Abduction, Uncertainty, and Probabilistic Reasoning

Chapter 15 and more

Introduction

- Abduction is a reasoning process that tries to form plausible explanations for abnormal observations
 - Abduction is distinct different from deduction and induction
 - Abduction is inherently uncertain
- Uncertainty becomes an important issue in AI research
- Some major formalisms for representing and reasoning about uncertainty
 - Mycin's certainty factor (an early representative)
 - Probability theory (esp. Bayesian belief networks)
 - Dempster-Shafer theory
 - Fuzzy logic
 - Truth maintenance systems

Abduction

- **Definition** (Encyclopedia Britannica): reasoning that derives an explanatory hypothesis from a given set of facts
 - The inference result is a *hypothesis*, which if true, could *explain* the occurrence of the given facts

• Examples

- Dendral, an expert system to construct 3D structure of chemical compounds
 - Fact: mass spectrometer data of the compound and its chemical formula
 - KB: chemistry, esp. strength of different types of bounds
 - Reasoning: form a hypothetical 3D structure which meet the given chemical formula, and would most likely produce the given mass spectrum if subjected to electron beam bombardment

– Medical diagnosis

- Facts: symptoms, lab test results, and other observed findings (called manifestations)
- KB: causal associations between diseases and manifestations
- Reasoning: one or more diseases whose presence would causally explain the occurrence of the given manifestations
- Many other reasoning processes (e.g., word sense disambiguation in natural language process, image understanding, detective's work, etc.) can also been seen as abductive reasoning.

Comparing abduction, deduction and induction

Deduction:	major premise:	All balls in the box are black	A => B
	minor premise:	These balls are from the box	A
	conclusion:	These balls are black	В
Abduction:	rule:	All balls in the box are black	
	observation:	These balls are black	$A \Longrightarrow B$
	explanation:	These balls are from the box	Possibly
Induction:	case: observation: hypothesized rule:	These balls are from the box These balls are black All ball in the box are black	Whenev A then I but not vice vers

γA ver R sa Possibly A => B

Induction: from specific cases to general rules Abduction and deduction:

> both from part of a specific case to other part of the case using general rules (in different ways)

Characteristics of abduction reasoning

- 1. Reasoning results are **hypotheses**, not theorems (may be false *even if* rules and facts are true),
 - e.g., misdiagnosis in medicine
- 2. There may be multiple plausible hypotheses
 - When given rules A => B and C => B, and fact B both A and C are plausible hypotheses
 - Abduction is inherently uncertain
 - Hypotheses can be ranked by their plausibility if that can be determined
- 3. Reasoning is often a Hypothesize- and-test cycle
 - hypothesize phase: postulate possible hypotheses, each of which could explain the given facts (or explain most of the important facts)
 - test phase: test the plausibility of all or some of these hypotheses

- One way to test a hypothesis H is to query if some thing that is currently unknown but can be predicted from H is actually true.
 - If we also know A => D and C => E, then ask if D and E are true.
 - If it turns out D is true and E is false, then hypothesis A becomes more plausible (support for A increased, support for C decreased)
- 4. Reasoning is non-monotonic
 - Plausibility of hypotheses can increase/decrease as new facts are collected (deductive inference determines if a sentence is true but would never change its truth value)
 - Some hypotheses may be discarded, and new ones may be formed when new observations are made

Source of Uncertainty

- Uncertain data (noise)
- Uncertain knowledge (e.g, causal relations)
 - A disorder may cause any and all POSSIBLE manifestations in a specific case
 - A manifestation can be caused by more than one POSSIBLE disorders
- Uncertain reasoning results
 - Abduction and induction are inherently uncertain
 - Default reasoning, even in deductive fashion, is uncertain
 - Incomplete deductive inference may be uncertain

Probabilistic Inference

- Based on probability theory (especially Bayes' theorem)
 - Well established discipline about uncertain outcomes
 - Empirical science like physics/chemistry, can be verified by experiments
- Probability theory is too rigid to apply directly in many applications
 - Some assumptions have to be made to simplify the reality
 - Different formalisms have been developed in which some aspects of the probability theory are changed/modified.
- We will briefly review the basics of probability theory before discussing different approaches to uncertainty
- The presentation uses diagnostic process (an abductive and evidential reasoning process) as an example

Probability of Events

- Sample space and events
 - Sample space S: (e.g., all people in an area)
 - Events $E1 \subseteq S$: (e.g., all people having cough)
 - $E2 \subseteq S$: (e.g., all people having cold)
- Prior (marginal) probabilities of events
 - -P(E) = |E| / |S| (frequency interpretation)
 - -P(E) = 0.1 (subjective probability)
 - $-0 \le P(E) \le 1$ for all events
 - Two special events: \emptyset and S: P(\emptyset) = 0 and P(S) = 1.0
- Boolean operators between events (to form compound events)
 - Conjunctive (intersection): E1 ^ E2 (E1 \cap E2)
 - Disjunctive (union):
 - Negation (complement): $\sim E (E^C = S E)$
- E1 v E2 (E1 \cup E2)

- Probabilities of compound events
 - $-P(\sim E) = 1 P(E)$ because $P(\sim E) + P(E) = 1$
 - $P(E1 v E2) = P(E1) + P(E2) P(E1 ^ E2)$
 - But how to compute the *joint probability* P(E1 ^ E2)?



• Conditional probability (of E1, given E2)

– How likely E1 occurs in the subspace of E2

$$P(E1 | E2) = \frac{|E1 \land E2|}{|E2|} = \frac{|E1 \land E2| / |S|}{|E2| / |S|} = \frac{P(E1 \land E2)}{P(E2)}$$
$$P(E1 \land E2) = P(E1 | E2)P(E2)$$

- Independence assumption
 - Two events E1 and E2 are said to be independent of each other if P(E1 | E2) = P(E1) (given E2 does not change the likelihood of E1)
 - It can simplify the computation

 $P(E1 \land E2) = P(E1 | E2)P(E2) = P(E1)P(E2)$

$$P(E1 \lor E2) = P(E1) + P(E2) - P(E1 \land E2)$$

= P(E1) + P(E2) - P(E1)P(E2)
= 1 - (1 - P(E1)(1 - P(E2)))

• Mutually exclusive (ME) and exhaustive (EXH) set of events

- ME:
$$E_i \wedge E_j = \emptyset \ (P(E_i \wedge E_j) = 0), i, j = 1, ..., n, i \neq j$$

- EXH:
$$E_1 \lor ... \lor E_n = S (P(E_1 \lor ... \lor E_n) = 1)$$

Bayes' Theorem

• In the setting of diagnostic/evidential reasoning



hypotheses

evidence/manifestations

- Know prior probability of hypothesis $P(H_i)$ conditional probability $P(E_j | H_i)$

– Want to compute the *posterior probability* $P(H_i | E_j)$

- Bayes' theorem (formula 1): $P(H_i | E_j) = P(H_i)P(E_j | H_i)/P(E_j)$
- If the purpose is to find which of the n hypotheses H₁,..., H_n is more plausible given E_j, then we can ignore the denominator and rank them use *relative likelihood*

 $rel(H_i | E_j) = P(E_j | H_i)P(H_i)$

• $P(E_j)$ can be computed from $P(E_j | H_i)$ and $P(H_i)$, if we assume all hypotheses $H_1, ..., H_n$ are ME and EXH

$$P(E_j) = P(E_j \land (H_1 \lor ... \lor H_n) \quad \text{(by EXH)}$$
$$= \sum_{i=1}^n P(E_j \land H_i) \quad \text{(by ME)}$$
$$= \sum_{i=1}^n P(E_j \mid H_i) P(H_i)$$

• Then we have another version of Bayes' theorem:

$$P(H_i | E_j) = \frac{P(E_j | H_i)P(H_i)}{\sum_{k=1}^{n} P(E_j | H_k)P(H_k)} = \frac{rel(H_i | E_j)}{\sum_{k=1}^{n} rel(H_k | E_j)}$$

where $\sum_{k=1}^{n} P(E_j | H_k) P(H_k)$, the sum of relative likelihood of all n hypotheses, is a normalization factor

Probabilistic Inference for simple diagnostic problems

• Knowledge base:

$$\begin{split} & E_1, \dots, E_m: \quad \text{evidence/manifestation} \\ & H_1, \dots, H_n: \quad \text{hypotheses/disorders} \\ & E_j \text{ and } H_i \text{ are binary and hypotheses form a ME & EXH set} \\ & P(E_j \mid H_i), i = 1, \dots, n, j = 1, \dots, m \quad \text{ conditional probabilities} \end{split}$$

- Case input: E_1, \dots, E_l
- Find the hypothesis H_i with the highest posterior probability $P(H_i | E_1, ..., E_l)$
- By Bayes' theorem $P(H_i | E_1, ..., E_l) = \frac{P(E_1, ..., E_l | H_i)P(H_i)}{P(E_1, ..., E_l)}$
- Assume all pieces of evidence are conditionally independent, given any hypothesis

 $\boldsymbol{P}(\boldsymbol{E}_1, \dots \boldsymbol{E}_l \mid \boldsymbol{H}_i) = \prod_{j=1}^l \boldsymbol{P}(\boldsymbol{E}_j \mid \boldsymbol{H}_i)$

• The relative likelihood

 $rel(H_i | E_1,...,E_l) = P(E_1,...,E_l | H_i)P(H_i) = P(H_i)\prod_{j=1}^l P(E_j | H_i)$

• The absolute posterior probability

$$P(H_i | E_1, ..., E_l) = \frac{rel(H_i | E_1, ..., E_l)}{\sum_{k=1}^n rel(H_k | E_1, ..., E_l)} = \frac{P(H_i) \prod_{j=1}^l P(E_j | H_i)}{\sum_{k=1}^n P(H_k) \prod_{j=1}^l P(E_j | H_k)}$$

• Evidence accumulation (when new evidence discovered)

$$rel(H_i | E_1,...,E_l,E_{l+1}) = P(E_{l+1} | H_i)rel(H_i | E_1,...,E_l)$$

 $rel(H_i | E_1, ..., E_l, \sim E_{l+1}) = (1 - P(E_{l+1} | H_i))rel(H_i | E_1, ..., E_l)$

Assessment of Assumptions

- Assumption 1: hypotheses are mutually exclusive and exhaustive
 - Single fault assumption (one and only hypothesis must true)
 - Multi-faults do exist in individual cases
 - Can be viewed as an approximation of situations where hypotheses are independent of each other and their prior probabilities are very small

 $P(H_1 \wedge H_2) = P(H_1)P(H_2) \approx 0$ if both $P(H_1)$ and $P(H_2)$ are very small

- Assumption 2: pieces of evidence are conditionally independent of each other, given any hypothesis
 - Manifestations themselves are not independent of each other, they are correlated by their common causes
 - Reasonable under single fault assumption
 - Not so when multi-faults are to be considered

Limitations of the simple Bayesian system

- Cannot handle well hypotheses of multiple disorders
 - Suppose $H_1, ..., H_n$ are independent of each other
 - Consider a composite hypothesis $H_1^A H_2$
 - How to compute the posterior probability (or relative likelihood)

 $P(H_1^{A}H_2 | E_1,...,E_l)?$

- Using Bayes' theorem

$$P(H_1^{A}H_2 | E_1,...,E_l) = \frac{P(E_1,...E_l | H_1^{A}H_2)P(H_1^{A}H_2)}{P(E_1,...E_l)}$$

 $P(H_1^{\wedge}H_2) = P(H_1)P(H_2)$ because they are independent

 $P(E_1,...E_i | H_1^A H_2) = \prod_{j=1}^{l} P(E_j | H_1^A H_2)$ assuming E_j are independent, given $H_1^A H_2$ How to compute $P(E_j | H_1^A H_2)$? - Assuming $H_1, ..., H_n$ are independent, given $E_1, ..., E_l$?

 $P(H_1 \land H_2 | E_1, ..., E_l) = P(H_1 | E_1, ..., E_l) P(H_2 | E_1, ..., E_l)$ but this is a very unreasonable assumption



- Cannot handle causal chaining
 - Ex. A: weather of the year
 - B: cotton production of the year
 - C: cotton price of next year
 - Observed: A influences C
 - The influence is not direct (A -> B -> C)
 - P(C|B, A) = P(C|B): instantiation of B blocks influence of A on C
- Need a better representation and a better assumption

E and B are independent But when A is given, they are (adversely) dependent because they become competitors to explain A P(B|A, E) << P(B|A)

Bayesian Belief Networks (BBN)

- Definition: A BBN = (DAG, CPD)
 - **DAG**: directed acyclic graph
 - nodes: random variables of interest (binary or multi-valued) arcs: direct causal/influential relations between nodes
 - **CPD**: conditional probability distribution at each node x_i

 $P(x_i | p_i)$ where p_i is the set of all parent nodes of x_i

- For root nodes $\mathbf{p}_i = \emptyset$, so $\mathbf{P}(\mathbf{x}_i | \mathbf{p}_i) = \mathbf{P}(\mathbf{x}_i)$ Since roots are not influenced by anyone, they are considered independent of each other
- Example BBN



P(A) = 0.001P(B|A) = 0.3P(C|A) = 0.2 $P(C|\sim A) = 0.005$ P(E|C) = 0.4

 $P(B|\sim A) = 0.001$ P(D|B,C) = 0.1 $P(D|B,\sim C) = 0.01$ $P(D|\sim B,C) = 0.01$ $P(D|\sim B,\sim C) = 0.00001$ $P(E|\sim C) = 0.002$

- Independence assumption
 - $P(x_i | p_i, q) = P(x_i | p_i)$ where q is any set of variables (nodes) other than x_i and its successors
 - \mathbf{p}_i blocks influence of other nodes on x_i and its successors (q influences x_i only through variables in \mathbf{p}_i)



 With this assumption, the complete joint probability distribution of all variables in the network can be represented by (recovered from) local CPD by chaining these CPD

$$\boldsymbol{P}(\boldsymbol{x}_1,...,\boldsymbol{x}_n) = \prod_{i=1}^n \boldsymbol{P}(\boldsymbol{x}_i \mid \boldsymbol{p}_i)$$

P(A, B, C, D, E)

- = P(E|A, B, C, D) P(A, B, C, D) by Bayes' theorem
- = P(E|C) P(A, B, C, D)

- by indep. assumption
- $= P(E|C) P(D|A, \boldsymbol{B}, \boldsymbol{C}) P(A, B, C)$
- = P(E|C) P(D|B, C) P(C|A, B) P(A, B)
- = P(E|C) P(D|B, C) P(C|A) P(B|A) P(A)

Inference with BBN

• Belief update

- Original belief (no variable is instantiated): the prior probability $P(x_i)$ If x_i is a root, then $P(x_i)$ is given in BBN. Otherwise, $P(x_i) = \sum_{p_i} P(x_i | p_i) P(p_i)$ $P(x_i | p_i)$ is given, but computer $P(p_i)$ is complicated Ex: $P(B,C) = P(A,B,C) + P(\sim A,B,C)$ $= P(B | A,C)P(A,C) + P(B | \sim A,C)P(\sim A,C)$ $= P(B | A)P(C | A)P(A) + P(B | \sim A)P(C | \sim A)P(\sim A)$
- When some variables are instantiated (say x_j has value X_j), beliefs on all other variable x_i is changed to $P(x_i | X_j)$

 $P(x_i | X_j)$ can be computed from the joint probability distribution Ex : d = D and e = E

$$P(A \mid D, E) = \frac{P(A, D, E)}{P(D, E)} = \frac{\sum_{b,c} P(A, b, c, D, E)}{\sum_{a,b,c} P(a, b, c, D, E)}$$

This approach is not computationally feasible with large network

-*Algorithmic approach* (Pearl and others)

- Singly connected network, SCN (also known as poly tree) there is at most one undirected path between any two nodes (i.e., the network is a tree if the direction of arcs are ignored)
- The influence of the instantiated variable spreads to the rest of the network along the arcs
 - The instantiated variable influences its predecessors and successors differently
 - Computation is linear to the diameter of the network (the longest undirected path)



- For non-SCN (network with general structure)
 - Conditioning: find the the network's smallest cutset *C* (a set of nodes whose removal will render the network singly connected)
 for each instantiation of *C*, compute the belief update with the SCN algorithm
 - Combine the results from all possible instantiation of *C*.
 - Computationally expensive (finding the smallest cutset is itself NP-hard, and total number of possible instantiations of *C* is $O(2^{A}|C|)$.)

- Stochastic simulation

- Randomly generate large number of instantiations of ALL variables $I_k^{(n)}$ according to CPD
- Only keep those instantiations $I_k^{(n)}$ which are consistent with the values of given instantiated variables
- Updated belief of those un-instantiated variables as their frequencies in the pool of recorded $I_{k}^{(n)}$
- The accuracy of the results depend on the size of the pool (asymptotically approaches the exact results)

• MAP problems

- Let X denote the set of all variables in a BBN, $V \subseteq X$ the set of instantiated variables, U = X - V the set of all un-instantiated varialbes. Then the MAP (maximum aposteriori probability) problem is to find the most probable instantiation of U, given V, i.e., $\max_{u}(P(U|V))$
- This is an optimization problem
- Algorithms developed for *exact* solutions for different special BBN (Peng, Cooper, Pearl) have exponential complexity
- Other techniques for *approximate* solutions
 - Genetic algorithms
 - Neural networks
 - Simulated annealing
 - Mean field theory

Noisy-Or BBN

- A special BBN of binary variables (Peng & Reggia, Cooper)
 - Each link $x_i \rightarrow x_j$ is associated with a probability value called *causal strength* c_{ij} that measures the strength of x_i *alone* may cause x_j , i.e., $c_{ij} = P(x_i | x_j \text{ is true and all others in } p_i \text{ are false})$
 - Causation independence: parent nodes influence a child independently
- Advantages:
 - One-to-one correspondence between causal links and causal strengths
 - Easy for humans to understand (acquire and evaluate KB)
 - Fewer # of probabilities needed in KB

Complete joint prob. distribution : 2ⁿ

General BBN : $\sum_{i=1}^{n} 2^{|\mathbf{p}_i|}$

Noisy - Or BBN :

$$\sum_{i=1}^{n} 2^{n} |\mathbf{p}_i|$$

- Computation is less expensive

• Disadvantage: less expressive (less general)

Learning BBN (from case data)

- Need for learning
 - Experts' opinions are often biased, inaccurate, and incomplete
 - Large databases of cases become available
- What to learn
 - Learning CPD when DAG is known (easy)
 - Learning DAG (hard)
- Difficulties in learning DAG from case data
 - There are too many possible DAG when # of variables is large (more than exponential)
 - n = 3, # of possible DAG = 25
 - n = 10, # of possible DAG = 4*10^18
 - Missing values in database
 - Noisy data

• Approaches

- Early effort: Based on variable dependencies (Pearl)
 - Find all pairs of variables that are dependent of each other (applying standard statistical method on the database)
 - Eliminate (as much as possible) indirect dependencies
 - Determine directions of dependencies
 - Learning results are often incomplete (learned BBN contains indirect dependencies and undirected links)

-Bayesian approach (Cooper)

- Find the most probable DAG, given database DB, i.e., max(*P*(*DAG*/*DB*)) or max(*P*(*DAG*, *DB*))
- Based on some assumptions, a formula is developed to compute *P*(*DAG*, *DB*) for a given pair of DAG and DB
- A hill-climbing algorithm (K2) is developed to search a (sub)optimal DAG
- Extensions to handle some form of missing values

- *Minimum description length* (MDL) (Lam)

- Sacrifices accuracy for simpler (less dense) structure
 - Case data not always accurate
 - Fewer links imply smaller CPD tables and less expensive inference
- L = L1 + L2 where
 - *L1*: the length of the encoding of DAG (smaller for simpler DAG)
 - L2: the length of the encoding of the difference between DAG and DB (smaller for better match of DAG with DB)
 - Smaller *L2* implies more accurate (and more complex) DAG, and thus larger *L1*
- Find DAG by heuristic best-first search, that Minimizes L

-Neural network approach (Neal, Peng)

- For noisy-or BBN
- Maximizing $L = \ln \prod_{V' \in D} P(\tilde{V} = V')$ where

D: case database; V^r : case in D; \tilde{V} : state vector of the learned network

L measures the similarity of the two distributions : one in D, another in the learned network

Dempster-Shafer theory

- A variation of Bayes' theorem to represent **ignorance**
- Uncertainty and ignorance
 - Suppose two events A and B are ME and EXH, given an evidence E

A: having cancer B: not having cancer E: smoking

- By Bayes' theorem: our beliefs on A and B, given E, are measured by P(A|E) and P(B|E), and P(A|E) + P(B|E) = 1
- In reality,

I may have some belief in A, given E

I may have some belief in B, given E

I may have some belief not committed to either one,

- The uncommitted belief (ignorance) should not be given to either A or B, even though I know one of the two must be true, but rather it should be given to "A or B", denoted {A, B}
- Uncommitted belief may be given to A and B when new evidence is discovered

- Representing ignorance
 - Frame of discernment : $\mathbf{q} = \{\mathbf{h}_1, ..., \mathbf{h}_n\}$, a set of ME and EXH hypotheses. The power set $2^{\mathbf{q}}$ is organized as a lattice of super/subset relations. Each node S is a subset of hypotheses ($S \subseteq \mathbf{q}$)
 - -Ex: $\theta = \{A,B,C\}$

Each node *S* is associated with a basic probability assignment m(S) $0 \le m(S) \le 1;$ $m(\emptyset) = 0;$ $\sum_{S \subseteq q} m(S) = 1$



• Belief function

$$Bel(S) = \sum_{S' \subseteq S} m(S'); \quad Bel(\emptyset) = 0; \quad Bel(\mathbf{q}) = 1$$

$$Bel(\{A, B\}) = m(\{A, B\}) + m(\{A\}) + m(\{B\}) + m(\emptyset)$$

$$= 0.1 + 0.1 + 0.2 + 0 = 0.4$$

$$Bel(\{A, B\}^{C}) = Bel(\{C\}) = 0.3$$

– Plausibility (upper bound of belief of a node) All belief not committed to S^{C} may be committed to S $Pls(S) = 1 - Bel(S^{C})$ $Pls(\{A, B\}) = 1 - Bel(\{C\}) = 1 - 0.3 = 0.7$ [Bel(S), Pls(S)] belief interval $\{A,B,C\}\ 0.15$ Lower Upper bound bound {A,B} 0.1 {B,C}0.05 {A,C} 0.1 (known (maximally belief) possible) $\{\bar{B}\}\ 0.2$ {A} 0.1 C}0.3

 $\{\} 0$

- Evidence combination (how to use D-S theory)
 - Each piece of evidence has its own m(.) function for the same θ $q = \{A, B\}$: A : having cancer; B : not having cancer





$$m(\{A\}) = \frac{m_1(\{A\})m_2(\{A\}) + m_1(\{A\})m_2(\{A,B\}) + m_1(\{A,B\})m_2(\{A\}))}{1 - [m_1(\{A\})m_2(\{B\}) + m_1(\{B\})m_2(\{A\})]}$$

$$= \frac{0.2 \cdot 0.7 + 0.2 \cdot 0.1 + 0.3 \cdot 0.7}{1 - [0.2 \cdot 0.2 + 0.5 \cdot 0.7]} = \frac{0.37}{0.61} = 0.607$$

$$m(\{B\}) = \frac{m_1(\{B\})m_2(\{B\}) + m_1(\{B\})m_2(\{A,B\}) + m_1(\{A,B\})m_2(\{B\}))}{1 - [m_1(\{A\})m_2(\{B\}) + m_1(\{B\})m_2(\{A\})]}$$

$$= \frac{0.5 \cdot 0.2 + 0.5 \cdot 0.1 + 0.3 \cdot 0.2}{1 - [0.2 \cdot 0.2 + 0.5 \cdot 0.7]} = \frac{0.21}{0.61} = 0.344$$

$$m(\{A,B\}) = \frac{m_1(\{A,B\})m_2(\{A,B\})}{0.61} = \frac{0.03}{0.61} = 0.049$$

– Ignorance is reduced

from m1($\{A,B\}$) = 0.3 to m($\{A,B\}$) = 0.049)

– Belief interval is narrowed

A: from [0.2, 0.5] to [0.607, 0.656]

- B: from [0.5, 0.8] to [0.344, 0.393]
- Advantage:
 - The only formal theory about ignorance
 - Disciplined way to handle evidence combination
- Disadvantages
 - Computationally very expensive (lattice size $2^{|\theta|}$)
 - Assuming hypotheses are ME and EXH
 - How to obtain m(.) for each piece of evidence is not clear, except subjectively

Fuzzy sets and fuzzy logic

• Ordinary set theory

 $f_{A}(x) = \begin{cases} 1 & \text{if } x \in A \\ \hline 0 & \text{otherwise} \end{cases}$ $f_{A}(x) \text{ is called the characteristic or membership function of set } A$ Predicate $A(x) = \begin{cases} \frac{1}{0} & \text{if } x \in A \\ \hline 0 & \text{otherwise} \end{cases}$

When it is uncertain if $x \in A$, use probability $P(x \in A)$

- There are sets that are described by vague linguistic terms (sets without hard, clearly defined boundaries), e.g., tall-person, fastcar
 - Continuous
 - Subjective (context dependent)
 - Hard to define a clear-cut 0/1 membership function

- Fuzzy set theory
 - Relax $f_A(x)$ from binary {0,1} to continuous [0,1] stands for the degree x is thought to belong to set A

height(john) = 6'5''Tall(john) = 0.9height(harry) = 5'8''Tall(harry) = 0.5height(joe) = 5'1''Tall(joe) = 0.1

- Examples of membership functions



- Fuzzy logic: many-value logic
 - Fuzzy predicates (degree of truth) $F_A(x) = y$ if $f_A(x) = y$
 - Connectors/Operators

negation: $\sim F_A(x) = 1 - F_A(x)$ conjunction: $F_A(x) \wedge F_B(x) = \min\{F_A(x), F_B(x)\}$ disjunction: $F_A(x) \vee F_B(x) = \max\{F_A(x), F_B(x)\}$

- Compare with probability theory
 - Prob. Uncertainty of outcome,
 - Based on large # of repetitions or instances
 - For each experiment (instance), the outcome is either true or false (without uncertainty or ambiguity) unsure before it happens but sure after it happens

Fuzzy: vagueness of conceptual/linguistic characteristics

 Unsure even after it happens whether a child of tall mother and short father is tall unsure before the child is born unsure after grown up (height = 5'6")

- Empirical vs subjective (testable vs agreeable)
- Fuzzy set connectors may lead to unreasonable results
 - Consider two events A and B with P(A) < P(B)
 - If A => B (or A ⊆ B) then
 P(A ^ B) = P(A) = min{P(A), P(B)}
 P(A v B) = P(B) = max{P(A), P(B)}
 - Not the case in general

 $P(A \land B) = P(A)P(B|A) \le P(A)$ $P(A \lor B) = P(A) + P(B) - P(A \land B) \ge P(B)$ (equality holds only if P(B|A) = 1, i.e., A => B)

– Something prob. theory cannot represent

- Tall(john) = 0.9, ~Tall(john) = 0.1 Tall(john) ^ ~Tall(john) = min{0.1, 0.9} = 0.1 john's degree of membership in the fuzzy set of "medianheight people" (both Tall and not-Tall)
- In prob. theory: $P(john \in Tall \land john \notin Tall) = 0$

Uncertainty in rule-based systems

- Elements in Working Memory (WM) may be uncertain because
 - Case input (initial elements in WM) may be uncertain

Ex: the CD-Drive does not work 70% of the time

- Decision from a rule application may be uncertain even if the rule's conditions are met by WM with certainty
 Ex: flu => sore throat with high probability
- Combining symbolic rules with numeric uncertainty: Mycin's Uncertainty Factor (CF)
 - An early attempt to incorporate uncertainty into KB systems
 - $-CF \in [-1, 1]$
 - Each element in WM is associated with a CF: certainty of that assertion
 - Each rule C1,..., Cn => Conclusion is associated with a CF:
 certainty of the association (between C1,...Cn and Conclusion).

- CF propagation:
 - Within a rule: each *Ci* has CFi, then the certainty of Action is *min{CF1,...CFn} * CF-of-the-rule*
 - When more than one rules can apply to the current WM for the same *Conclusion* with different CFs, the *largest of these CFs* will be assigned as the CF for *Conclusion*
 - Similar to fuzzy rule for conjunctions and disjunctions
- Good things of Mycin's CF method
 - Easy to use
 - CF operations are reasonable in many applications
 - Probably the only method for uncertainty used in real-world rule-base systems
- Limitations
 - It is in essence an ad hoc method (it can be viewed as a probabilistic inference system with some strong, sometimes unreasonable assumptions)
 - May produce counter-intuitive results.