

Summary

Summary

- This week:
 - Classification



Classification

1. Decision Trees based Classification
2. Naive Bayesian Classification
3. Support Vector Machines (SVM)

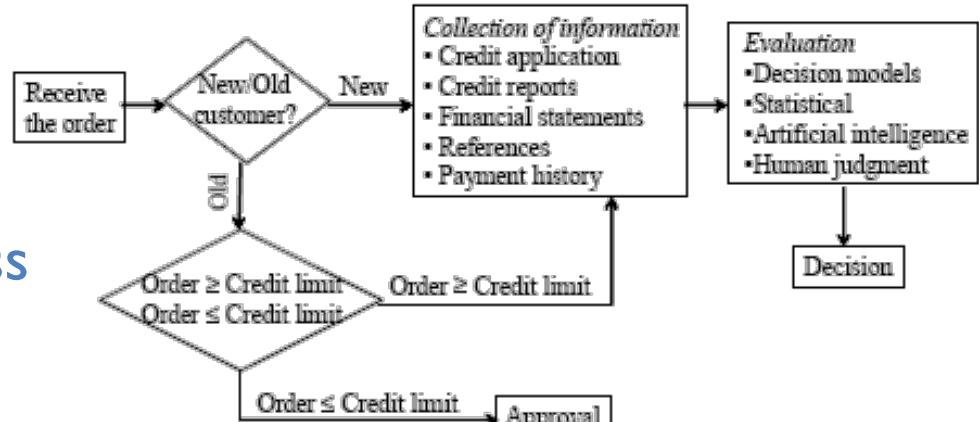


Classification

- What is **classification**?
 - Given is a collection of records (**training set**)
 - Each record consists of a set of attributes, plus a specific **class attribute**
 - Find a model for the class attribute as a **function** of the values of other attributes
 - **Goal:** new records should be assigned to some class as accurately as possible
 - A test set is used to determine the accuracy of the model
 - Usually, the given data set is divided into **training and test sets**, with training set used to **build** the model and test set used to **validate** it

Classification

- How is this useful?
 - Financial field
 - Credit approval process
 - Targeted marketing
 - Classifying credit card transactions as legitimate or fraudulent
 - Categorizing news stories as finance, weather, entertainment, sports, etc.



Classification

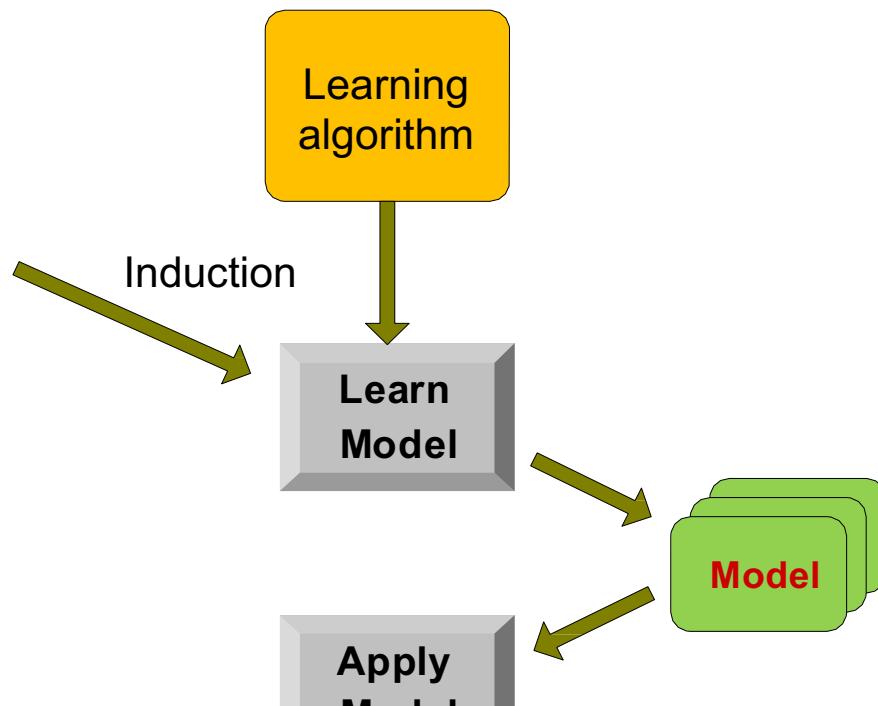
- How does it work?

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

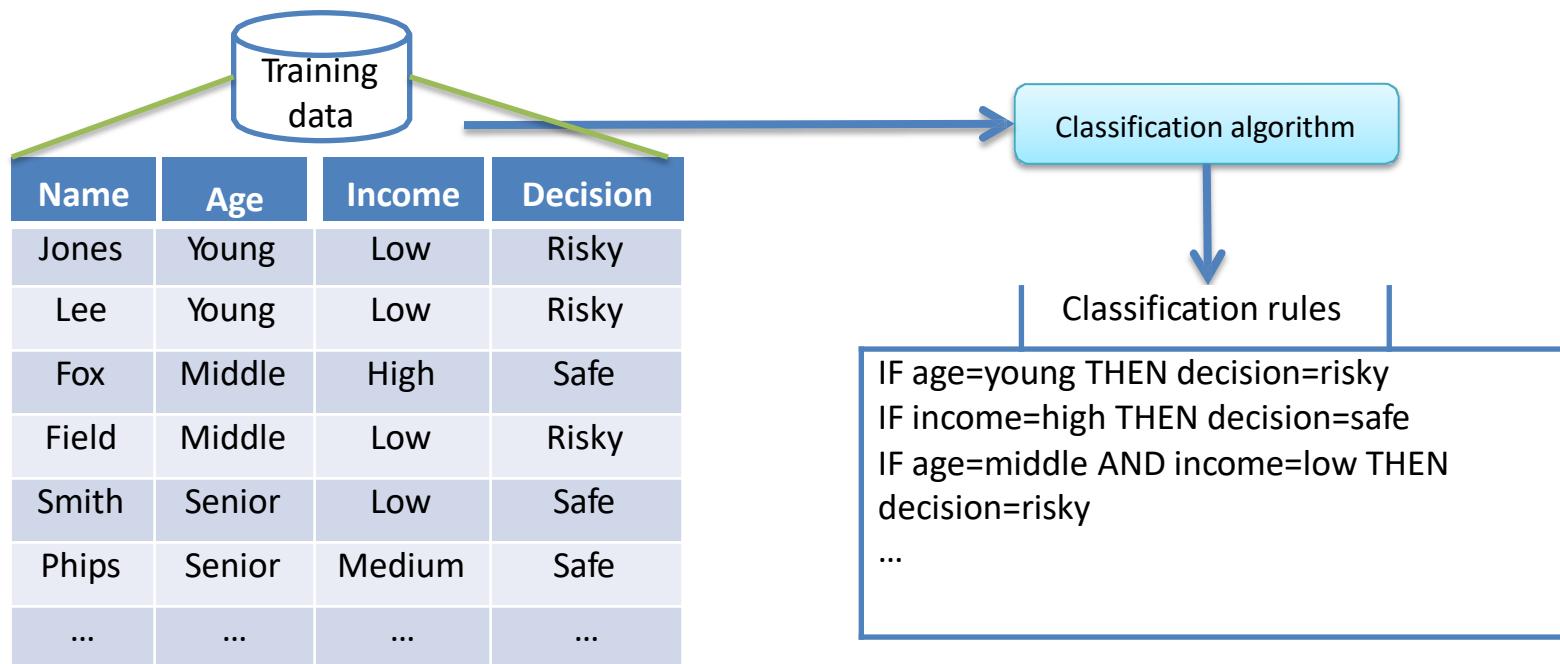
Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



Classification

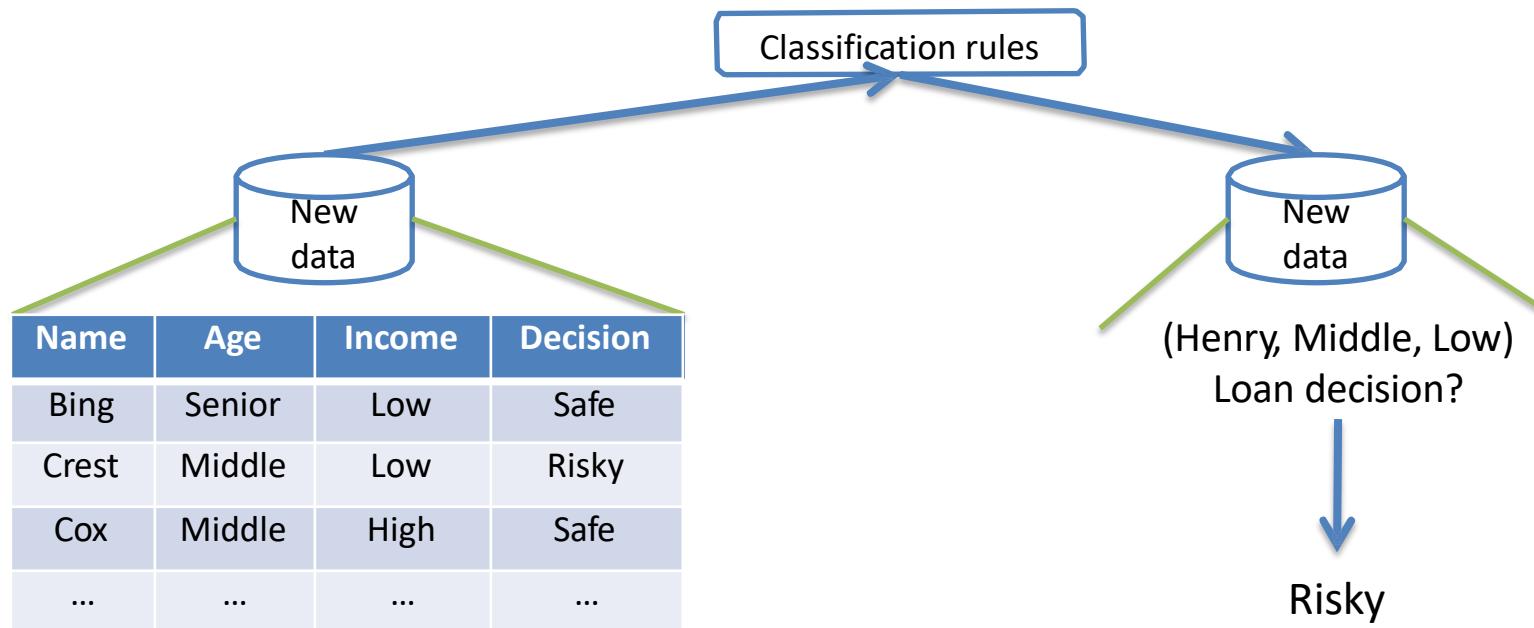
- Example: credit approval
 - Step I: learning (induction)
 - Training data is analyzed by some classification algorithm and the learned model is coded into **classification rules**



Classification

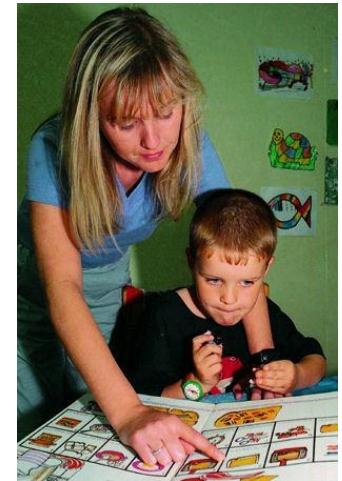
– Step 2: classification (deduction)

- Test data validates the accuracy of the classification rules
- If the accuracy is considered acceptable, then the rules can be applied to the classification of new records



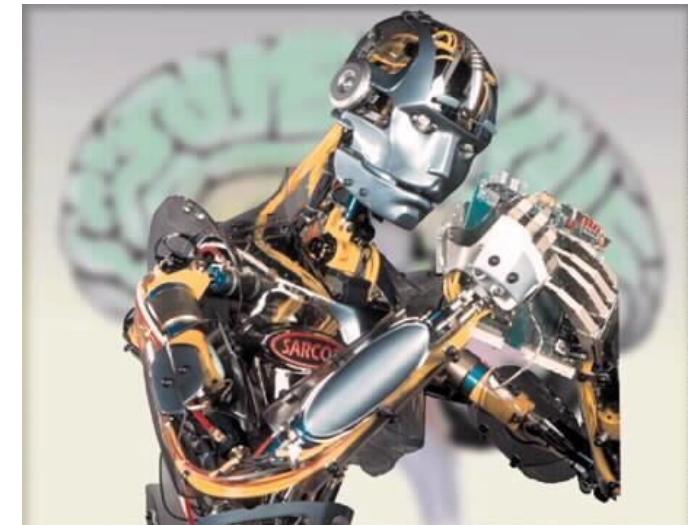
Classification

- **Supervised learning**
 - The training data (observations, measurements, etc.) is accompanied by labels indicating the class of the observations
 - New data is classified based on the training set
- **Unsupervised learning**
 - The class labels of training data is unknown
 - Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data



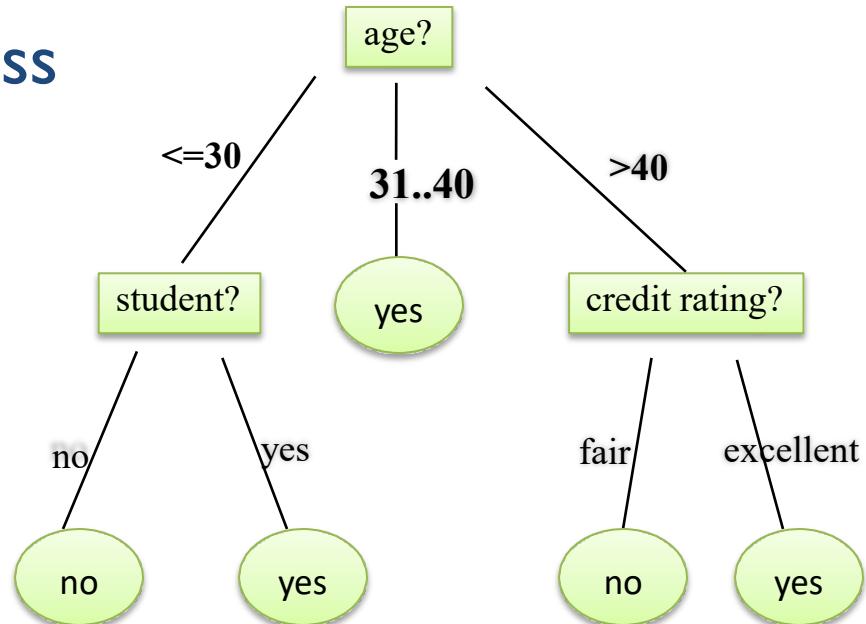
Classification

- Prominent **classification techniques**
 - Decision Tree based Methods
 - Rule-based Methods
 - Naive Bayes and Bayesian Belief Networks
 - Support Vector Machines (SVM)
 - Neural Networks



Decision Tree

- **Decision tree**
 - A flow-chart-like **tree** structure
 - **Internal node** denotes a test on an attribute
 - **Branch** represents an outcome of the test
 - **Leaf nodes** represent class labels or class distribution
 - E.g., decision making
 - Who buys a computer?



Decision Trees

- Classification based on decision tree example

- Step 1 – Induction

Generate the decision tree

- Input: **training data set**, attribute list used for classification, attribute selection method
 - Output: the decision tree

- Step 2 – Deduction

Predict the classes of the **new data**

- Input: the decision tree from step 1 and the new data
 - Output: classified new data

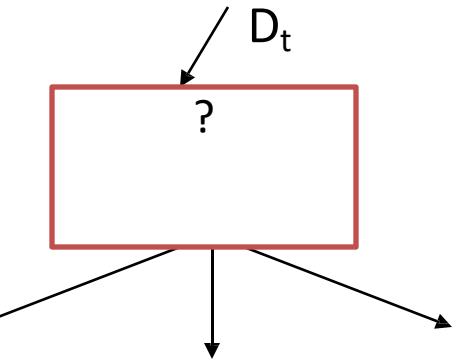
Decision Tree Induction

- Decision tree Algorithms
 - Many Algorithms:
 - Hunt's Algorithm
 - One of the earliest methods
 - ID3 and its successor, C4.5
 - Represents a benchmark to supervised learning algorithms
 - Based on the Hunt algorithm
 - Classification and Regression Trees (CART)
 - Similar to ID3
 - SLIQ, SPRINT

Decision Tree Induction

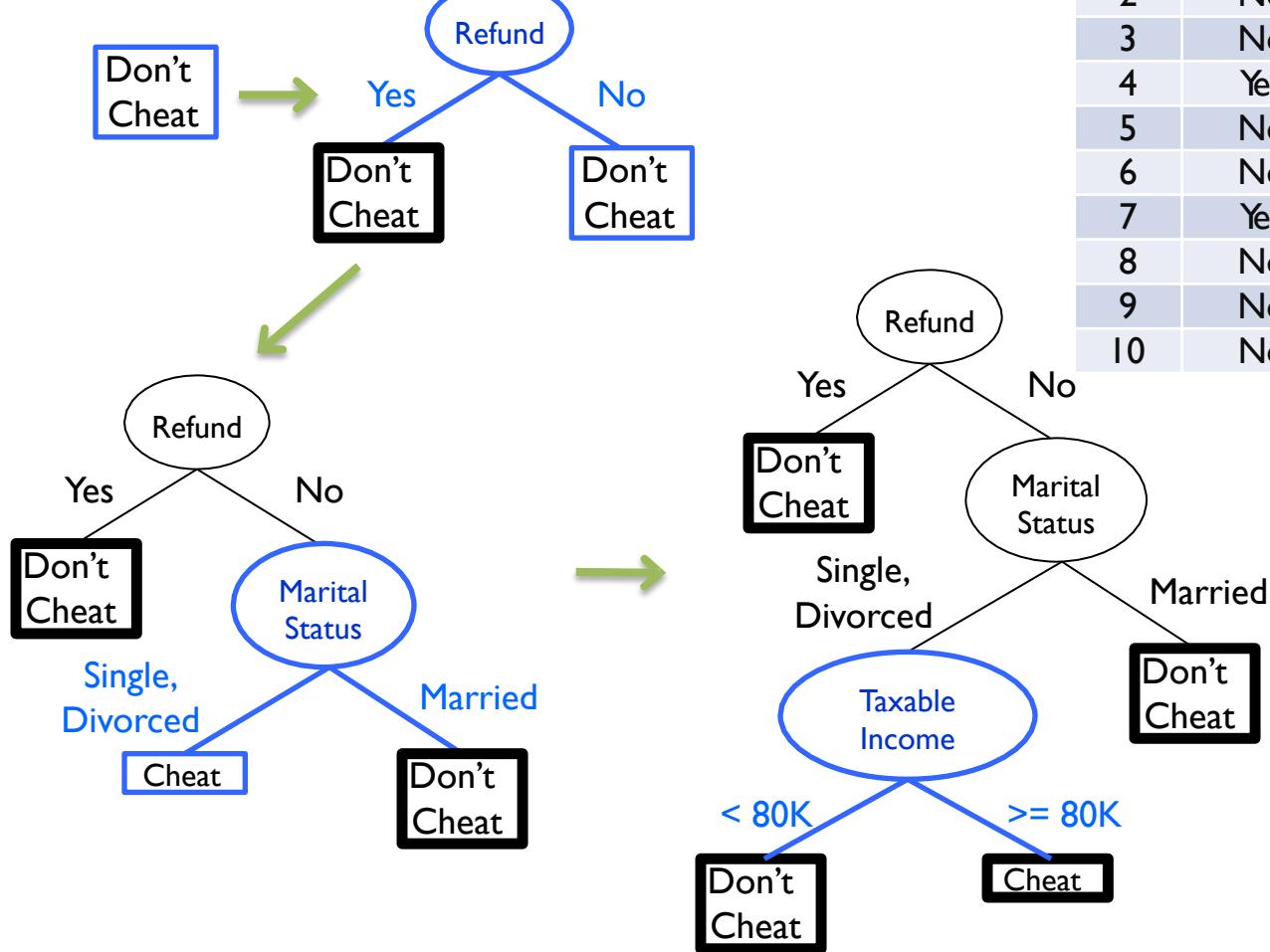
- **Hunt's algorithm, general structure**

- Let D_t be the set of training records that reach a node t
- General Procedure:
 - If D_t contains records that **belong the same class** y_t , then t is a **leaf node** labeled as y_t
 - If D_t is an empty set, then t is a leaf node labeled by the default class, y_d
 - If D_t contains records that belong to **more than one class**, use an attribute test to split the data into smaller subsets: **recursively apply the procedure to each subset**



Hunts Algorithm

- Example, VAT refund



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

Hunts Algorithm

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



Attribute test condition

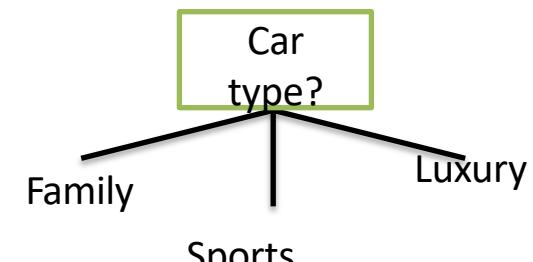
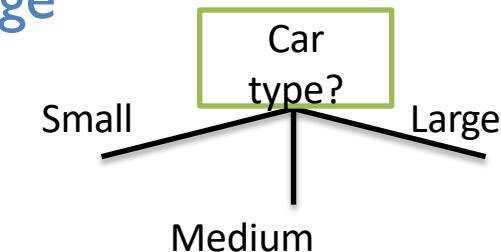
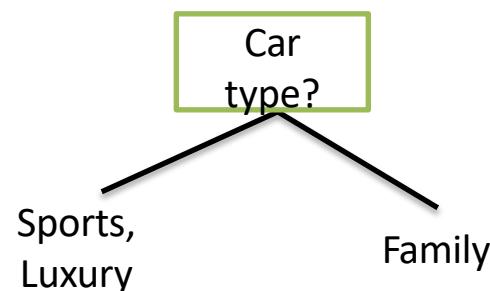
- **Splitting:** how to specify the attribute test condition?

- Depends on **attribute types**

- Nominal (discrete values) e.g., car:sports, luxury, family
 - Ordinal (categories) e.g., small, medium, large
 - Continuous (numbers) e.g., age

- Depends on **number of ways** to split

- Binary split
 - Multi-way split

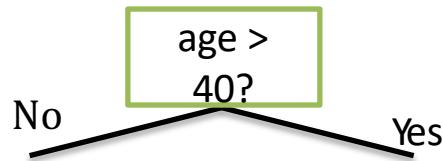


Attribute test condition

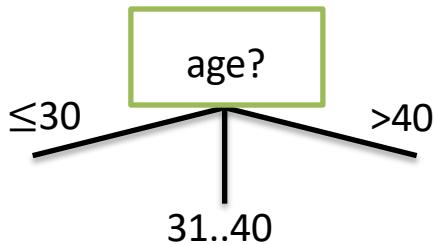
- What about splitting **continuous attributes**?
 - Discretization to form an ordinal categorical attribute
 - Static – discretize once at the beginning
 - Dynamic – ranges can be found by equal interval bucketing, equal frequency bucketing (percentiles), clustering, or supervised clustering
 - Binary decision
 - Consider all possible splits and finds the best cut
 - Can be quite computationally expensive

Attribute test condition

- Splitting **continuous attributes** (e.g., age)
 - Binary split



- Multi-way split



age	income	student	Credit rating	Buys computer
27	high	no	fair	no
28	high	no	excellent	no
31	high	no	fair	yes
45	medium	no	excellent	yes
43	low	yes	excellent	yes
56	low	yes	fair	no
37	low	yes	excellent	yes
20	medium	no	fair	no
20	low	yes	fair	yes
60	medium	yes	excellent	yes
24	medium	yes	excellent	yes
36	medium	no	excellent	yes
31	high	yes	fair	yes
41	medium	no	fair	no

Determine the best split

- How do we determine the best split?

- We can split on **any** of the 4 attributes!

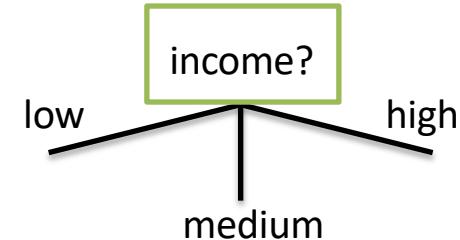
- E.g., income

- Low – yes:3, no:1
 - Medium – yes:4, no:2
 - High – yes:2, no:2

- E.g., student

- No – yes:3, no:4
 - Yes – yes:6, no:1

- Which split is **better**?



age	income	student	Credit rating	Buys computer
27	high	no	fair	no
28	high	no	excellent	no
31	high	no	fair	yes
45	medium	no	excellent	yes
43	low	yes	excellent	yes
56	low	yes	fair	no
37	low	yes	excellent	yes
20	medium	no	fair	no
20	low	yes	fair	yes
60	medium	yes	excellent	yes
24	medium	yes	excellent	yes
36	medium	no	excellent	yes
31	high	yes	fair	yes
41	medium	no	fair	no

Determine the best split

- What does better mean?
 - Nodes with **homogeneous** class distribution (pure nodes) are preferred
 - E.g., homogeneous nodes
 - Student attribute, Yes – yes:6, no:1
 - E.g., heterogeneous nodes
 - Income attribute, High – yes:2, no:2
- How do we measure node impurity?



Determine the best split

- Methods to measure impurity
 - **Information gain** (used in C4.5)
 - All attributes are assumed to be categorical
 - Can be modified for continuous-valued attributes
 - Also Kullback–Leibler divergence
 - **Gini index**
 - All attributes are assumed continuous-valued
 - Assume there exist several possible split values for each attribute
 - May need other tools, such as clustering, to get the possible split values

Decision Tree Induction

- **Information gain**

- Method

- Assume there are two classes, P and N

- Let the set of examples S contain p elements of class P and n elements of class N
 - The amount of information, needed to decide whether an arbitrary example in S belongs to P or N is defined as

$$I(p, n) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

- Select the attribute with the **highest** information gain

estimates the probability that label is p

estimates the probability that label is n

Information Gain

- Information gain in decision tree induction
 - Assume that using attribute A a set S will be partitioned into sets $\{S_1, S_2, \dots, S_v\}$
 - If S_i contains p_i examples of P and n_i examples of N, the **entropy**, or the expected information needed to classify objects in all subtrees S_i is

$$E(A) = \sum_{i=1}^v \frac{p_i + n_i}{p + n} I(p_i, n_i)$$

- The encoding information that would be gained by branching on A is

$$Gain(A) = I(p, n) - E(A)$$

Stop Condition

- Stop conditions
 - All the records **belong to the same class**
 - E.g.,

income	credit	class
high	fair	no
high	excellent	no
medium	fair	no

- In this case a leaf node is created with the corresponding class label (here “no”)



Stop Condition

- Stop conditions
 - All the records **have similar attribute values**
 - E.g., perform split by student but all records are students

student	credit	class
yes	fair	no
yes	fair	no
yes	fair	no
yes	fair	yes
yes	excellent	yes

- In this case instead of performing the split, a leaf node is created with the **majority class** as label (here “no”)

Classification Example using Decision Tree

Decision Tree Induction

- Step I
 - Input
 - Training set data
 - Use all attributes for classification
 - Use information gain as selection method

age	income	student	Credit rating	Buys computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	excellent	yes
>40	low	yes	excellent	yes
>40	low	yes	fair	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	excellent	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	fair	no

Decision Tree Induction

- Attribute selection by gain computation, example:
 - Class P: `buys_computer = “yes”`
 - Class N: `buys_computer = “no”`

$$I(p, n) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

$$- I(p, n) = I(9, 5) = 0.94$$

age	income	student	Credit rating	Buys computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	excellent	yes
>40	low	yes	excellent	yes
>40	low	yes	fair	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	excellent	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	fair	no

Decision Tree Induction

- Compute the entropy for the following 3 age partitions

age ≤ 30 , $30 < \text{age} \leq 40$

and $40 < \text{age}$

$$E(A) = \sum_{i=1}^v \frac{p_i + n_i}{p + n} I(p_i, n_i) \quad \rightarrow$$

$$\begin{aligned} E(\text{age}) &= \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) \\ &\quad + \frac{5}{14} I(3,2) = 0.694 \end{aligned}$$

- Gain(age) = $I(p, n) - E(\text{age}) \Rightarrow 0.94 - 0.694 = 0.246$

- The same we can calculate also the gains for

- Gain(income) = 0.029
- Gain(student) = 0.151
- Gain(credit_rating) = 0.048

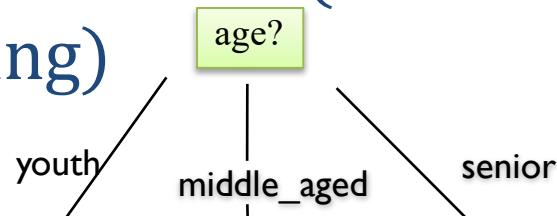
age	p_i	n_i	$I(p_i, n_i)$
≤ 30	2	3	0,971
31...40	4	0	0
>40	3	2	0,971

Decision Tree Induction

- First node was already calculated

– Gain(age) = 0.246 > Gain(income), Gain(student), Gain(credit_rating)

income	student	credit	class
high	no	fair	no
high	no	excellent	no
medium	no	fair	no
low	yes	fair	yes
medium	yes	excellent	yes



income	student	credit	class
medium	no	excellent	yes
low	yes	excellent	yes
low	yes	fair	no
medium	yes	excellent	yes
medium	no	fair	no

income	student	credit	class
high	no	fair	yes
low	yes	excellent	yes
medium	no	excellent	yes
high	yes	fair	yes

age	p _i	n _i	I(p _i , n _i)
<=30	2	3	0,971
30...40	4	0	0
>40	3	2	0,971

Decision Tree Induction

age?
youth

- Subnode age = youth
 - For the income attribute
 - $I(2, 3) = 0.97$
 - $E(\text{income}) = 1/5 * I(1, 0) + 2/5 * I(1, 1) + 2/5 * I(0, 2) = 0.4$
 - Thus $\text{Gain}(\text{youth}, \text{income}) = 0.97 - 0.4 = 0.57$
 - For the student attribute
 - I is the same
 - $E(\text{student}) = 2/5 * I(2, 0) + 3/5 * I(0, 3) = 0$
 - Thus $\text{Gain}(\text{youth}, \text{student}) = 0.97$

income	student	credit	class
high	no	fair	no
high	no	excellent	no
medium	no	fair	no
low	yes	fair	yes
medium	yes	excellent	yes

Class P: buys_computer = "yes"

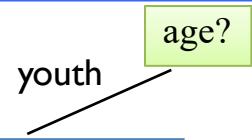
Class N: buys_computer = "no"

$$I(p, n) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

$$E(A) = \sum_{i=1}^v \frac{p_i + n_i}{p + n} I(p_i, n_i)$$

Decision Tree Induction

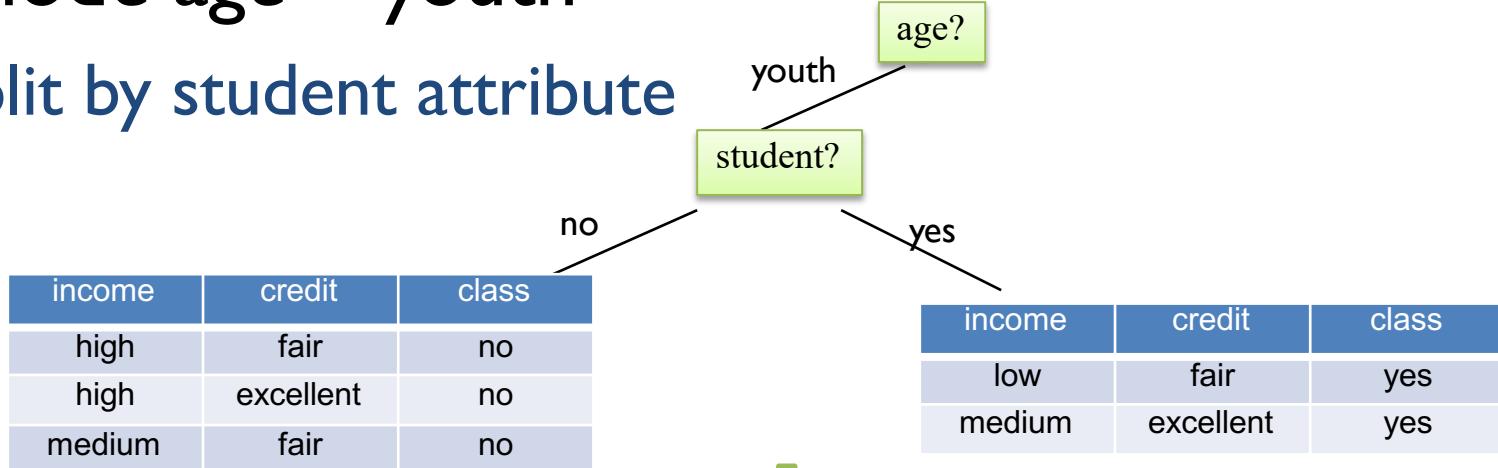
- Subnode age = youth
 - For the credit attribute
 - $I(2, 3) = 0.97$
 - $E(\text{credit}) = 3/5 * I(1, 2) + 2/5 * I(1, 1) = 0.95$
 - Thus $\text{Gain}(\text{youth}, \text{credit}) = 0.97 - 0.95 = 0.02$
 - Largest gain was of 0.97 for the **student** attribute



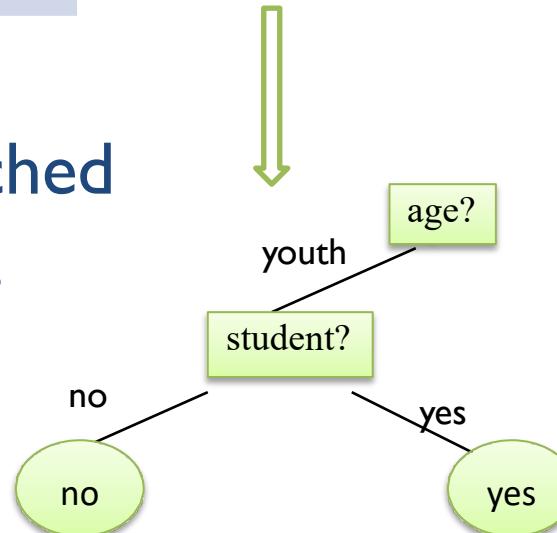
income	student	credit	class
high	no	fair	no
high	no	excellent	no
medium	no	fair	no
low	yes	fair	yes
medium	yes	excellent	yes

Decision Tree Induction

- Subnode age = youth
 - Split by student attribute



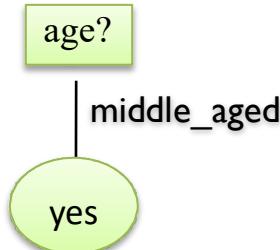
- Stop condition reached
 - Resulting subtree is



Decision Tree Induction

- Subnode age = middle_aged
 - **Stop condition reached**
 - We have just one class

age?			
middle_aged			
income	student	credit	class
high	no	fair	yes
low	yes	excellent	yes
medium	no	excellent	yes
high	yes	fair	yes



Decision Tree Induction

- Subnode age = senior

- For the income attribute

- $I(3, 2) = 0.97$

- $E(\text{income}) = 2/5 * I(1, 1) + 3/5 * I(2, 1) = 0.95$

- Thus $\text{Gain}(\text{senior}; \text{income}) = 0.97 - 0.95 = 0.02$

- For the student attribute

- I is the same

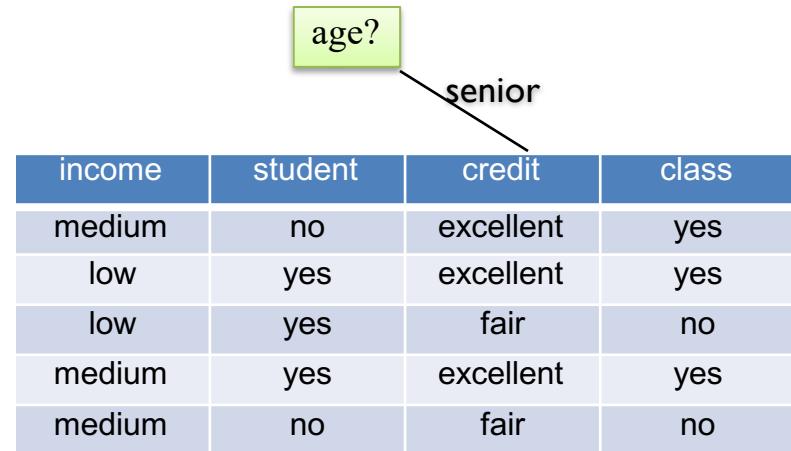
- $E(\text{student}) = 3/5 * I(2, 1) + 2/5 * I(1, 1) = 0.95$

- Thus $\text{Gain}(\text{senior}, \text{student}) = 0.02$

income	student	credit	class
medium	no	excellent	yes
low	yes	excellent	yes
low	yes	fair	no
medium	yes	excellent	yes
medium	no	fair	no

Decision Tree Induction

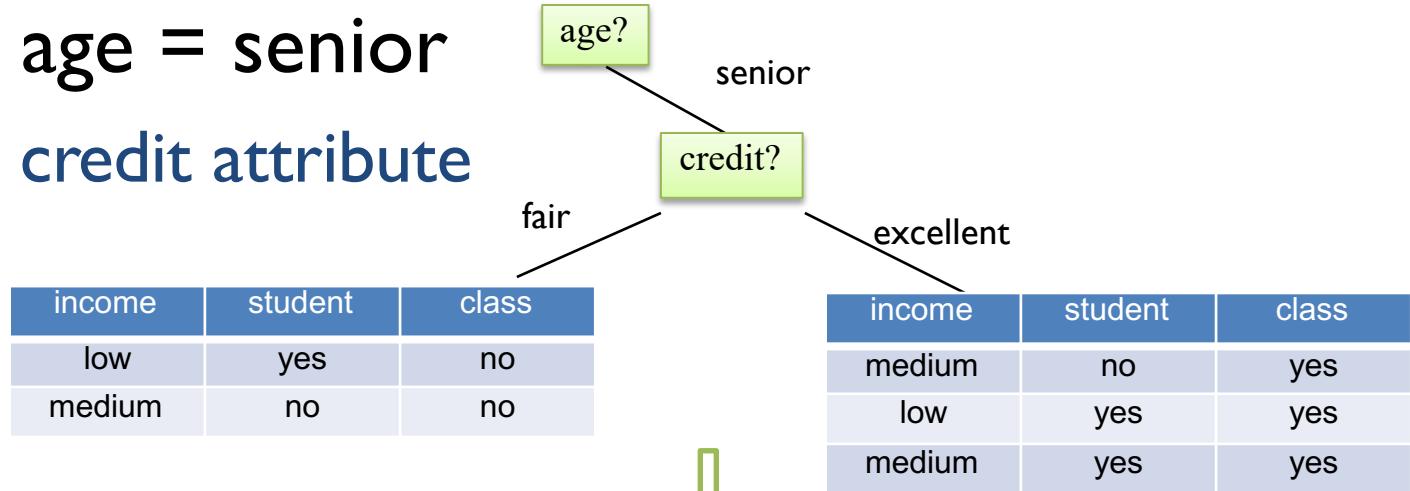
- Subnode age = senior
 - For the credit attribute
 - I is the same
 - $E(\text{income}) = 2/5 * I(0, 2) + 3/5 * I(3, 0) = 0$
 - Thus $\text{Gain}(\text{senior}; \text{credit}) = 0.97$
 - Thus **split by credit attribute**



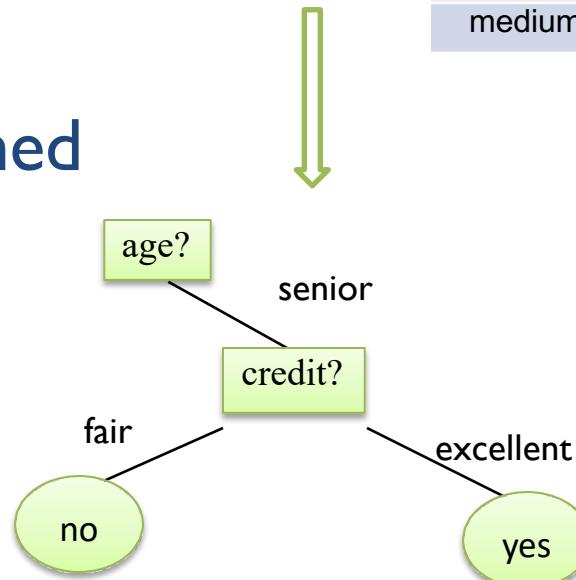
income	student	credit	class
medium	no	excellent	yes
low	yes	excellent	yes
low	yes	fair	no
medium	yes	excellent	yes
medium	no	fair	no

Decision Tree Induction

- Subnode age = senior
 - Split by credit attribute

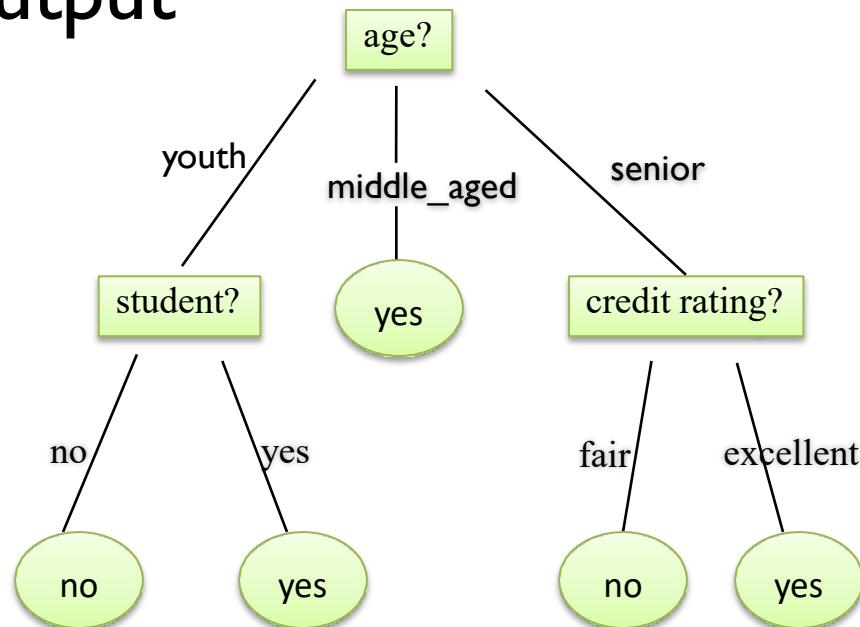


- Stop condition reached



Decision Tree Induction

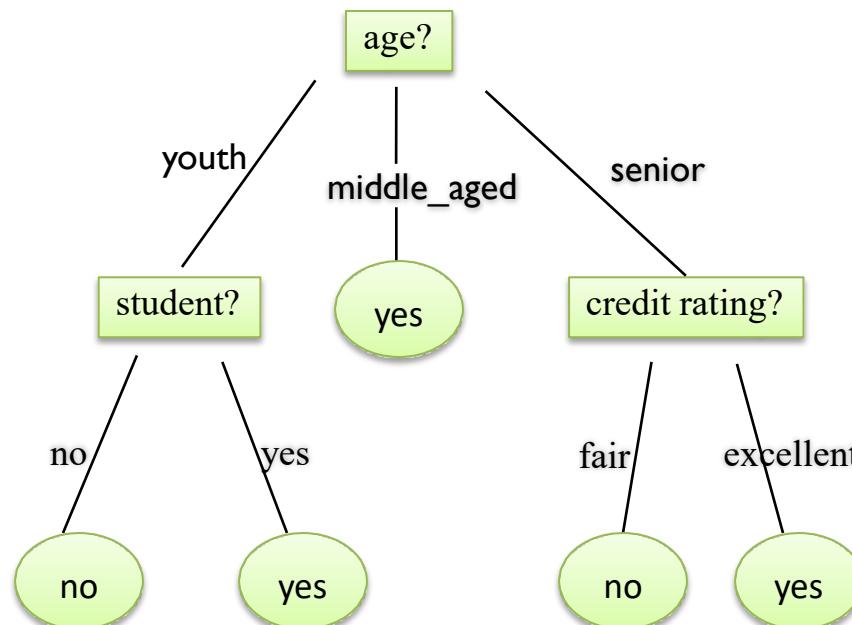
- Step I has finished with the following decision tree as output



Decision Tree Deduction

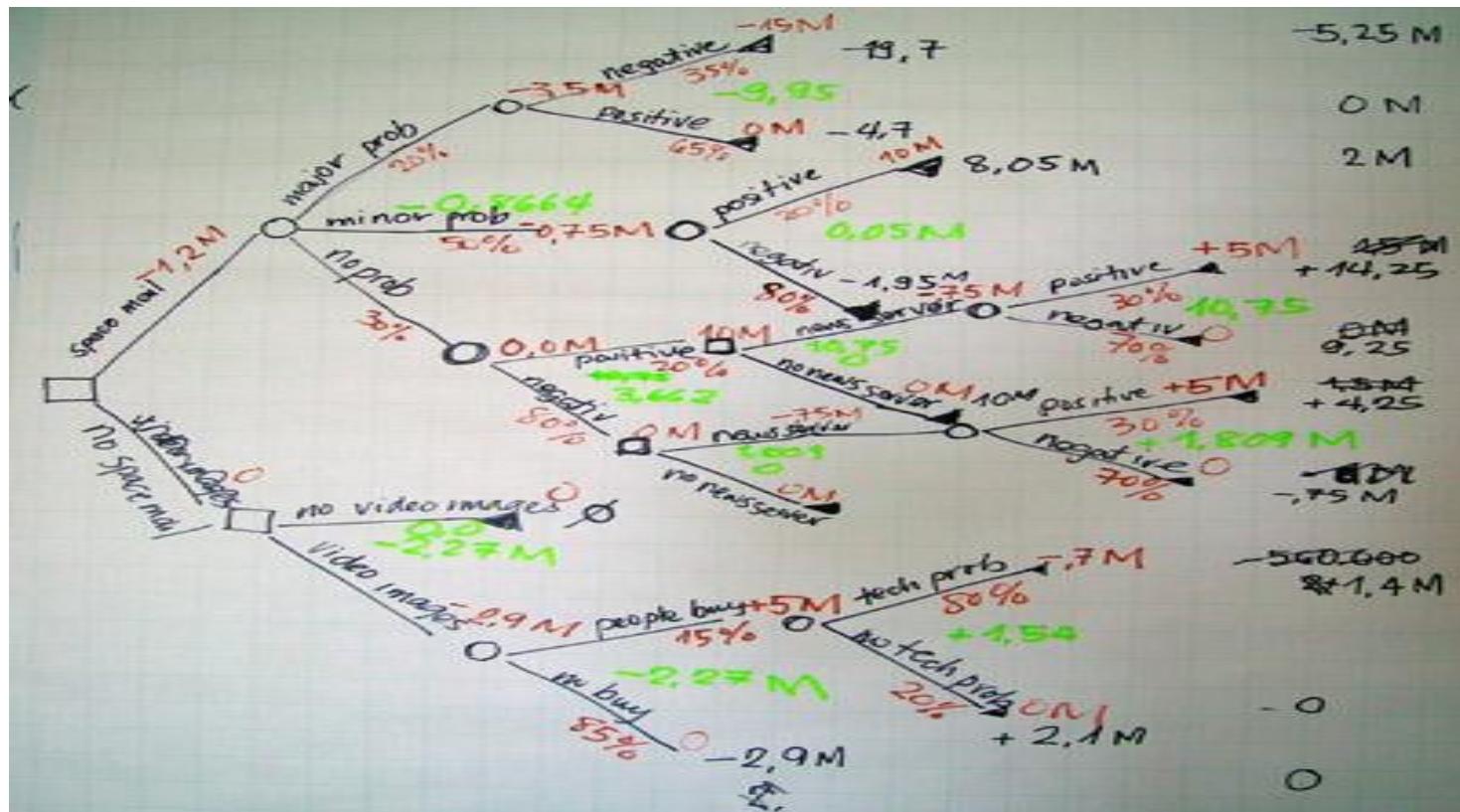
- Step 2
 - New data

age	income	student	Credit rating	Buy computer
31...40	low	yes	fair	yes
<=30	low	yes	fair	yes
<=30	low	yes	excellent	yes
>40	low	no	fair	no



Decision Tree Deduction

- Decision tree deduction
 - Use the decision tree rules to classify new data



Decision Trees

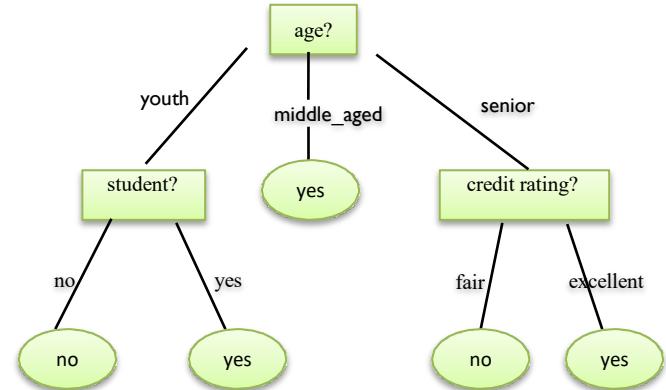
- Extracting classification rules from trees
 - Represent the knowledge in the form of **IF-THEN rules**
 - **One rule** is created for **each path** from the root to a leaf
 - Each attribute-value pair along a path forms a **conjunction**
 - The leaf node holds the class prediction
 - Rules are easier for humans to understand



Decision Trees

– Example

- IF age = “ ≤ 30 ” AND student = “no”
THEN buys_computer = “no”
- IF age = “ ≤ 30 ” AND student = “yes”
THEN buys_computer = “yes”
- IF age = “31...40” THEN buys_computer = “yes”
- IF age = “>40” AND credit_rating = “excellent” THEN
buys_computer = “yes”
- IF age = “>40” AND credit_rating = “fair” THEN
buys_computer = “no”



Decision Trees

- 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
 - Very good average performance over many datasets

Decision Trees

- Avoid over-fitting in classification
 - The generated tree may over-fit the training data
 - Too many branches, some may reflect anomalies due to noise or outliers e.g. a person with blue eye colour will buy a computer.
 - Result is in poor accuracy for unseen samples
 - Two approaches to avoid over-fitting
 - Prepruning
 - Postpruning



Decision Trees

- Prepruning
 - Halt tree construction early
 - Do not split a node if this would result in the information gain falling below a threshold
 - Difficult to choose an appropriate threshold
- Postpruning
 - Remove branches from a “fully grown” tree
 - Get a sequence of progressively pruned trees
 - Use a set of data different from the training data to decide which is the “best pruned tree”

Decision Trees

- Enhancements
 - Allow for continuous-valued attributes
 - Dynamically define new discrete-valued attributes that partition the continuous attribute value into a discrete set of intervals
 - Handle missing attribute values
 - Assign the most common value of the attribute
 - Assign probability to each of the possible values
 - Attribute construction
 - Create new attributes based on existing ones that are sparsely represented
 - This reduces fragmentation, repetition, and replication

Classification

- Decision Trees based Classification
 - Induction
 - Deduction

Next lecture

- Naïve Bayesian Classification