26:198:722 Expert Systems

- Knowledge representation
- Knowledge acquisition
- Machine learning
- ID3 & C4.5

Recall:

- Knowledge engineering
 - * Knowledge acquisition
 - Knowledge elicitation
 - * Knowledge representation
 - Production rules
 - Semantic networks
 - Frames

- Representation is more than just encoding (encrypting)
- Coding preserves structural ambiguity
- Communication assumes prior knowledge
- Representation implies organization

- Representation
 - * A set of syntactic and semantic conventions that make it possible to describe things (Winston)
- Description
 - * makes use of the conventions of a representation to describe some particular thing
- Syntax v. semantics

STRIPS

- Predicate-argument expressions
 - at (robot, roomA)
- * World models
- * Operator tables
 - ◆ push (X, Y, Z)
 - → Preconditions
 - → Delete list
 - → Add list

STRIPS

- * maintained lists of goals
- * selected goal to work on next
- * searched for applicable operators
- * matched goals against formulas in add lists
- * set up preconditions as sub-goals
- * used means-end analysis

- STRIPS lessons
 - * Heuristic search
 - Uniform representation
 - * Problem reduction

Procedural semantics

MYCIN

- * Assists physicians who are not experts in the field of antibiotics in treating blood infections
- Consists of
 - Knowledge base
 - Dynamic patient database
 - Consultation program
 - Explanation program
 - Knowledge acquisition program

MYCIN

- * Production rules
 - Premises
 - Conjunctions of conditions
 - Actions
 - Conclusions or instructions
- * Patient information stored in context tree
- Certainty factors for uncertain reasoning
- * Backward chaining control structure (based on AND/OR tree)

MYCIN

* Evaluation

- Panel of experts approved 72% of recommendations
- Good as experts
- Better than non-experts
- Knowledge base incomplete (400 rules)
- Required more computing power than available in hospitals
- Doctors did not like the user interface

- Stages
 - * Identification
 - * Conceptualization
 - * Formalization
 - * Implementation
 - * Testing
- KADS
- Ontological analysis

- Expert system shells
 - * EMYCIN
 - * TEIRESIAS
 - Rule models (meta-rules)
 - Schemas for data types
 - Domain-specific knowledge
 - Representation-specific knowledge
 - Representation-independent knowledge
 - Explain-Test-Review

Methods and tools

- * Structured interview
- * Unstructured interview
- * Case studies
 - Retrospective v. observational
 - Familiar v. unfamiliar
- * Concurrent protocols
 - Verbalization, "thinking aloud"
- * Tape recording
- * Video recording

- Methods and tools
 - * Automated knowledge acquisition
 - Domain models
 - Graphical interfaces
 - Visual programming language

Different types of knowledge

- Procedural knowledge
 - Rules, strategies, agendas, procedures
- * Declarative knowledge
 - Concepts, objects, facts
- * Meta-knowledge
 - Knowledge about other types of knowledge and how to use them
- * Structural knowledge
 - Rules sets, concept relationships, concept to object relationships

- Sources of knowledge
 - * Experts
 - * End-users
 - * Multiple experts (panels)
 - * Reports
 - * Books
 - * Regulations
 - * Guidelines

- Major difficulties with elicitation
 - * Expert may
 - be unaware of the knowledge used
 - be unable to verbalize the knowledge used
 - provide irrelevant knowledge
 - provide incomplete knowledge
 - provide incorrect knowledge
 - provide inconsistent knowledge

"The more competent domain experts become, the less able they are to describe the knowledge they used to solve problems" (Waterman)

Detailed guidelines for conducting structured and unstructured interviews and both retrospective and observational case studies are given in Durkin (Chapter 17)

Technique Capabilities

	<u>Interv</u>	<u>views</u>	<u>Case Studies</u>			
			Retrospective		Observational	
Knowledge	Unstructured	Structured	Familiar	Unfamiliar	Familiar	<u>Unfamiliar</u>
Facts	Poor	Good	Fair	Average	Good	Excellent
Concepts	Excellent	Excellent	Average	Average	Good	Good
Objects	Good	Excellent	Average	Average	Good	Good
Rules	Fair	Average	Average	Average	Good	Excellent
Strategies	Average	Average	Good	Good	Excellent	Excellent
Heuristics	Fair	Average	Excellent	Good	Good	Poor
Structures	Fair	Excellent	Average	Average	Average	Average

- Analyzing the knowledge collected
 - Producing transcripts
 - Interpreting transcripts
 - Chunking
 - * Analyzing transcripts
 - Knowledge dictionaries
 - Graphical techniques
 - Cognitive maps
 - → Inference networks
 - → Flowcharts
 - Decision trees

- Rote learning
- Supervised learning
 - * Induction
 - Concept learning
 - Descriptive generalization
- Unsupervised learning

META-DENDRAL

* RULEMOD

- Removing redundancy
- Merging rules
- Making rules more specific
- Making rules more general
- Selecting final rules

META-DENDRAL

- Version spaces
 - Partial ordering
 - Boundary sets
 - Candidate elimination algorithm
 - Monotonic, non-heuristic
 - Results independent of order of presentation
 - Each training instance examine only once
 - Discarded hypotheses never reconsidered
 - Learning is properly incremental

Decision trees and production rules

- Decision trees are an alternative way of structuring rules
- Efficient algorithms exist for constructing decision trees
- There is a whole family of such learning systems:
 - ◆ CLS (1966)
 - ◆ ID3 (1979)
 - ◆ ACLS (1982)
 - ◆ ASSISTANT (1984)
 - ◆ IND (1990)
 - ◆ C4.5 (1993) and C5.0
- Decision trees can be converted to rules later

Entropy

- * Let X be a variable with states x_1 - x_n
- * Define the entropy of X by

$$H(X) = -\sum_{i=1}^{n} p(x_i) \log_2(p(x_i))$$

* N.B.
$$\log_2(x) = \frac{\log_{10}(x)}{\log_{10}(2)} = \frac{\ln(x)}{\ln(2)}$$

Entropy

* Consider flipping a perfect coin:

e.g.,
$$n = 2$$

$$X: x_1, x_2$$

$$p(x_1) = p(x_2) = 1/2$$

Entropy

$$H(X) = -\sum_{i=1}^{n} p(x_i) \log_2(p(x_i))$$

$$= -\left[\frac{1}{2}\log_2(\frac{1}{2}) + \frac{1}{2}\log_2(\frac{1}{2})\right]$$

$$= -\left[\frac{1}{2}(-1) + \frac{1}{2}(-1)\right] = 1$$

Entropy

* Consider n equiprobable outcomes

$$H(X) = -\sum_{i=1}^{n} p(x_i) \log_2(p(x_i))$$

$$= -\sum_{i=1}^{n} \frac{1}{n} \log_2(\frac{1}{n})$$

$$= \sum_{i=1}^{n} \frac{1}{n} \log_2(n) = \log_2(n)$$

Entropy

* Consider flipping a totally biased coin:

e.g.,
$$n = 2$$

$$X: x_1, x_2$$

$$p(x_1) = 1$$
 $p(x_2) = 0$

Entropy

$$H(X) = -\sum_{i=1}^{n} p(x_i) \log_2(p(x_i))$$

$$= -\left[\log_2(1) + 0\log_2(0)\right]$$

$$= -\left[0 + 0\log_2(0)\right] = 0$$

(by L'Hopital's rule)

- Entropy
 - *Entropy is a measure of chaos or disorder
 - * H(X) is maximum for equiprobable outcomes

Entropy

* X: x_1 - - - x_m and Y: y_1 - - - y_n be two variables

$$H(X,Y) = -\sum_{i=1}^{m} \sum_{j=1}^{n} p(x_{i}, y_{j}) \log_{2}(p(x_{i}, y_{j}))$$

* If X and Y are independent

$$H(X,Y) = H(X) + H(Y)$$

- Conditional Entropy
 - * Partial conditional entropy of Y given X is in state x_i :

$$H(Y|x_i) = -\sum_{j=1}^{n} p(y_j|x_i) \log_2(p(y_j|x_i))$$

* Full conditional entropy of Y given X

$$H(Y|X) = \sum_{i=1}^{m} p(x_i) \cdot H(Y|x_i)$$

Binary Logarithms

```
1 0.0000
```

- 2 1.0000
- 3 1.5850
- 4 2.0000
- 5 2.3219
- 6 2.5850
- 7 2.8074
- 8 3.0000

ID3

- * Builds a decision tree first, then rules
- * Given a set of attributes, and a decision, recursively selects attributes to be the root of the tree based on Information Gain:
 - H(decision) H(decision | attribute)
- * Favors attributes with many outcomes
- * Is not guaranteed to find the simplest decision tree
- * Is not incremental

- **C4.5**
 - * Selects attributes based on Information gain ratio:
 - (H(decision) H(decision | attribute)) / H(attribute)
 - * Uses pruning heuristics to simplify decision trees
 - to simplify
 - to reduce dependence on training set
 - * Tunes the resulting rule(s)

C4.5 rule tuning

- * Derive initial rules by enumerating paths through the decision tree
- Generalize the rules by possibly deleting unnecessary conditions
- * Group rules according to target classes and delete any that do not contribute to overall performance on the class
- * Order the sets of rules for the target classes and choose a default class

- Rule tuning
 - * Rule tuning may be useful for rules derived by a variety of other means besides C4.5
 - Evaluate the contribution of individual rules
 - Evaluate the performance of the rule set as a whole

A data set for classification (Quinlan)

		Decision		
	<u>Height</u>	<u>Hair</u>	<u>Eyes</u>	<u>Attractiveness</u>
1	short	blond	blue	+
2	tall	blond	brown	-
3	tall	red	blue	+
4	short	dark	blue	-
5	tall	dark	blue	-
6	tall	blond	blue	+
7	tall	dark	brown	-
8	short	blond	brown	-

- A data set for classification (Quinlan)
 - * H(decision) = H(Attractiveness)

$$= -\frac{3}{8}\log_2\left(\frac{3}{8}\right) - \frac{5}{8}\log_2\left(\frac{5}{8}\right) = 0.955$$

- A data set for classification (Quinlan)
 - * Height:

◆ short: 1, 4, 8

p(+|short) = 1/3 p(-|short) = 2/3

◆ tall: 2, 3, 5, 6, 7

p(+|tall) = 2/5 p(-|tall) = 3/5

* H(decision|attribute) = H(Attractiveness|Height) =

$$\frac{3}{8} \left[-\frac{1}{3} \log_2 \left(\frac{1}{3} \right) - \frac{2}{3} \log_2 \left(\frac{2}{3} \right) \right] + \frac{5}{8} \left[-\frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right] = 0.951$$

* Information gain = 0.955 - 0.951 = 0.004

- A data set for classification (Quinlan)
 - * Hair:

• blond: 1, 2, 6, 8
$$p(+|blond) = 2/4$$
 $p(-|blond) = 2/4$

• red: 3

$$p(+|red) = 1/1$$
 $p(-|red) = 0/1$

$$p(-||eu|) = 0/1$$

dark: 4, 5, 7

$$p(+|dark) = 0/3$$
 $p(-|dark) = 3/3$

$$p(-|dark) = 3/3$$

* H(decision|attribute) = H(Attractiveness|Hair) =

$$\frac{4}{8}[1] + \frac{1}{8}[0] + \frac{3}{8}[0] = 0.500$$

* Information gain = 0.955 - 0.500 = 0.455

- A data set for classification (Quinlan)
 - * Eyes:
 - blue: 1, 3, 4, 5, 6 p(+|b|ue) = 3/5 p(-|b|ue) = 2/5
 - brown: 2, 7, 8 p(+|brown) = 0/3 p(-|brown) = 3/3
 - * H(decision|attribute) = H(Attractiveness|Eyes) =

$$\frac{5}{8} \left[-\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{5} \right) \right] + \frac{3}{8} [0] = 0.607$$

* Information gain = 0.955 - 0.607 = 0.348

- A data set for classification (Quinlan)
 - * Hence Hair is chosen as the best choice for the root of the tree
 - Now we recursively repeat this process for the (three) resulting branches
 - * In this case, the branches for Hair: red and Hair: dark are already completely classified, and we need to work only on the sub-table for Hair: blond

A data set for classification (Quinlan)

-----Attributes-----
Height Eyes ----- Attractiveness

1 short blue +

2 tall brown -

6 tall blue +

8 short brown -

* H(decision) = H(Attractiveness)

$$= -\frac{2}{4}\log_2\left(\frac{2}{4}\right) - \frac{2}{4}\log_2\left(\frac{2}{4}\right) = 1$$

- A data set for classification (Quinlan)
 - * Height:

◆ short: 1, 8

p(+|short) = 1/2 p(-|short) = 1/2

◆ tall: 2, 6

p(+|tall) = 1/2 p(-|tall) = 1/2

* H(decision|attribute) = H(Attractiveness|Height) =

$$\frac{2}{4}[1] + \frac{2}{4}[1] = 1$$

★ Information gain = 1-1=0

- A data set for classification (Quinlan)
 - * Eyes:

♦ blue: 1, 6

p(+|b|ue) = 2/2 p(-|b|ue) = 0/2

◆ brown: 2, 8

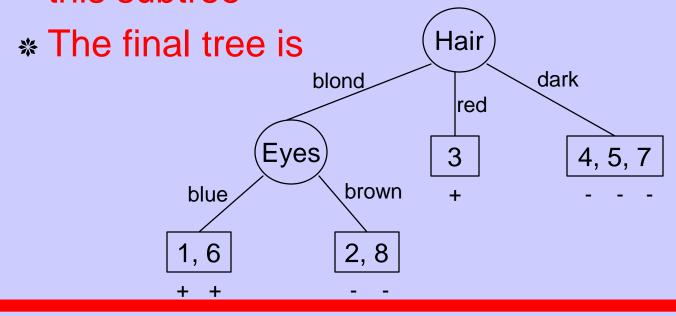
p(+|brown) = 0/2 p(-|brown) = 2/2

* H(decision|attribute) = H(Attractiveness|Eyes) =

$$\frac{2}{4}[0] + \frac{2}{4}[0] = 0$$

★ Information gain = 1-0=1

- A data set for classification (Quinlan)
 - * Hence Eyes is chosen as the best root of this subtree



- A data set for classification (Quinlan)
 - * We may now build rules from this decision tree
 - ◆ R1: (Hair, dark) --> (Attractiveness, -)
 - ◆ R2: (Hair, red) --> (Attractiveness, +)
 - ◆ R3: (Hair, blond) & (Eyes, blue) --> (Attractiveness, +)
 - ◆ R4: (Hair, blond) & (Eyes, brown) --> (Attractiveness, -)
 - Note that height is irrelevant

- A data set for classification (Quinlan)
 - * Dropping conditions from rules
 - Rules 1 and 2 have only one condition
 - Rule 3: neither condition can be dropped (case 5 needs the first condition and case 2 needs the second condition)
 - Rule 4: we can drop the first condition
 - ◆ R4': (Eyes, brown) --> (Attractiveness, -)

- A data set for classification (Quinlan)
 - Dropping conditions from rules
 - Linear
 - → Scan rule left to right
 - → Try to drop conditions one at a time
 - → If possible, drop for good
 - → Iterate (n conditions, n attempts)
 - Exponential
 - → Scan rule left to right
 - → Try to drop conditions one at a time
 - → Then try to drop pairs, triples, etc. (n conditions, 2ⁿ-2 attempts)

- A data set for classification (Quinlan)
 - Now consider Information gain ratio
 - * For initial root of tree we already know
 - * H(decision) = H(Attractiveness)

$$= -\frac{3}{8}\log_2\left(\frac{3}{8}\right) - \frac{5}{8}\log_2\left(\frac{5}{8}\right) = 0.955$$

- A data set for classification (Quinlan)
 - * H(decision|attribute) = H(Attractiveness|Height) =

$$\frac{3}{8} \left[-\frac{1}{3} \log_2 \left(\frac{1}{3} \right) - \frac{2}{3} \log_2 \left(\frac{2}{3} \right) \right] + \frac{5}{8} \left[-\frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right] = 0.951$$

* H(decision|attribute) = H(Attractiveness|Hair) =

$$\frac{4}{8}[1] + \frac{1}{8}[0] + \frac{3}{8}[0] = 0.500$$

* H(decision|attribute) = H(Attractiveness|Eyes) =

$$\frac{5}{8} \left[-\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{5} \right) \right] + \frac{3}{8} [0] = 0.607$$

A data set for classification (Quinlan)

$$-\frac{3}{8}\log_2\left(\frac{3}{8}\right) - \frac{5}{8}\log_2\left(\frac{5}{8}\right) = 0.955$$

* H(attribute) = H(Hair) =

$$-\frac{4}{8}\log_2\left(\frac{4}{8}\right) - \frac{1}{8}\log_2\left(\frac{1}{8}\right) - \frac{3}{8}\log_2\left(\frac{3}{8}\right) = 1.406$$

* H(attribute) = H(Eyes) =

$$-\frac{5}{8}\log_2\left(\frac{5}{8}\right) - \frac{3}{8}\log_2\left(\frac{3}{8}\right) = 0.955$$

- A data set for classification (Quinlan)
 - * Hence the Information gain ratios are

◆ Height: 0.004

◆ Hair: 0.324

◆ Eyes: 0.364

- * By this criterion, Eyes is chosen as the best root available
- * The branch for Eyes: brown is already completely classified, and we need to work only on the sub-table for Eyes: blue

A data set for classification (Quinlan)

	Attribut	es	Decision
	<u>Height</u>	<u>Hair</u>	<u>Attractiveness</u>
1	short	blond	+
3	tall	red	+
4	short	dark	-
5	tall	dark	-
6	tall	blond	+
			_

* H(decision) = H(Attractiveness)

$$= -\frac{3}{5}\log_2\left(\frac{3}{5}\right) - \frac{2}{5}\log_2\left(\frac{2}{5}\right) = 0.971$$

- A data set for classification (Quinlan)
 - * Height:

◆ short: 1, 4

p(+|short) = 1/2 p(-|short) = 1/2

◆ tall: 3, 5, 6

p(+|tall) = 2/3 p(-|tall) = 1/3

* H(decision|attribute) = H(Attractiveness|Height) =

$$\frac{2}{5} \left[-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right] + \frac{3}{5} \left[-\frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right] = 0.951$$

* H(Height) =
$$-\frac{2}{5}\log_2\left(\frac{2}{5}\right) - \frac{3}{5}\log_2\left(\frac{3}{5}\right) = 0.971$$

- A data set for classification (Quinlan)
 - * Hair:

$$p(+|blond) = 2/2$$
 $p(-|blond) = 0/2$

$$p(+|red) = 1/1$$
 $p(-|red) = 0/1$

$$p(-|red) = 0/1$$

$$p(+|dark) = 0/2$$

$$p(+|dark) = 0/2$$
 $p(-|dark) = 2/2$

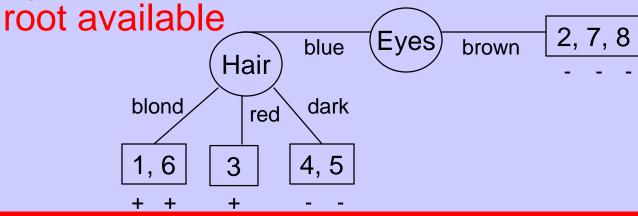
* H(decision|attribute) = H(Attractiveness|Hair) =

$$\frac{2}{5}[0] + \frac{1}{5}[0] + \frac{2}{5}[0] = 0$$

* H(Hair) =
$$-\frac{2}{5}\log_2\left(\frac{2}{5}\right) - \frac{1}{5}\log_2\left(\frac{1}{5}\right) - \frac{2}{5}\log_2\left(\frac{2}{5}\right) = 0.793$$

- A data set for classification (Quinlan)
 - * Hence the Information gain ratios are
 - ◆ Height: 0.021
 - ◆ Hair: 1.224

* By this criterion, Hair is chosen as the best



- A data set for classification (Quinlan)
 - * We may now build rules from this decision tree
 - ◆ R1: (Eyes, brown) --> (Attractiveness, -)
 - ◆ R2: (Eyes, blue) & (Hair, blond) --> (Attractiveness, +)
 - ◆ R3: (Eyes, blue) & (Hair, red) --> (Attractiveness, +)
 - ◆ R4: (Eyes, blue) & (Hair, dark) --> (Attractiveness, -)
 - * These are different rules
 - Note that after dropping conditions, however, they are the same - this is NOT generally true

- ID3 & C4.5
 - * What if too many cases?
 - Windowing
 - * What if the data is incomplete?
 - * What if the data is inconsistent?
 - * What if the data is continuous?
 - Binarization
 - Discretization
 - * Incremental algorithms?
 - * Pruning?