### **Chapter#2 Supervised Learning**

#### Learning from Observations

## What is Learning?

- Herbert Simon: "Learning is any process by which a system improves performance from experience."
- "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 Mitchell

# Learning

- Learning is essential for unknown environments,
   i.e., when designer lacks omniscience
- Learning is useful as a system construction method,
  - i.e., expose the agent to reality rather than trying to write it down
- Learning modifies the agent's decision mechanisms to improve performance

## **Machine Learning**

- Machine learning: how to acquire a model on the basis of data / experience
  - Learning parameters (e.g. probabilities)
  - Learning structure (e.g. BN graphs)
  - Learning hidden concepts (e.g. clustering)

## **Machine Learning Areas**

- Supervised Learning: Data and corresponding labels are given
- Unsupervised Learning: Only data is given, no labels provided
- Semi-supervised Learning: Some (if not all) labels are present
- Reinforcement Learning: An agent interacting with the world makes observations, takes actions, and is rewarded or punished; it should learn to choose actions in such a way as to obtain a lot of reward

## **Supervised Learning : Important Concepts**

- Data: labeled instances <*x<sub>i</sub>*, *y*>, e.g. emails marked spam/not spam
  - Training Set
  - Held-out Set
  - Test Set
- Features: attribute-value pairs which characterize each x
- Experimentation cycle
  - Learn parameters (e.g. model probabilities) on training set
  - (Tune hyper-parameters on held-out set)
  - Compute accuracy of test set
  - Very important: never "peek" at the test set!
- Evaluation
  - Accuracy: fraction of instances predicted correctly
- Overfitting and generalization
  - Want a classifier which does well on test data
  - Overfitting: fitting the training data very closely, but not generalizing well

# **Example: Spam Filter**

Input: email Output: spam/ham Setup:



- Get a large collection of 0 example emails, each labeled "spam" or "ham"
- Note: someone has to hand label all this data!



Want to learn to predict labels of new, future emails

Features: The attributes used to make the ham / spam decision

Words: FREE! 0

0

...

- Text Patterns: \$dd, CAPS 0
- Non-text: SenderInContacts



Dear Sir.



First, I must solicit your confidence in this transaction, this is by virture of its nature as being utterly confidencial and top secret. ...

TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

99 MILLION EMAIL ADDRESSES FOR ONLY \$99

Ok, Iknow this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

# **Example: Digit Recognition**

Input: images / pixel grids Output: a digit 0-9 Setup:

- Get a large collection of example images, each labeled with a digit
- Note: someone has to hand label all this data!
- Want to learn to predict labels of new, future digit images

Features: The attributes used to make the digit decision

• Pixels: (6,8)=ON

° ...

 Shape Patterns: NumComponents, AspectRatio, NumLoops



ŏ

## **Classification Examples**

- In classification, we predict labels y (classes) for inputs x
- Examples:
  - OCR (input: images, classes: characters)
  - Medical diagnosis (input: symptoms, classes: diseases)
  - Automatic essay grader (input: document, classes: grades)
  - Fraud detection (input: account activity, classes: fraud / no fraud)
  - Customer service email routing
  - Recommended articles in a newspaper, recommended books
  - DNA and protein sequence identification
  - Categorization and identification of astronomical images
  - Financial investments
  - … many more

## **Inductive learning**

- Simplest form: learn a function from examples
- •
- *f* is the target function
- An example is a pair (x, f(x))

```
    Pure induction task:
```

- Given a collection of examples of *f*, return a function *h* that approximates *f*.
- find a hypothesis h, such that  $h \approx f$ , given a training set of examples

(This is a highly simplified model of real learning)
 – Ignores prior knowledge

- Construct/adjust *h* to agree with *f* on training set
- (*h* is consistent if it agrees with *f* on all examples)
- •
- E.g., curve fitting:



- Construct/adjust *h* to agree with *f* on training set
- (*h* is consistent if it agrees with *f* on all examples)
- •
- E.g., curve fitting:



- Construct/adjust *h* to agree with *f* on training set
- (*h* is consistent if it agrees with *f* on all examples)
- •
- E.g., curve fitting:
  - f(x)

- Construct/adjust *h* to agree with *f* on training set
- (*h* is consistent if it agrees with *f* on all examples)
- •
- E.g., curve fitting:



- Construct/adjust h to agree with f on training set
- (*h* is consistent if it agrees with *f* on all examples)



- Construct/adjust *h* to agree with *f* on training set
- (*h* is consistent if it agrees with *f* on all examples)
- •



 Ockham's razor: prefer the simplest hypothesis consistent with data

## Generalization

- Hypotheses must generalize to correctly classify instances not in the training data.
- Simply memorizing training examples is a consistent hypothesis that does not generalize.
- Occam's razor:
  - Finding a *simple* hypothesis helps ensure generalization.

## **Training Error vs Test Error**



# **Supervised Learning**

- Learning a discrete function: Classification
  - Boolean classification:
    - Each example is classified as true(positive) or false(negative).
- Learning a continuous function: Regression

#### **Classification—A Two-Step Process**

- Model construction: describing a set of predetermined classes
  - Each tuple/sample is assumed to belong to a predefined class, as determined by the class label
  - The set of tuples used for model construction is training set
  - The model is represented as classification rules, decision trees, or mathematical formulae
- Model usage: for classifying future or unknown objects
  - Estimate accuracy of the model
    - The known label of test sample is compared with the classified result from the model
    - Test set is independent of training set, otherwise overfitting will occur
  - If the accuracy is acceptable, use the model to classify data tuples whose class labels are not known

20

## **Illustrating Classification Task**

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	?



## **Issues: Data Preparation**

- Data cleaning
  - Preprocess data in order to reduce noise and handle missing values
- Relevance analysis (feature selection)
  - Remove the irrelevant or redundant attributes
- Data transformation
  - Generalize data to (higher concepts, discretization)
  - Normalize attribute values

## **Classification Techniques**

- Decision Tree based Methods
- Rule-based Methods
- Naïve Bayes and Bayesian Belief Networks
- Neural Networks
- Support Vector Machines
- and more...

# Learning decision trees

- **Example Problem:** decide whether to wait for a table at a restaurant, based on the following attributes:
  - **1. Alternate:** is there an alternative restaurant nearby?
  - 2. Bar: is there a comfortable bar area to wait in?
  - 3. Fri/Sat: is today Friday or Saturday?
  - 4. Hungry: are we hungry?
  - 5. Patrons: number of people in the restaurant (None, Some, Full)
  - 6. Price: price range (\$, \$\$, \$\$\$)
  - **7. Raining**: is it raining outside?
  - **8. Reservation**: have we made a reservation?
  - **9. Type**: kind of restaurant (French, Italian, Thai, Burger)
  - 10. WaitEstimate: estimated waiting time (0-10, 10-30, 30-60, >60)

#### Feature(Attribute)-based representations

- Examples described by feature(attribute) values
  - (Boolean, discrete, continuous)
- E.g., situations where I will/won't wait for a table:

Example					At	tributes	;				Target
Linampro	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	Wait
$X_1$	Т	F	F	Т	Some	\$\$\$	F	Т	French	0–10	Т
$X_2$	Т	F	F	Т	Full	\$	F	F	Thai	30–60	F
$X_3$	F	Т	F	F	Some	\$	F	F	Burger	0–10	Т
$X_4$	Т	F	Т	Т	Full	\$	F	F	Thai	10–30	Т
$X_5$	Т	F	Т	F	Full	\$\$\$	F	Т	French	>60	F
$X_6$	F	Т	F	Т	Some	\$\$	Т	Т	Italian	0–10	Т
$X_7$	F	Т	F	F	None	\$	Т	F	Burger	0–10	F
$X_8$	F	F	F	Т	Some	\$\$	Т	Т	Thai	0–10	Т
$X_9$	F	Т	Т	F	Full	\$	Т	F	Burger	>60	F
$X_{10}$	Т	Т	Т	Т	Full	\$\$\$	F	Т	Italian	10–30	F
$X_{11}$	F	F	F	F	None	\$	F	F	Thai	0–10	F
$X_{12}$	Т	Т	Т	Т	Full	\$	F	F	Burger	30–60	Т

• Classification of examples is positive (T) or negative (F)

#### **Decision trees**

- One possible representation for hypotheses
- E.g., here is the "true" tree for deciding whether to wait:



## **Expressiveness**

- Decision trees can express any function of the input attributes.
- E.g., for Boolean functions, truth table row  $\rightarrow$  path to leaf:



- Trivially, there is a consistent decision tree for any training set with one path to leaf for each example (unless *f* nondeterministic in *x*) but it probably won't generalize to new examples
- Prefer to find more compact decision trees

## **Decision tree learning**

- Aim: find a small tree consistent with the training examples
- Idea: (recursively) choose "most significant" attribute as root of (sub)tree

```
function DTL(examples, attributes, default) returns a decision tree

if examples is empty then return default

else if all examples have the same classification then return the classification

else if attributes is empty then return MODE(examples)

else

best \leftarrow CHOOSE-ATTRIBUTE(attributes, examples)

tree \leftarrow a new decision tree with root test best

for each value v_i of best do

examples_i \leftarrow \{elements of examples with best = v_i\}

subtree \leftarrow DTL(examples_i, attributes - best, MODE(examples))

add a branch to tree with label v_i and subtree subtree

return tree
```

#### **Decision Tree Construction Algorithm**

#### Principle

- Basic algorithm (adopted by ID3, C4.5 and CART): a greedy algorithm
- Tree is constructed in a *top-down recursive divide-and-conquer* manner
- Iterations
  - At start, all the training tuples are at the root
  - Tuples are partitioned recursively based on selected attributes
  - Test attributes are selected on the basis of a heuristic or statistical measure (e.g, information gain)
- Stopping conditions
  - All samples for a given node belong to the same class
  - There are no remaining attributes for further partitioning –
     majority voting is employed for classifying the leaf
  - There are no samples left

#### **Decision Tree Induction: Training Dataset**

This follows an example of Quinlan's ID3 (Playing Tennis)

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
3140	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
3140	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
3140	medium	no	excellent	yes
3140	high	yes	fair	yes
>40	medium	no	excellent	no



RID	age	student	credit-rating	Class: buys_computer
1	youth	yes	fair	yes
2	youth	yes	fair	yes
3	youth	yes	fair	no
4	youth	no	fair	no
5	middle-aged	no	excellent	yes
6	senior	yes	fair	no
7	senior	yes	excellent	yes



RID	age	student	credit-rating	Class: buys_computer
1	youth	yes	fair	yes
2	youth	yes	fair	yes
3	youth	yes	fair	no
4	youth	no	fair	no
5	middle-aged	no	excellent	yes
6	senior	yes	fair	no
7	senior	yes	excellent	yes

32



RID	age	student	credit-rating	Class: buys_computer
1	youth	yes	fair	yes
2	youth	yes	fair	yes
3	youth	yes	fair	no
4	youth	no	fair	no
5	middle-aged	no	excellent	yes
6	senior	yes	fair	no
7	senior	yes	excellent	yes



RID	age	student	credit-rating	Class: buys_computer
1	youth	yes	fair	yes
2	youth	yes	fair	yes
3	youth	yes	fair	no
4	youth	no	fair	no
5	middle-aged	no	excellent	yes
6	senior	yes	fair	no
7	senior	yes	excellent	yes



RID	age	student	credit-rating	Class: buys_computer
1	youth	yes	fair	yes
2	youth	yes	fair	yes
3	youth	yes	fair	no
4	youth	no	fair	no
5	middle-aged	no	excellent	yes
6	senior	yes	fair	no
7	senior	yes	excellent	yes

35



RID	age	student	credit-rating	Class: buys_computer
1	youth	yes	fair	yes
2	youth	yes	fair	yes
3	youth	yes	fair	no
4	youth	no	fair	no
5	middle-aged	no	excellent	yes
6	senior	yes	fair	no
7	senior	yes	excellent	yes



RID	age	student	credit-rating	Class: buys_computer
1	youth	yes	fair	yes
2	youth	yes	fair	yes
3	youth	yes	fair	no
4	youth	no	fair	no
5	middle-aged	no	excellent	yes
6	senior	yes	fair	yes
7	senior	yes	excellent	no

## **Tree Induction**

- Greedy strategy.
  - Split the records based on an attribute test that optimizes 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

## **Choosing an attribute**

 Idea: a good attribute splits the examples into subsets that are (ideally) "all positive" or "all negative"



• *Patrons?* is a better choice

## How to determine the Best Split

- Greedy approach:
  - Nodes with homogeneous class distribution are preferred
- Need a measure of node impurity:

Non-homogeneous, High degree of impurity C0: 9 C1: 1

Homogeneous,

Low degree of impurity

## **Measures of Node Impurity**

Information Gain

• Gini Index

Misclassification error

Choose attributes to split to achieve minimum impurity

## **Attribute Selection Measure: Information Gain (ID3/C4.5)**

- Select the attribute with the highest information gain
- Let p<sub>i</sub> be the probability that an arbitrary tuple in D belongs to class C<sub>i</sub>, estimated by |C<sub>i, D</sub>|/|D|
- Expected information (entropy) needed to classify a tuple in D:  $I(D) = -\sum_{i=1}^{m} p_i \log_2(p_i)$
- Information needed (after using A to split D into v partitions) to classify D:  $Info_A(D) = \sum_{i=1}^{\nu} \frac{|D_j|}{|D|} \times I(D_j)$
- Information gained by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$

42

## **Information gain**

For the training set, p = n = 6, I(6/12, 6/12) = 1 bit

Consider the attributes *Patrons* and *Type* (and others too):

$$IG(Patrons) = 1 - \left[\frac{2}{12}I(0,1) + \frac{4}{12}I(1,0) + \frac{6}{12}I(\frac{2}{6},\frac{4}{6})\right] = .0541 \text{ bits}$$
$$IG(Type) = 1 - \left[\frac{2}{12}I(\frac{1}{2},\frac{1}{2}) + \frac{2}{12}I(\frac{1}{2},\frac{1}{2}) + \frac{4}{12}I(\frac{2}{4},\frac{2}{4}) + \frac{4}{12}I(\frac{2}{4},\frac{2}{4})\right] = 0 \text{ bits}$$

Patrons has the highest IG of all attributes and so is chosen by the DTL algorithm as the root



## **Example contd.**

• Decision tree learned from the 12 examples:



 Substantially simpler than "true" tree---a more complex hypothesis isn't justified by small amount of data

#### Measure of Impurity: GINI (CART, IBM IntelligentMiner)

• Gini Index for a given node t :

$$GINI(t) = 1 - \sum_{j} [p(j | t)]^{2}$$

(NOTE: p(j | t) is the relative frequency of class j at node t).

- Maximum  $(1 1/n_c)$  when records are equally distributed among all classes, implying least interesting information
- Minimum (0.0) when all records belong to one class, implying most interesting information

C1	0	C1	1	C1	2	C1	3
C2	6	C2	5	C2	4	C2	3
Gini=0.000		Gini=	0.278	Gini=	0.444	Gini=	0.500

# **Splitting Based on GINI**

- Used in CART, SLIQ, SPRINT.
- When a node p is split into k partitions (children), the quality of split is computed as,

$$GINI_{split} = \sum_{i=1}^{k} \frac{n_i}{n} GINI(i)$$

where,  $n_i =$  number of records at child i, n = number of records at node p.

#### **Comparison of Attribute Selection Methods**

- The three measures return good results but
  - Information gain:
    - biased towards multivalued attributes
  - Gain ratio:
    - tends to prefer unbalanced splits in which one partition is much smaller than the others

– Gini index:

- biased to multivalued attributes
- has difficulty when # of classes is large
- tends to favor tests that result in equal-sized partitions and purity in both partitions

47

# Example Algorithm: C4.5

- Simple depth-first construction.
- Uses Information Gain
- Sorts Continuous Attributes at each node.
- Needs entire data to fit in memory.
- Unsuitable for Large Datasets.
- You can download the software from Internet

## **Decision Tree Based Classification**

- Advantages:
  - Easy to construct/implement
  - Extremely fast at classifying unknown records
  - Models are easy to interpret for small-sized trees
  - Accuracy is comparable to other classification techniques for many simple data sets
  - Tree models make no assumptions about the distribution of the underlying data : nonparametric
  - Have a built-in feature selection method that makes them immune to the presence of useless variables

## **Decision Tree Based Classification**

- Disadvantages
  - Computationally expensive to train
  - Some decision trees can be overly complex that do not generalise the data well.
  - Less expressivity: There may be concepts that are hard to learn with limited decision trees

# **Overfitting and Tree Pruning**

- Overfitting: An induced tree may overfit the training data
  - Too many branches, some may reflect anomalies due to noise or outliers
  - Poor accuracy for unseen samples
- Two approaches to avoid overfitting
  - Prepruning: Halt tree construction early—do not split a node if this would result in the goodness measure 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"

51

## **Rule-Based Classifier**

- Classify records by using a collection of "if...then..." rules
- Rule: (Condition)  $\rightarrow y$ 
  - where
    - *Condition* is a conjunctions of attributes
    - y is the class label
  - LHS: rule antecedent or condition
  - *RHS*: rule consequent
  - Examples of classification rules:
    - (Blood Type=Warm)  $\land$  (Lay Eggs=Yes)  $\rightarrow$  Birds
    - (Taxable Income < 50K)  $\land$  (Refund=Yes)  $\rightarrow$  Evade=Ng<sub>2</sub>

## **Rule-based Classifier (Example)**

Name	Blood Type	Give Birth	Can Fly	Live in Water	Class
human	warm	yes	no	no	mammals
python	cold	no	no	no	reptiles
salmon	cold	no	no	yes	fishes
whale	warm	yes	no	yes	mammals
frog	cold	no	no	sometimes	amphibians
komodo	cold	no	no	no	reptiles
bat	warm	yes	yes	no	mammals
pigeon	warm	no	yes	no	birds
cat	warm	yes	no	no	mammals
leopard shark	cold	yes	no	yes	fishes
turtle	cold	no	no	sometimes	reptiles
penguin	warm	no	no	sometimes	birds
porcupine	warm	yes	no	no	mammals
eel	cold	no	no	yes	fishes
salamander	cold	no	no	sometimes	amphibians
gila monster	cold	no	no	no	reptiles
platypus	warm	no	no	no	mammals
owl	warm	no	yes	no	birds
dolphin	warm	yes	no	yes	mammals
eagle	warm	no	yes	no	birds

R1: (Give Birth = no)  $\land$  (Can Fly = yes)  $\rightarrow$  Birds

R2: (Give Birth = no)  $\land$  (Live in Water = yes)  $\rightarrow$  Fishes

R3: (Give Birth = yes)  $\land$  (Blood Type = warm)  $\rightarrow$  Mammals

R4: (Give Birth = no)  $\land$  (Can Fly = no)  $\rightarrow$  Reptiles

R5: (Live in Water = sometimes)  $\rightarrow$  Amphibians

#### **Rule Extraction from a Decision Tree**

- Rules are easier to understand than large trees
- 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 holds the class prediction





54

#### **Extra Slides**

## **Learning agents**

#### Performance standard



# Classification(Sınıflandırma)

- IDEA: Build a model based on past data to predict the class of the new data
- Given a collection of records (*training set*)
  - Each record contains a set of *attributes*, one of the attributes is the *class*.
- Find a *model* for class attribute as a function of the values of other attributes.
- **Goal:** <u>previously unseen</u> records should be assigned a class as accurately as possible.

## Hypothesis spaces

How many distinct decision trees with *n* Boolean attributes? = number of Boolean functions = number of distinct truth tables with 2<sup>n</sup> rows = 2<sup>2<sup>n</sup></sup>

• E.g., with 6 Boolean attributes, there are 18,446,744,073,709,551,616 trees

How many purely conjunctive hypotheses (e.g., *Hungry*  $\land \neg Rain$ )?

- Each attribute can be in (positive), in (negative), or out
  - $\Rightarrow$  3<sup>n</sup> distinct conjunctive hypotheses
- More expressive hypothesis space
  - increases chance that target function can be expressed
  - increases number of hypotheses consistent with training set
    - $\Rightarrow$  may get worse predictions

# **Using information theory**

- To implement Choose-Attribute in the DTL algorithm
- Information Content (Entropy):

$$I(P(v_1), ..., P(v_n)) = \Sigma_{i=1} - P(v_i) \log_2 P(v_i)$$

 For a training set containing p positive examples and n negative examples:

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

## **Information gain**

 A chosen attribute A divides the training set E into subsets E<sub>1</sub>, ..., E<sub>v</sub> according to their values for A, where A has v distinct values.

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

 Information Gain (IG) or reduction in entropy from the attribute test:

$$IG(A) = I(\frac{p}{p+n}, \frac{n}{p+n}) - remainder(A)$$

• Choose the attribute with the largest IG

#### **Performance measurement**

- How do we know that  $h \approx f$ ?
  - 1. Use theorems of computational/statistical learning theory
  - 2. Try *h* on a new test set of examples

(use same distribution over example space as training set)

Learning curve = % correct on test set as a function of training set size

