



ENEE459M: Topics in Computer Engineering: Machine Learning and Data Mining

- Course Objectives:
 - ♦ Introduction of basic machine learning and data mining techniques.
 - ♦ Get a practical feel through the Weka software environment.
- Will more or less follow textbook
- Weka Software: download Weka3.6 (see syllabus)
- Course Grade:
 - ♦ Three midterms (20% each)
 - ♦ Final (30%)
 - ♦ Homework (10%)



Data Mining

Practical Machine Learning Tools and Techniques

I. H. Witten and E. Frank

Chapter 1



What's it all about?

- Data vs information
- Data mining and machine learning
- Structural descriptions
 - ♦ Rules: classification and association
 - ♦ Decision trees
- Datasets
 - ♦ Weather, contact lens, CPU performance, labor negotiation data, soybean classification
- Fielded applications
 - ♦ Loan applications, screening images, load forecasting, machine fault diagnosis, market basket analysis
- Generalization as search
- Data mining and ethics



Examples: Classification

- Example 1: Spam Detection
 - ♦ Given: an email message X
 - ♦ Problem: Determine whether X is a spam
 - ♦ Data: A labeled set of email messages
- Example 2: Hand-Written Character Recognition
 - ♦ Given: hand-written character
 - ♦ Problem: determine the corresponding letter
 - ♦ Data: a set of training examples



Examples: Regression

- Given a customer (annual income, family size, location, job, ...), determine the annual credit card spending based on a sample of customers data.
- Same type as classification except the output is a numerical value.
- General setup: Given a set of training examples (\mathbf{x} , $f(\mathbf{x})$), determine an approximation of the function f .



Examples: Clustering

- Group a large collection of Web documents into a small number of groups, each of which contains “similar” documents.
- Cluster a market into distinct clusters, where a cluster represents a group of similar customers to be targeted with a distinct marketing mix.



Machine learning techniques

- *Algorithms for acquiring structural descriptions from examples – Supervised Learning*
- Structural descriptions represent patterns explicitly
 - Can be used to predict outcome in new situation
 - Can be used to understand and explain how prediction is derived (*may be even more important*)
- Methods originate from artificial intelligence, statistics, and research on databases



Structural descriptions

- Example: if-then rules

```
If tear production rate = reduced
then recommendation = none
Otherwise, if age = young and astigmatic = no
then recommendation = soft
```



Age	Spectacle prescription	Astigmatism	Tear production rate	Recommended lenses
Young	Myope	No	Reduced	None
Young	Hypermetrope	No	Normal	Soft
Pre-presbyopic Presbyopic	Hypermetrope	No	Reduced	None
	Myope	Yes	Normal	Hard
...

The weather problem

• Conditions for playing a certain game

Outlook	Temperature	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	Normal	False	Yes
...

```

If outlook = sunny and humidity = high then play = no
If outlook = rainy and windy = true then play = no
If outlook = overcast then play = yes
If humidity = normal then play = yes
If none of the above then play = yes

```

Classification vs. association rules

• Classification rule:

predicts value of a given attribute (the classification of an example)

```

If outlook = sunny and humidity = high
then play = no

```

• Association rule:

predicts value of arbitrary attribute (or combination)

```

If temperature = cool then humidity = normal
If humidity = normal and windy = false
then play = yes
If outlook = sunny and play = no
then humidity = high
If windy = false and play = no
then outlook = sunny and humidity = high

```

Weather data with mixed attributes

• Some attributes have numeric values

Outlook	Temperature	Humidity	Windy	Play
Sunny	85	85	False	No
Sunny	80	90	True	No
Overcast	83	86	False	Yes
Rainy	75	80	False	Yes
...

```

If outlook = sunny and humidity > 83 then play = no
If outlook = rainy and windy = true then play = no
If outlook = overcast then play = yes
If humidity < 85 then play = yes
If none of the above then play = yes

```

The contact lenses data

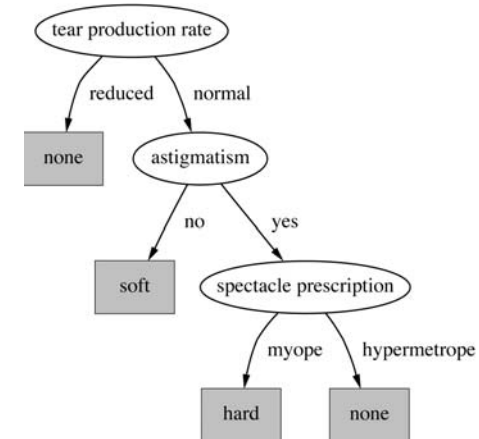
Age	Spectacle prescription	Astigmatism	Tear production rate	Recommended lenses
Young	Myope	No	Reduced	None
Young	Myope	No	Normal	Soft
Young	Myope	Yes	Reduced	None
Young	Myope	Yes	Normal	Hard
Young	Hypermetrope	No	Reduced	None
Young	Hypermetrope	No	Normal	Soft
Young	Hypermetrope	Yes	Reduced	None
Young	Hypermetrope	Yes	Normal	hard
Pre-presbyopic	Myope	No	Reduced	None
Pre-presbyopic	Myope	No	Normal	Soft
Pre-presbyopic	Myope	Yes	Reduced	None
Pre-presbyopic	Myope	Yes	Normal	Hard
Pre-presbyopic	Hypermetrope	No	Reduced	None
Pre-presbyopic	Hypermetrope	No	Normal	Soft
Pre-presbyopic	Hypermetrope	Yes	Reduced	None
Pre-presbyopic	Hypermetrope	Yes	Normal	None
Presbyopic	Myope	No	Reduced	None
Presbyopic	Myope	No	Normal	None
Presbyopic	Myope	Yes	Reduced	None
Presbyopic	Myope	Yes	Normal	Hard
Presbyopic	Hypermetrope	No	Reduced	None
Presbyopic	Hypermetrope	No	Normal	Soft
Presbyopic	Hypermetrope	Yes	Reduced	None
Presbyopic	Hypermetrope	Yes	Normal	None

A complete and correct rule set

```

If tear production rate = reduced then recommendation = none
If age = young and astigmatic = no
  and tear production rate = normal then recommendation = soft
If age = pre-presbyopic and astigmatic = no
  and tear production rate = normal then recommendation = soft
If age = presbyopic and spectacle prescription = myope
  and astigmatic = no then recommendation = none
If spectacle prescription = hypermetrope and astigmatic = no
  and tear production rate = normal then recommendation = soft
If spectacle prescription = myope and astigmatic = yes
  and tear production rate = normal then recommendation = hard
If age young and astigmatic = yes
  and tear production rate = normal then recommendation = hard
If age = pre-presbyopic
  and spectacle prescription = hypermetrope
  and astigmatic = yes then recommendation = none
If age = presbyopic and spectacle prescription = hypermetrope
  and astigmatic = yes then recommendation = none
  
```

A decision tree for this problem



Classifying iris flowers



	Sepal length	Sepal width	Petal length	Petal width	Type
1	5.1	3.5	1.4	0.2	Iris setosa
2	4.9	3.0	1.4	0.2	Iris setosa
...					
51	7.0	3.2	4.7	1.4	Iris versicolor
52	6.4	3.2	4.5	1.5	Iris versicolor
...					
101	6.3	3.3	6.0	2.5	Iris virginica
102	5.8	2.7	5.1	1.9	Iris virginica
...					

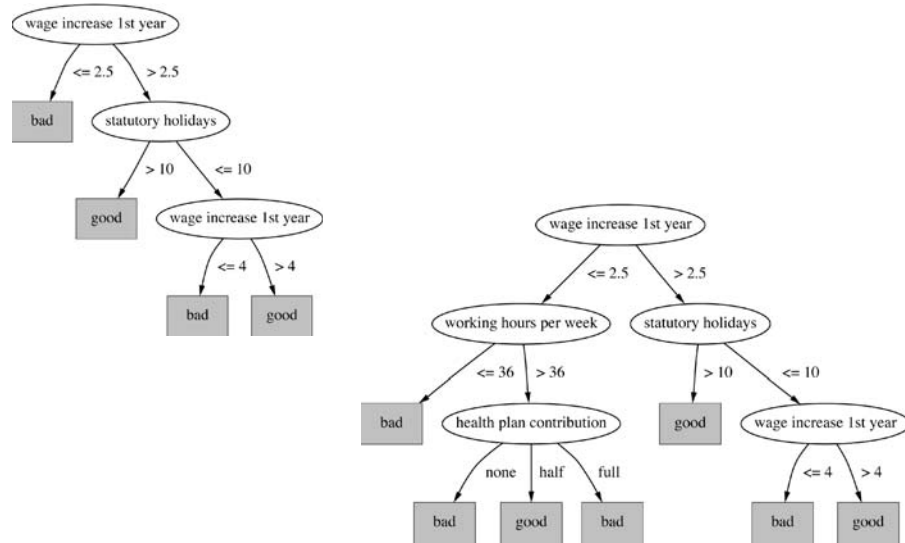
```

If petal length < 2.45 then Iris setosa
If sepal width < 2.10 then Iris versicolor
...
  
```

Data from labor negotiations

Attribute	Type	1	2	3	...	40
Duration	(Number of years)	1	2	3	...	2
Wage increase first year	Percentage	2%	4%	4.3%	...	4.5
Wage increase second year	Percentage	?	5%	4.4%	...	4.0
Wage increase third year	Percentage	?	?	?	...	?
Cost of living adjustment	{none,tcf,tc}	none	tcf	?	...	none
Working hours per week	(Number of hours)	28	35	38	...	40
Pension	{none,ret-allw, empl- pact}	none	?	?	...	?
Standby pay	Percentage	?	13%	?	...	?
Shift-work supplement	Percentage	?	5%	4%	...	4
Education allowance	{yes,no}	yes	?	?	...	?
Statutory holidays	(Number of days)	11	15	12	...	12
Vacation	{below-avg,avg,gen}	avg	gen	gen	...	avg
Long-term disability	{yes,no}	no	?	?	...	yes
Disability plan contribution	{none,half,full}	none	?	full	...	full
Bereavement assistance	{yes,no}	no	?	?	...	yes
Health plan contribution	{none,half,full}	none	?	full	...	half
Acceptability of contract	{good,bad}	bad	good	good	...	good

Decision trees for the labor data



Data Mining: Practical Machine Learning Tools and Techniques (Chapter 1)

17

Fielded applications

- The result of learning—or the learning method itself—is deployed in practical applications
 - Processing loan applications
 - Screening images for oil slicks
 - Electricity supply forecasting
 - Diagnosis of machine faults
 - Marketing and sales
 - Separating crude oil and natural gas
 - Reducing banding in rotogravure printing
 - Finding appropriate technicians for telephone faults
 - Scientific applications: biology, astronomy, chemistry
 - Automatic selection of TV programs
 - Monitoring intensive care patients

Data Mining: Practical Machine Learning Tools and Techniques (Chapter 1)

18

Processing loan applications (American Express)

- Given: questionnaire with financial and personal information
- Question: should money be lent?
- Simple statistical method covers 90% of cases
- Borderline cases referred to loan officers
- But: 50% of accepted borderline cases defaulted!
- Solution: reject all borderline cases?
 - No! Borderline cases are most active customers

Data Mining: Practical Machine Learning Tools and Techniques (Chapter 1)

19

Enter machine learning

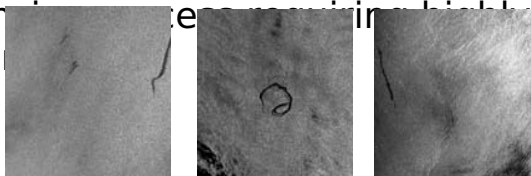
- 1000 training examples of borderline cases
- 20 attributes:
 - age
 - years with current employer
 - years at current address
 - years with the bank
 - other credit cards possessed,...
- Learned rules: correct on 70% of cases
 - human experts only 50%
- Rules could be used to explain decisions to customers

Data Mining: Practical Machine Learning Tools and Techniques (Chapter 1)

20

Screening images

- Given: radar satellite images of coastal waters
- Problem: detect oil slicks in those images
- Oil slicks appear as dark regions with changing size and shape
- Not easy: lookalike dark regions can be caused by weather conditions (e.g. high wind)
- Expertise in screening requires highly trained personnel



Enter machine learning

- Extract dark regions from normalized image
- Attributes:
 - ♦ size of region
 - ♦ shape, area
 - ♦ intensity
 - ♦ sharpness and jaggedness of boundaries
 - ♦ proximity of other regions
 - ♦ info about background
- Constraints:
 - ♦ Few training examples—oil slicks are rare!
 - ♦ Unbalanced data: most dark regions aren't slicks
 - ♦ Regions from same image form a batch
 - ♦ Requirement: adjustable false-alarm rate

Load forecasting

- Electricity supply companies need forecast of future demand for power
- Forecasts of min/max load for each hour ⇒ significant savings
- Given: manually constructed load model that assumes “normal” climatic conditions
- Problem: adjust for weather conditions
- Static model consist of:
 - ♦ base load for the year
 - ♦ load periodicity over the year
 - ♦ effect of holidays



Enter machine learning

- Prediction corrected using “most similar” days
- Attributes:
 - ♦ temperature
 - ♦ humidity
 - ♦ wind speed
 - ♦ cloud cover readings
 - ♦ plus difference between actual load and predicted load
- Average difference among three “most similar” days added to static model
- Linear regression coefficients form attribute weights in similarity function

Marketing and sales I

- Companies precisely record massive amounts of marketing and sales data
- Applications:
 - ♦ Customer loyalty: identifying customers that are likely to defect by detecting changes in their behavior (e.g. banks/phone companies)
 - ♦ Special offers: identifying profitable customers (e.g. reliable owners of credit cards that need extra money during the holiday season)

Marketing and sales II

- Market basket analysis
 - ♦ Association techniques find groups of items that tend to occur together in a transaction (used to analyze checkout data)
- Historical analysis of purchasing patterns
- Identifying prospective customers
 - ♦ Focusing promotional mailouts (targeted campaigns are cheaper than mass-marketed ones)



Machine learning and statistics

- Historical difference (grossly oversimplified):
 - ♦ Statistics: testing hypotheses
 - ♦ Machine learning: finding the right hypothesis
- But: huge overlap
 - ♦ Decision trees (C4.5 and CART)
 - ♦ Nearest-neighbor methods
- Today: perspectives have converged
 - ♦ Most ML algorithms employ statistical techniques

Generalization as search

- Inductive learning: find a concept description that fits the data
- Example: rule sets as description language
 - ♦ Enormous, but finite, search space
- Simple solution:
 - ♦ enumerate the concept space
 - ♦ eliminate descriptions that do not fit examples
 - ♦ surviving descriptions contain target concept

Data mining and ethics I



- Ethical issues arise in practical applications
- Data mining often used to discriminate
 - ♦ E.g. loan applications: using some information (e.g. sex, religion, race) is unethical
- Ethical situation depends on application
 - ♦ E.g. same information ok in medical application
- Attributes may contain problematic information
 - ♦ E.g. area code may correlate with race

Data mining and ethics II

- Important questions:
 - ♦ Who is permitted access to the data?
 - ♦ For what purpose was the data collected?
 - ♦ What kind of conclusions can be legitimately drawn from it?
- Caveats must be attached to results
- Purely statistical arguments are never sufficient!
- Are resources put to good use?