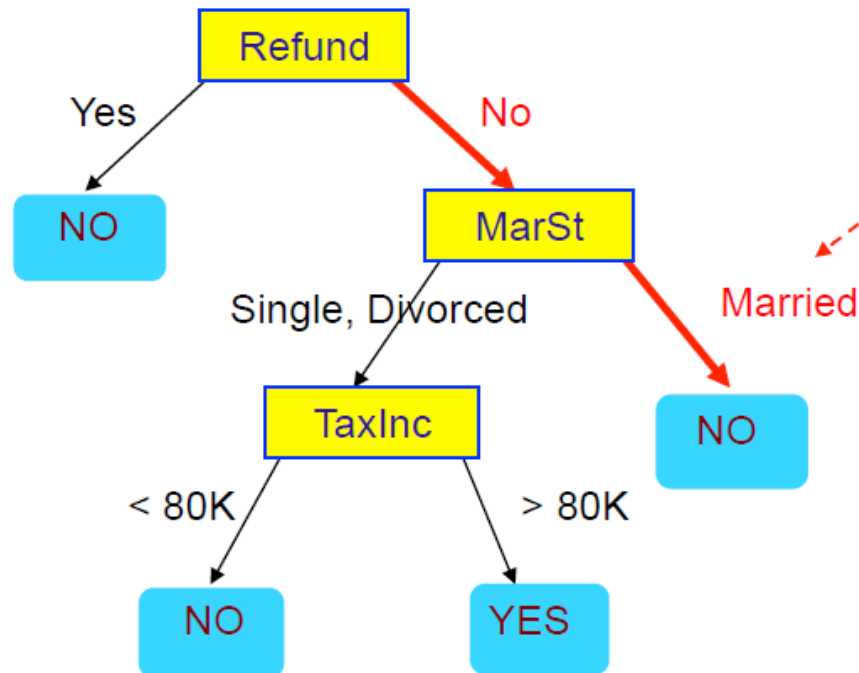


Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

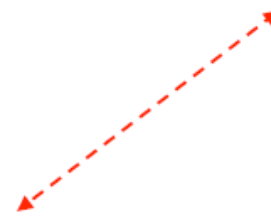
Supervised Learning :

Decision Tree : Model application



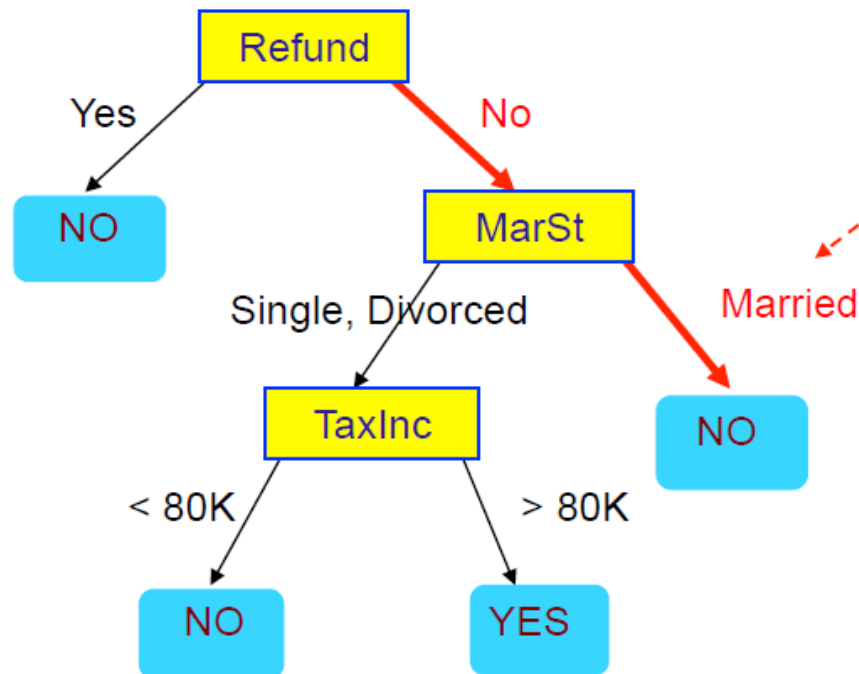
Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



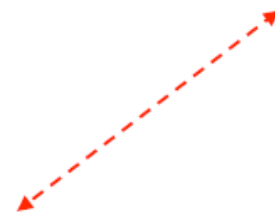
Supervised Learning :

Decision Tree : Model application



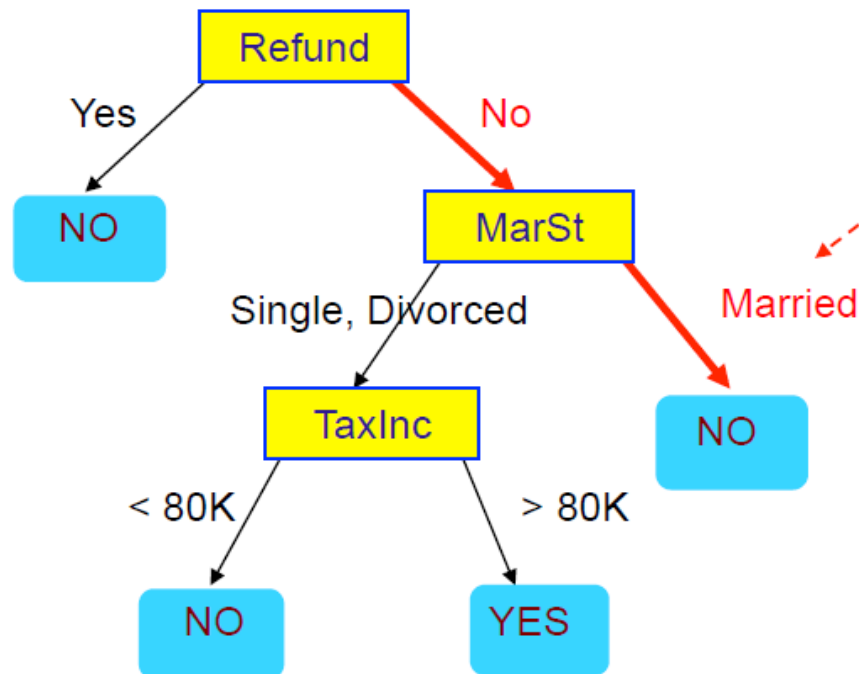
Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



Supervised Learning :

Decision Tree : Model application



Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

Assign Cheat to "No"

Decision Tree : Learning algorithms

- Hunt's Algorithm (one of the earliest)
- CART
- ID3, C4.5
- SLIQ,SPRINT

Decision Tree Induction

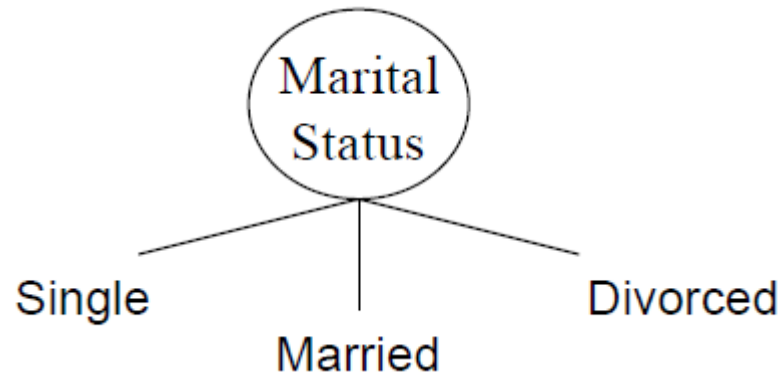
- Greedy strategy : Split the records based on an attribute test that optimizes certain criterion
- Determine how to split the records :
 - How to specify the attribute test condition?
 - How to determine the best split?
- Determine when to stop splitting

Decision Tree Induction : How to Specify Test Condition?

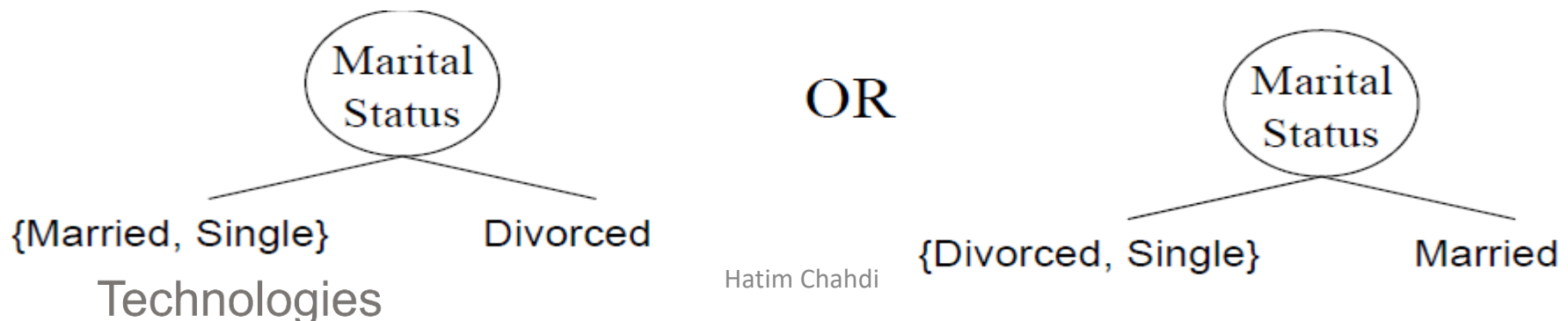
- Depends on attribute types :
Nominal, ordinal or continuous
- Depends on number of ways to split
2-way split, Multi-way split

Decision Tree Induction : Splitting Based on Nominal Attributes

- Multi-way split: Use as many partitions as distinct values

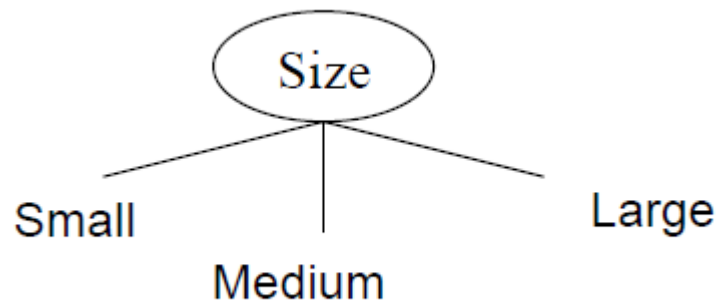


- Binary split: Divides values into two subsets. Need to find optimal partitioning



Decision Tree Induction : Splitting Based on Ordinal Attributes

- We can imagine an attribute SIZE defined over the ordered set {Small, Medium, Large}
- Multi-way split: Use as many partitions as distinct values.

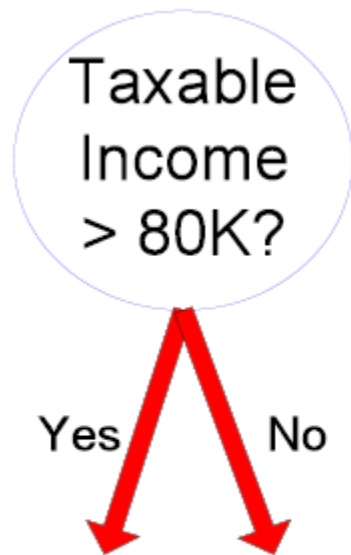


Decision Tree Induction : Splitting Based on Continuous Attributes

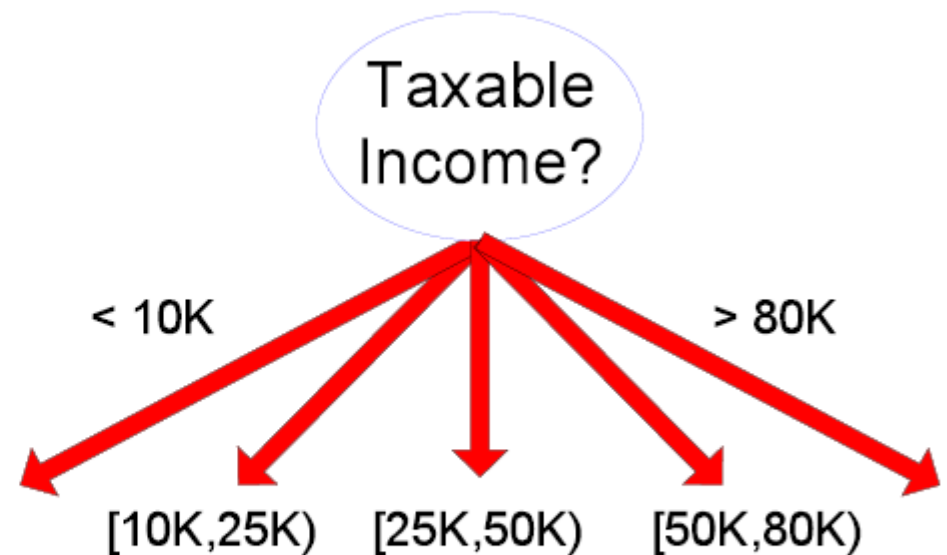
Different ways of handling

- Discretization to form an ordinal categorical attribute
 - Static : discretize once at the beginning
 - Dynamic : ranges can be found by equal interval bucketing, equal frequency bucketing (percentiles), or clustering
- Binary Decision: $(A < v)$ or $(A \geq v)$
 - consider all possible splits and finds the best cut
 - can be more compute intensive

Decision Tree Induction : Splitting Based on Continuous Attributes



(i) Binary split



(ii) Multi-way split

Supervised Learning :

Decision Tree Induction : Best Split

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

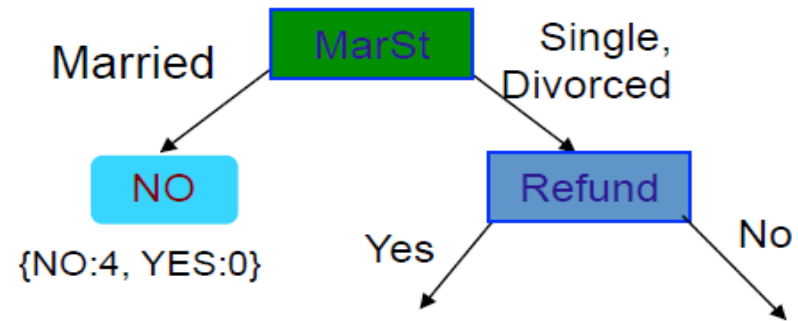


	Yes	No
Married	0	4
Single, Divorced	3	3

Supervised Learning :

Decision Tree Induction : Best Split

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

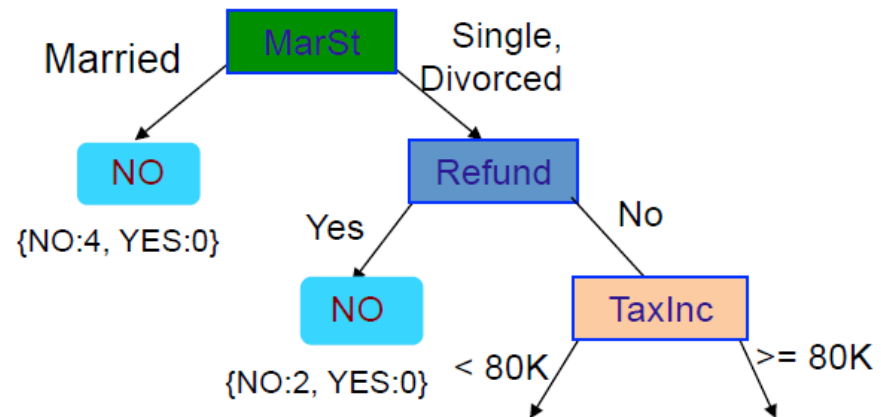


		YES	NO
Single, Divorced	Refund = NO	3	1
Single, Divorced	Refund = Yes	0	2

Supervised Learning :

Decision Tree Induction : Best Split

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

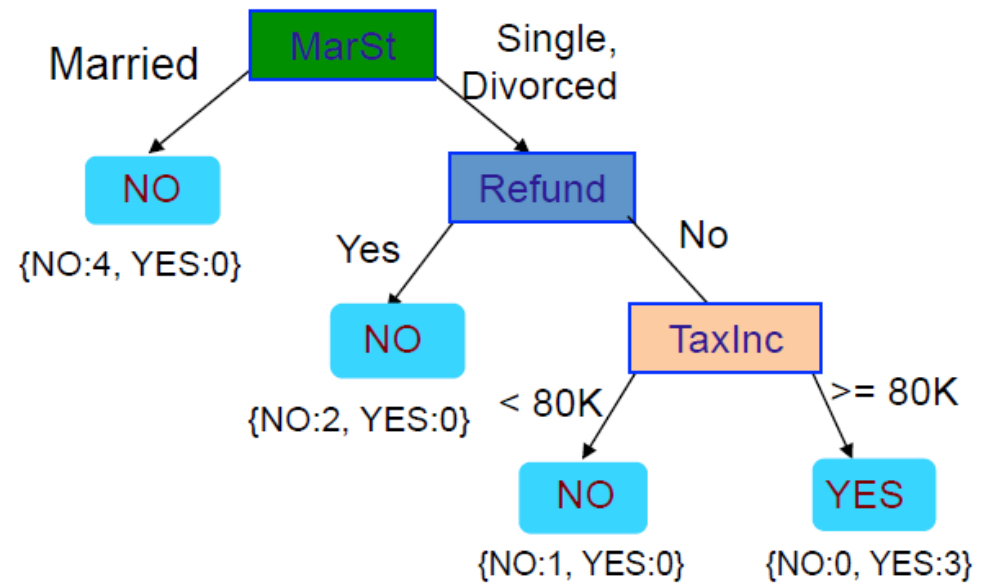


			YES	NO
Single, Divorced	Refund = NO	TaxInc = < 80k	0	1
Single, Divorced	Refund = NO	TaxInc = >= 80k	3	0

Supervised Learning :

Decision Tree Induction : Best Split

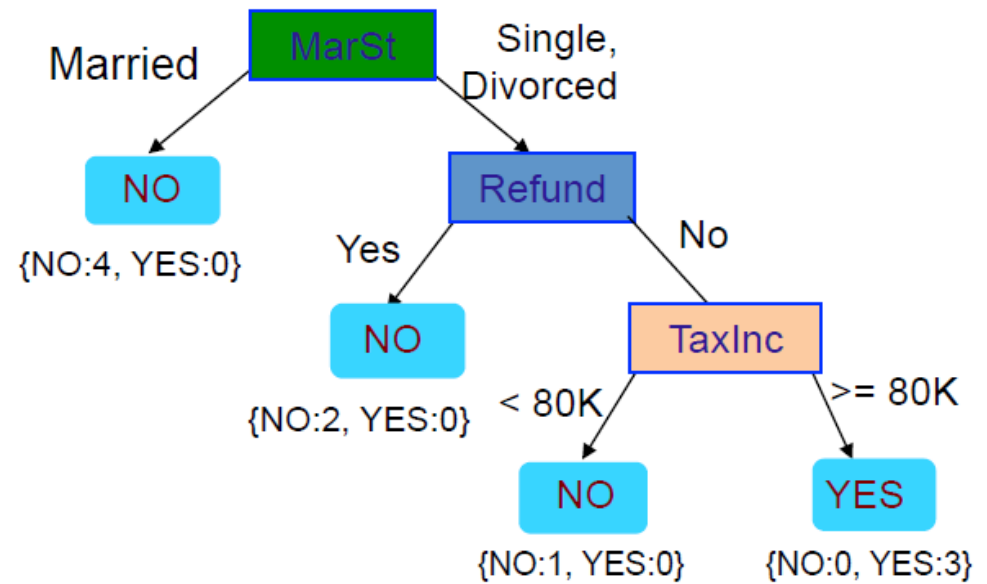
Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Supervised Learning :

Decision Tree Induction : Best Split

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Decision Tree Induction : Stopping criteria

- Stop expanding a node when all the records belong to the same class
- Stop expanding a node when all the records have similar attribute values

Decision Tree Induction : Advantages

- Inexpensive to construct
- Extremely fast at classifying unknown records
- Easy to interpret for small-sized trees
- Accuracy is comparable to other classification
- Techniques for many simple data sets

Naïve Bayes

- Uses probability theory to model the training set
- Assumes independence between attributes
- Produces a model for each class

Naïve Bayes : Principles

- Conditional Probability:
$$P(C | A) = \frac{P(A, C)}{P(A)}$$

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

- Bayes theorem:

$$P(C | A) = \frac{P(A | C)P(C)}{P(A)}$$

Naïve Bayes : Example

Given:

A doctor knows that meningitis causes stiff neck 50% of the time

Prior probability of any patient having meningitis is $1/50,000$

Prior probability of any patient having stiff neck is $1/20$

If a patient has stiff neck, what's the probability he/she has meningitis?

$$P(M | S) = \frac{P(S | M)P(M)}{P(S)} = \frac{0.5 \times 1/50000}{1/20} = 0.0002$$

Supervised Learning :

Naïve Bayes :

Consider each attribute and class label as random variables

- Given a record with attributes (A_1, A_2, \dots, A_n)
 - Goal is to predict class C
 - Specifically, we want to find the value of C that maximizes $P(C | A_1, A_2, \dots, A_n)$
- Can we estimate $P(C | A_1, A_2, \dots, A_n)$ directly from data?

Supervised Learning :

Naïve Bayes : Approach

- Compute the posterior probability $P(C \mid A_1, A_2, \dots, A_n)$ for all values of C using the Bayes theorem

$$P(C \mid A_1 A_2 \mathbf{K} A_n) = \frac{P(A_1 A_2 \mathbf{K} A_n \mid C)P(C)}{P(A_1 A_2 \mathbf{K} A_n)}$$

- Choose value of C that maximizes $P(C \mid A_1, A_2, \dots, A_n)$
- Equivalent to choosing value of C that maximizes $P(A_1, A_2, \dots, A_n \mid C)P(C)$
- How to estimate $P(A_1, A_2, \dots, A_n \mid C)$?

Naïve Bayes : Approach

- Assume independence among attributes A_i when class is given:

$$P(A_1, A_2, \dots, A_n | C) = P(A_1 | C_j) P(A_2 | C_j) \dots P(A_n | C_j)$$

Can estimate $P(A_i | C_j)$ for all A_i and C_j

New point is classified to C_j if $P(C_j) \prod P(A_i | C_j)$ is maximal.

Supervised Learning :

Naïve Bayes : Estimates Probabilities from data

- Class: $P(C) = N_c/N$

e.g., $P(\text{No}) = 7/10$, $P(\text{Yes}) = 3/10$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- For discrete attributes:

$$P(A_i | C_k) = |A_{ik}| / N_c$$

Where : $|A_{ik}|$ is number of instances having attribute A_i and belongs to class C_k

- Examples:

$$P(\text{Status}=\text{Married} | \text{No}) = 4/7$$

$$P(\text{Refund}=\text{Yes} | \text{Yes})=0$$

Supervised Learning :

Naïve Bayes : Estimates Probabilities from data

For continuous attributes:

- Discretize the range into bins :
 - One ordinal attribute per bin
 - Violates independence assumption
- Two-way split: $(A < v)$ or $(A > v)$
 - choose only one of the two splits as new attribute
- Probability density estimation:
 - Assume attribute follows a normal distribution
 - Use data to estimate parameters of distribution (e.g., mean and standard deviation)
 - Once probability distribution is known, can use it to estimate the conditional probability $P(A_i | c)$

Supervised Learning :

Naïve Bayes : Estimates Probabilities from data

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Compute :

- $P(\text{Status}=\text{Married} | \text{Yes}) = ?$
- $P(\text{Refund}=\text{Yes} | \text{No}) = ?$
- $P(\text{Status}=\text{Divorced} | \text{Yes}) = ?$
- $P(\text{TaxableInc} > 80K | \text{Yes}) = ?$
- $P(\text{TaxableInc} > 80K | \text{NO}) = ?$

Supervised Learning :

Naïve Bayes : Estimates Probabilities from data

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Compute :

- $P(\text{Status}=\text{Married} | \text{Yes}) = 0/3$
- $P(\text{Refund}=\text{Yes} | \text{No}) = 3/7$
- $P(\text{Status}=\text{Divorced} | \text{Yes}) = 1/3$
- $P(\text{TaxableInc} > 80K | \text{Yes}) = 3/3$
- $P(\text{TaxableInc} > 80K | \text{NO}) = 4/7$

Supervised Learning :

Naïve Bayes : Estimates Probabilities from data (Exercise)

Given a Test Record: $X = (\text{Refund} = \text{No}, \text{Married}, \text{Income} \geq 80\text{K})$

REFUND	$P(\text{Refund}=\text{Yes} \text{No}) = 3/7$ $P(\text{Refund}=\text{No} \text{No}) = 4/7$ $P(\text{Refund}=\text{Yes} \text{Yes}) = 0$ $P(\text{Refund}=\text{No} \text{Yes}) = 1$
MARITAL STATUS	$P(\text{Marital Status}=\text{Single} \text{No}) = 2/7$ $P(\text{Marital Status}=\text{Divorced} \text{No}) = 1/7$ $P(\text{Marital Status}=\text{Married} \text{No}) = 4/7$ $P(\text{Marital Status}=\text{Single} \text{Yes}) = 2/7$ $P(\text{Marital Status}=\text{Divorced} \text{Yes}) = 1/7$ $P(\text{Marital Status}=\text{Married} \text{Yes}) = 0$
TAXABLE INCOMING	$P(\text{TaxableInc} \geq 80\text{K} \text{Yes}) = 3/3$ $P(\text{TaxableInc} \geq 80\text{K} \text{NO}) = 4/7$ $P(\text{TaxableInc} < 80\text{K} \text{Yes}) = 0/3$ $P(\text{TaxableInc} < 80\text{K} \text{NO}) = 3/7$

Class=No	7/10
Class=Yes	3/10

$$P(C_j | A_1, \dots, A_n) = \frac{P(A_1, \dots, A_n | C_j) P(C_j)}{P(A_1 | C_j) \dots P(A_n | C_j) P(C_j)}$$

Supervised Learning :

Naïve Bayes : Estimates Probabilities from data (Solution)

- $P(X|\text{Class}=\text{No}) = P(\text{Refund}=\text{No}|\text{Class}=\text{No})$
 $\times P(\text{Married}|\text{Class}=\text{No})$
 $\times P(\text{Income} \geq 80\text{K}|\text{Class}=\text{No})$
 $= 4/7 \times 4/7 \times 4/7 = \mathbf{0.1865}$

$$P(X|\text{No})P(\text{No}) = \\ 0.1865 * 0.7 = 0.1306$$

- $P(X|\text{Class}=\text{Yes}) = P(\text{Refund}=\text{No}|\text{Class}=\text{Yes})$
 $\times P(\text{Married}|\text{Class}=\text{Yes})$
 $\times P(\text{Income} \geq 80\text{K}|\text{Class}=\text{Yes})$
 $= 1 \times 0 \times 1 = \mathbf{0}$

$$P(X|\text{Yes})P(\text{Yes}) = \\ 0 * 0.3 = 0$$

Since $P(X|\text{No})P(\text{No}) > P(X|\text{Yes})P(\text{Yes})$

Therefore $P(\text{No}|X) > P(\text{Yes}|X)$
 $\Rightarrow \mathbf{\text{Class} = \text{No}}$

Naïve Bayes : Summary

- Robust to isolated noise points
- Model each class separately
- Robust to irrelevant attributes
- Use the whole set of attribute to perform classification
- Independence assumption may not hold for some attributes

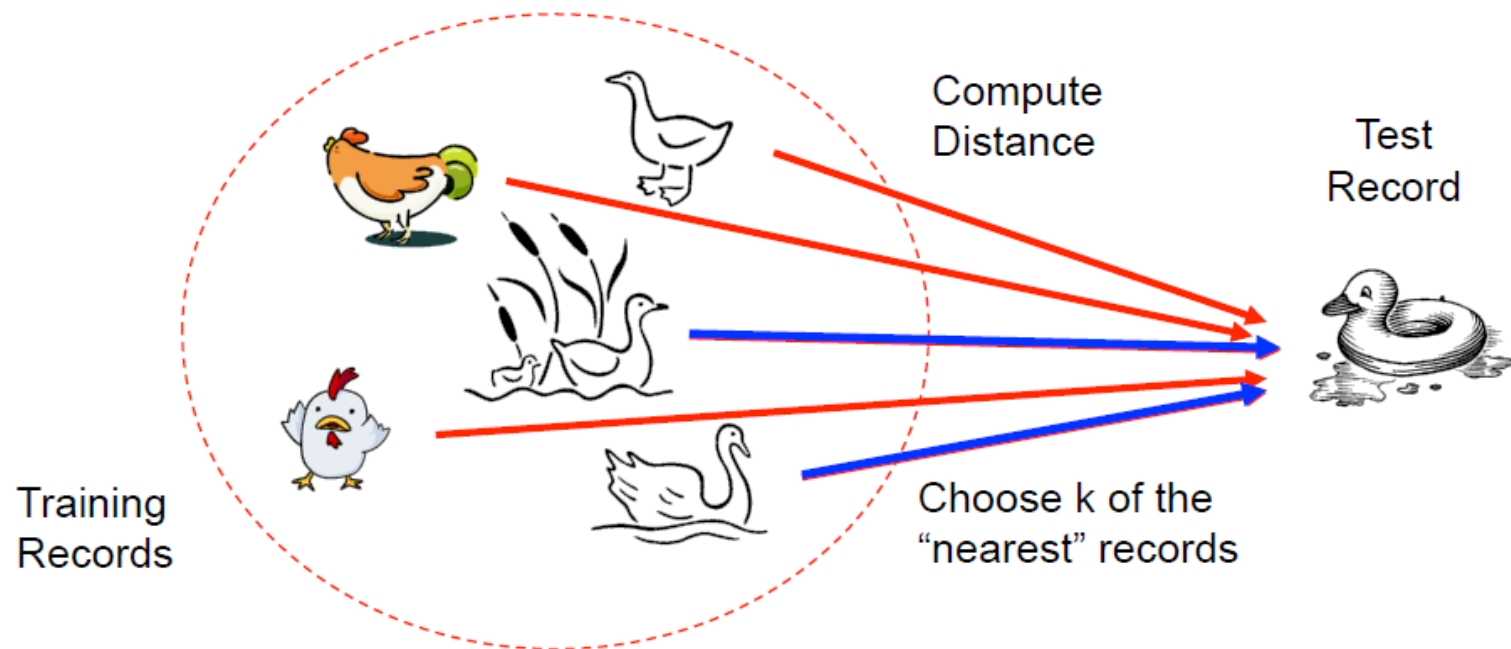
- Lazy approach to classification
- Uses all the training set to perform classification
- Uses distances between training and test records

- Lazy approach to classification
- Uses all the training set to perform classification
- Uses distances between training and test records

Supervised Learning : Instance-Based Classifier (KNN)

Principle :

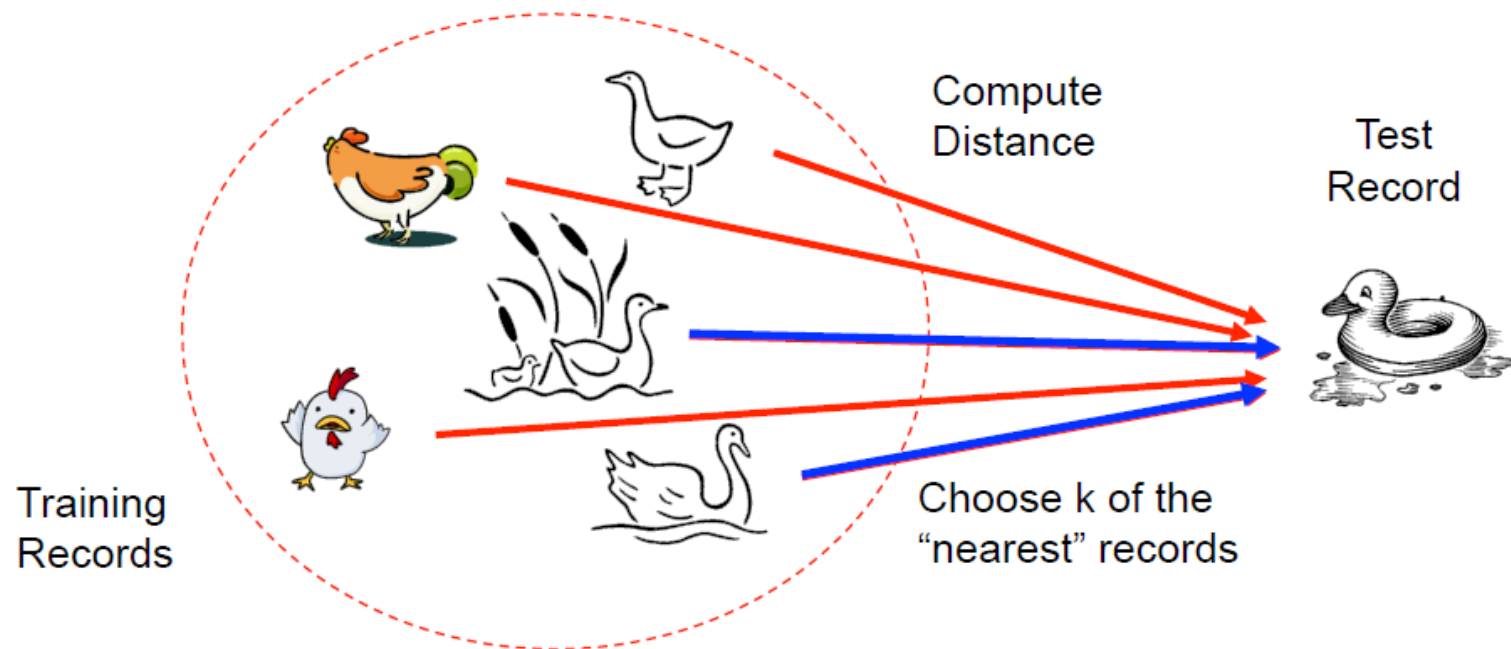
If it walks like a duck, quacks like a duck, then it's probably a duck



Supervised Learning : Instance-Based Classifier (KNN)

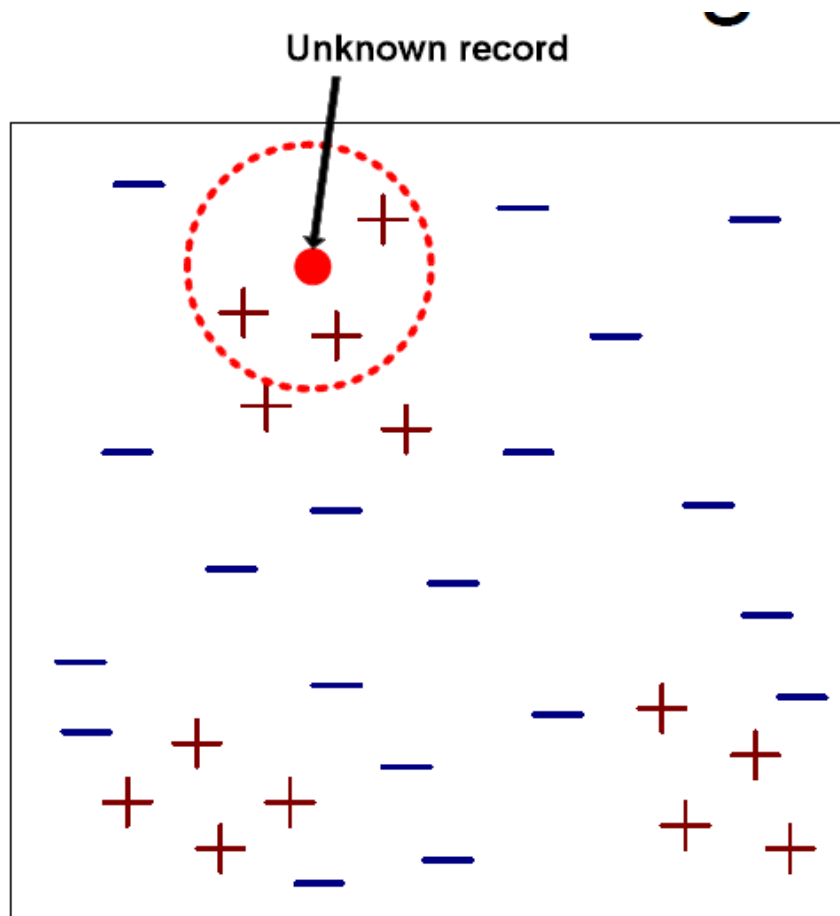
Principle :

If it walks like a duck, quacks like a duck, then it's probably a duck



Supervised Learning : Instance-Based Classifier (KNN)

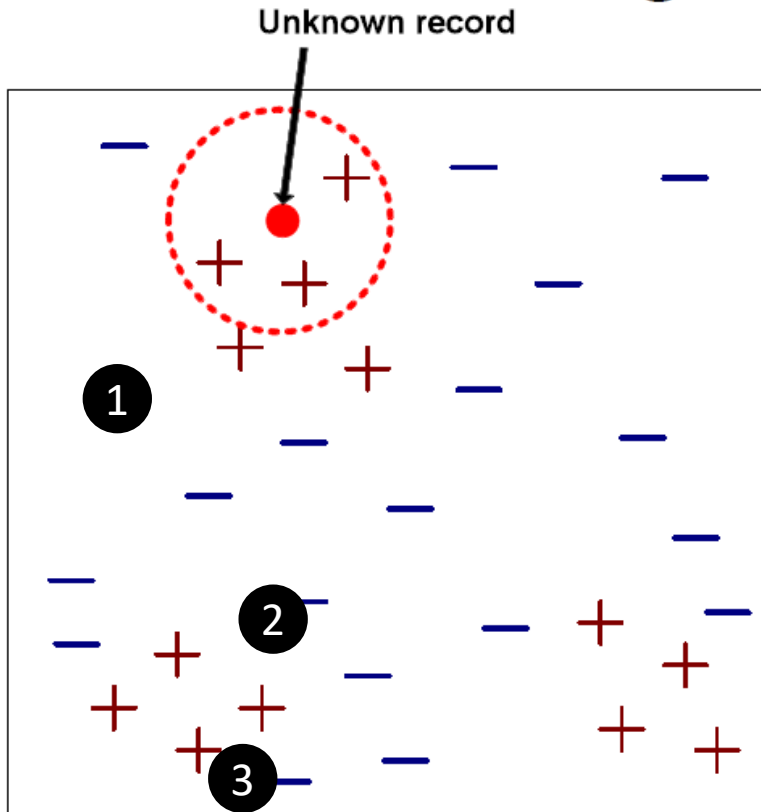
Principle :



- Requires three things
 - The set of stored records
 - Distance Metric to compute distance between records
 - The value of k , the number of nearest neighbors to retrieve
- To classify an unknown record:
 - Compute distance to other training records
 - Identify k nearest neighbors
 - Use class labels of nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote)

Supervised Learning : Instance-Based Classifier (KNN)

Principle :



- Requires three things
 - The set of stored records
 - Distance Metric to compute distance between records
 - The value of k , the number of nearest neighbors to retrieve
- To classify an unknown record:
 - Compute distance to other training records
 - Identify k nearest neighbors
 - Use class labels of nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote)

Exercise : Give the class attributed to each instance for $k = 1, 2$, or 5