

Summary: Knowledge Representation

1. Knowledge Representation and Propositional Logics p.1
2. DPLL+Logical reasoning for Sudokus p.7
3. Logic, Hardness, GSAT, MaxSat p.21
4. Clause Learning, Planning and Applications p.29
5. Propositional Logic, MaxSat + Exam Preparation p.36
6. Probabilistic Graphical Models 1 p.40
7. Probabilistic Graphical Models 2 p.55
8. Probabilistic Graphical Models 3 p.74
- 8a: Canvas Quiz with Answers p.87
9. Probabilistic Graphical Models 4 p.97
10. Ontologies, and Ontology Languages
11. Description Logics
12. Description Logic Reasoning
13. Knowledge Engineering

Summary Knowledge Representation

1. Knowledge Representation and Propositional Logics
2. DPLL+Logical reasoning for Sudokus
3. Logic, Hardness, GSAT, MaxSat
4. Clause Learning, Planning and Applications
5. Propositional Logic, MaxSat + Exam Preparation
6. Probabilistic Graphical Models 1
7. Probabilistic Graphical Models 2
8. Probabilistic Graphical Models 3
9. Probabilistic Graphical Models 4
10. Ontologies, and Ontology Languages
11. Description Logics
12. Description Logic Reasoning
13. Knowledge Engineering

Knowledge Representation and Propositional Logics

“Definition of intelligence”

- carry out complex reasoning (solve physics problems, prove)
- draw plausible inferences (diagnose cars, solve a murder)
- use natural language (read stories and answer questions about them, carry out extended conversation)
- solving novel complex problems (generating plans, designing)
- social activities that require a theory of mind

Two main lines of development in AI

- symbolic representations
- statistical representation

Statistical vs. symbolic AI:
very different types of applications

statistical:

- pattern recognition (images, sound, shapes)
- motor skills (robots)
- speech generation (sound)
- search engines

symbolic:

- planning (autonomous space missions)
- reasoning (diagnosis, design, decision support)
- language generation (conversations)
- search engines

Strengths & Weaknesses

	Symbolic	Connectionist
Construction	Human effort	Data hunger
Scaleable	+/- (worse with more data)	+/- (worse with less data)
Explainable	+	-
Generalisable	Performance cliff	Performance cliff

The goal of logic in KR

- To state statements which are known to be true (the knowledge base)
- Some statements that describe the current state of the world (the premises)
- To state statements for which we want to check if they are true (the conclusions ”)
- To see if the conclusions can be derived from the knowledge base + the premises through logical reasoning
- A variety of related tasks

The ingredients of a logic in KR

- How to formulate the statements (Syntax)
- Assign meaning to the statements (Semantics)
- Assign what can be derived (Calculus)

And all of this differs from application domain to application domain and even application to application.

3) PROPOSITIONAL LOGIC

The knowledge base: simple statements

- In propositional logic, a knowledge base can have simple statements , and complex statements.
- A simple statement states one fact about the world :
 - Examples: it is hot , it is raining , it is humid
 - In our examples we often just use a single letter for a simple statement:
 - Examples: P, Q, R.

The knowledge base: complex statements

- Complex statements are made by combining other statements . If S and T are statements (simple or complex), then the following are also complex statements:

These can be any combination

- $S \wedge T$ (“S and T”) (conjunction)
- $S \vee T$ (“S or T”) (disjunction)
- $\neg S$ (“not S”) (negation)
- $S \rightarrow T$ (“S implies T”) (implication)
- These can be any combination
- Examples: *it-is-hot* \wedge *it-is-raining*
 - \neg *it-is-humid*
 - *it-is-hot* \wedge *it-is-raining* \rightarrow *it-is-humid*
 - $(P_1 \vee \neg P_2) \wedge P_3$
 - $((P \vee H) \wedge \neg H) \rightarrow P$

The meaning of simple statements

- In propositional logic, the meaning of a statement is either true or false.
- We call this the “interpretation function” I :
it maps any statement p to true or false: $I(p)=0$ or $I(p)=1$ (0 means false, 1 means true)
- Example:
 - $I(\text{it is hot})=1$ means that it is hot today
 - $I(\text{it is hot})=0$ means that it is not hot today
 - $I(Q)=1$ means that Q is true

The meaning of complex statements

- The meaning of complex statements is calculated from the meaning of the parts:

- The meaning of $P \wedge Q$ is calculated by combining the meaning of P and Q :
 - $P \wedge Q$ is only true if **both** P and Q are true, otherwise $P \wedge Q$ is false.
 - Mathematically: $I(P \wedge Q)=1$ only if $I(P)=1$ and $I(Q)=1$
 - We can also write this in a table, which tells us the meaning of $P \wedge Q$ for each combination of P and Q :

P	Q	$P \wedge Q$
1	1	1
1	0	0
0	1	0
0	0	0

- In the same way, the meaning of $P \vee Q$ is calculated by combining the meaning of P and Q :
 - $P \vee Q$ is only true if **at least one of** P and Q are true, otherwise $P \vee Q$ is false.
 - Mathematically: $I(P \vee Q)=1$ only if $I(P)=1$ or $I(Q)=1$, or both
 - We can also write this in a table, which tells us the meaning of $P \vee Q$ for each combination of P and Q :

P	Q	$P \vee Q$
1	1	1
1	0	1
0	1	1
0	0	0

- In the same way, the meaning of $\neg P$ is calculated from the meaning of P
 - $\neg P$ is only true if P is false, otherwise $\neg P$ is false.
 - Mathematically: $I(\neg P)=1$ only if $I(P)=0$
 - We can also write this in a table, which tells us the meaning of $\neg P$ for each meaning of P :

P	$\neg P$
1	0
0	1

- The meaning of $P \rightarrow Q$ also computed from the meanings of P and Q , but in a bit of a strange way:

P	Q	$P \rightarrow Q$
1	1	1
1	0	0
0	1	1
0	0	1

- (Exercise: make the truth table of $\neg P \vee Q$)

We can use these tables to determine the meaning (the “truth”) of **complex statements**:

To calculate the meaning of $(P_1 \vee \neg P_2) \wedge P_3$ we do the following (read the table left-to-right):

- take all 8 combinations of the possible values for P_1 , P_2 and P_3
- Calculate the value for $\neg P_2$ from the values of P_2
- Calculate the value for $(P_1 \vee \neg P_2)$ by combining the values for P_1 and $\neg P_2$
- Calculate the value for $(P_1 \vee \neg P_2) \wedge P_3$ by combining the values for $(P_1 \vee \neg P_2)$ and P_3

p_1	p_2	p_3	$\neg p_2$	$(p_1 \vee \neg p_2)$	P
1	1	1	0	1	1
1	1	0	0	1	0
1	0	1	1	1	1
1	0	0	1	1	0
0	1	1	0	0	0
0	1	0	0	0	0
0	0	1	1	1	1
0	0	0	1	1	0

Special statements: tautologies

Let's look at the complex statement
 $((P \vee H) \wedge \neg H) \rightarrow P$.

P	H	$P \vee H$	$(P \vee H) \wedge \neg H$	$((P \vee H) \wedge \neg H) \Rightarrow P$
False	False	False	False	True
False	True	True	False	True
True	False	True	True	True
True	True	True	False	True

There is something special with the result in the final column: this statement is **always true!**

Such statements, which are true in any world that we can possibly imagine, are called **tautologies (or: valid)**

Special statements: inconsistencies

There are also statements which are always false.

They are false for any of the possible combinations of truth values of their parts, in other words, they are false in every world we can possibly imagine.

Such statements are called inconsistent (or: invalid).

The statements which are not a tautology and also not inconsistent (in other words: they are true for some combination of truth values and false for others) are called satisfiable

**Example: $(A \vee B) \wedge (\neg A \vee \neg B)$ is satisfiable,
but not a tautology.**

Some terms

- The meaning or semantics of a sentence is determined by its interpretation .
- Given the truth values of all symbols in a sentence, it can be “evaluated” to determine its truth value (True or False).
- A model for a KB is a “possible world” (assignment of truth values to propositional symbols) in which each sentence in the KB is True.

More terms

• A satisfiable sentence is a sentence that can be made true in at least one world (i.e. is not inconsistent)

• $P \models Q$, written $P \vDash Q$, means that whenever P is True, so is Q . In other words, all models of P are also models of Q

In Propositional Logics (and in other Logics) those notions can often be reduced to each other.

- E.g. P is satisfiable iff $\neg P$ is not valid
- Or $P \models Q$ iff $P \wedge \neg Q$ is unsatisfiable

So why are we doing all this?

1. We now have language in which we can make simple and complex statements about the world (“propositional
2. We can formulate problems that we need to solve as complex statements in propositional logic
3. We have a formal and unambiguous semantics to define if a statement can be satisfied
4. We have a method (e.g. Truth Tables) to determine which valuations for propositions satisfies the problem statements (a calculus)
5. Such a “satisfying truth assignment” is then the solution to our problem.

Practice Problem

Problem: A Diplomatic Problem

You are chief of protocol for the embassy ball. The crown prince instructs you either to invite Peru or to exclude Qatar. The queen asks you to invite either Qatar or Romania or both. The king, in a spiteful mood, wants to snub either Romania or Peru or both. Who do you invite?

- Three propositional symbols P , Q and R :

$P \equiv$	invite Peru	$\neg P \equiv$	exclude Peru
$Q \equiv$	invite Qatar	$\neg Q \equiv$	exclude Qatar
$R \equiv$	invite Romania	$\neg R \equiv$	exclude Romania

- Three operators: negation \neg , disjunction \vee , conjunction $\&$, \wedge
- Exercise: 5 minutes to formalise the problem in PL.

Representing the Diplomatic Problem in PL

- The problem can be formalized as

$$\begin{array}{lll} \text{prince:} & \text{invite Peru or exclude Qatar} & \equiv P \vee \neg Q \\ \text{queen:} & \text{invite Qatar or Romania or both} & \equiv Q \vee R \\ \text{king:} & \text{snub Romania or Peru or both} & \equiv \neg R \vee \neg P \end{array}$$

- $\Sigma = (P \vee \neg Q) \& (Q \vee R) \& (\neg R \vee \neg P)$
- Any assignment (choice of truth or falsity of the symbols) is a solution for the diplomatic problem.
- For example, let P and Q be true, and R false, then Σ is true.
- The assignment $P=T$, $Q=T$ and $R=F$ satisfies Σ .
- Therefore: invite Peru and Qatar, exclude Romania is solution!

Reasoning: Solving the diplomatic problem?

- Reminder: $\Sigma = (P \vee \neg Q) \& (Q \vee R) \& (\neg R \vee \neg P)$
- We can use truth tables

P	Q	R	$P \vee \neg Q$	$Q \vee R$	$\neg R \vee \neg P$	Σ
T	T	T	T	T	F	F
T	T	F	T	T	T	T
T	F	T	T	T	F	F
T	F	F	T	F	T	F
F	T	T	F	T	T	F
F	T	F	F	T	T	F
F	F	T	T	T	T	T
F	F	F	T	F	T	F

- Problem: we need 2^n solutions, i.e. with 3 variables 8 rows, with 20 variables 1048576, and with 50 variables 1125899906842624.

Problem solving by Satisfiability testing

1. We have to answer a question
2. We formulate question as a complex statement S in propositional logic
3. We try to find a “satisfying assignment” of truth values (in other words: a combination of truth values that makes S true)
4. That combination of truth values gives us the solution to the original problem.

Problem : if S contains n symbols, we need to check 2^n combinations of truth values.
Suppose we want to solve Sudoku with 721 variables.

DPLL+Logical reasoning for Sudokus

Propositional Logic

The language of propositional logic (Syntax)

- $S \wedge T$ (“ S and T ”) (conjunction, sometimes “and”, $\&$)
- $S \vee T$ (“ S or T ”) (disjunction, sometimes “or”)
- $\neg S$ (“not S ”) (negation)
- $S \rightarrow T$ (“ S implies T ”) (implication)

One special syntax: Clausal Normal Form:

$(P_1 \vee P_2 \vee \dots \vee P_k) \wedge (Q_1 \vee Q_2 \vee \dots \vee Q_n) \wedge \dots \wedge (R_1 \vee R_2 \vee \dots \vee R_m)$
will be explained in some minutes

The meaning of propositional logic (Semantics) 1

- Propositions P and Q are interpreted to be either false or true (truth assignment, e.g. $I(P)=\text{True}$ or $I(P)=1$)
- Complex Sentences S can be calculated to be true or false according to an assignment e.g. $I(\text{Rains} \rightarrow \text{Wet})$ is True iff $I(\text{Rains})=\text{False}$ (0) or $I(\text{Wet})=\text{True}$ (1)

The meaning of propositional logic (Semantics) 2

- Inconsistent sentence S : always false (boring) $I(S)=0$ for all assignment functions I
- Valid sentences S (tautology): always true (boring) $I(S)=1$ for all assignment functions I
- Satisfiable sentence S : interesting!

There is an assignment I such that $I(S)=1$. This is also called a model .
This is the task of SAT solving: finding a “model”

SAT solving for Problem Solving

- Represent your problem in Propositional Logic
- Is a problem solvable? = is the sentence satisfiable?

- Finding a solution = finding a model (= find a truth assignment I that makes the sentence S true, i.e. $I(S)=1$)
- Truth assignment is the solution to the problem

A calculus for propositional logic

► The problem can be formalized as

$$\begin{array}{lll} \text{prince:} & \text{invite Peru or exclude Qatar} \equiv & P \text{ or } \neg Q \\ \text{queen:} & \text{invite Qatar or Romania or both} \equiv & Q \text{ or } R \\ \text{king:} & \text{snub Romania or Peru or both} \equiv & \neg R \text{ or } \neg P \end{array}$$

► $\Sigma = (P \text{ or } \neg Q) \& (Q \text{ or } R) \& (\neg R \text{ or } \neg P)$

P	Q	R	$P \text{ or } \neg Q$	$Q \text{ or } R$	$\neg R \text{ or } \neg P$	Σ
F	T	T	T	T	F	F
T	T	F	T	T	T	T
T	F	T	T	F	F	F
T	F	F	T	F	T	F
F	T	T	F	T	T	F
F	T	F	F	T	T	F
F	F	T	T	T	T	T
F	F	F	T	F	T	F

A calculus for propositional logic: aka How to find a satisfying model?

- Truth table method
- But: size of the table is 2^N (N = number of variables)
- $2^{729} = 10^{219}$

And the problem can be even worse!
Luckily there are more efficient algorithms ...

Clause Normal Form

A special form of statement: the clausal normal form

Complex statements can be very complex, with many combinations of \neg , \vee , \wedge , \rightarrow

A special kind of complex statement (the '**clausal normal form**') that is easy for Computer programmes to process

The clausal normal form looks like;

$$(P_1 \vee P_2 \vee \dots \vee P_k) \wedge (Q_1 \vee Q_2 \vee \dots \vee Q_n) \wedge \dots \wedge (R_1 \vee R_2 \vee \dots \vee R_m)$$

where each of the letters can be **positive** (P_k) or **negative** ($\neg P_k$) atoms, and called **Literals**.

In other words: it is a conjunction of disjunctions of literals.

Theorem: Any complex statement can be rewritten to a logically equivalent statement in clausal normal form

Examples of clausal normal form

- $p_1, \neg p_1, \neg p_5$ are literals. They are also in CNF
- $(p \vee q)$ is in CNF
- $((p_1 \vee p_2) \wedge (\neg p_1 \vee p_3))$ and $((p_1 \vee p_2) \wedge \neg p_1)$ are in CNF
- $((p_1 \wedge p_2) \vee (\neg p_1 \wedge p_3))$ and $(p \vee (p \wedge \neg p))$ are not in CNF.

How to rewrite a statement into Clausal Form

- ① Remove implications $(p \rightarrow q) = (\neg p \vee q)$
- ② Move negations in front of propositional symbols
 - $\neg(p \vee q) = \neg p \wedge \neg q$
 - $\neg(p \wedge q) = \neg p \vee \neg q$
- ③ Move conjunctions “outside” over disjunctions
 - $r \vee (p \wedge q) = (r \vee p) \wedge (r \vee q)$ (result: a conjunction of clauses)
- ④ You can read the separate clauses:
 - $(p_1 \vee p_2) \wedge (q_1 \vee q_2) = (p_1 \vee p_2), (q_1 \vee q_2)$
- ⑤ Sometimes we write a disjunction of clauses as $\{\{p_1, p_2\}, \{q_1, q_2\}\}$
The outer comma means conjunction (“and”)
the inner comma means disjunction (“or”)

This can be done for any propositional formula. It doesn't matter in which order you apply the rule, the result is always the same

15

How to rewrite a statement into Clausal Form

- You may also have to use **simplification rules**:
 - $(p \vee p) = (p \wedge p) = \neg \neg p = p$
 - $(p \vee q) = (q \vee p), (p \wedge q) = (q \wedge p)$
 - $(p \vee q) \vee r = p \vee (q \vee r), (p \wedge q) \wedge r = p \wedge (q \wedge r)$
- **Repeatedly apply rules** to the result:
$$\begin{aligned} & \neg(\neg p \wedge (q \vee \neg r)) = \neg \neg p \vee \neg(q \vee \neg r) = \\ & \quad p \vee (\neg q \wedge \neg \neg r) = p \vee (\neg q \wedge r) \end{aligned}$$
- Sometimes **multiple rules apply** (always same result)

$$\begin{aligned} & ((p \wedge q) \vee r) \vee s = \\ & ((p \vee r) \wedge (q \vee r)) \vee s = \\ & ((p \vee r \vee s) \wedge (q \vee r \vee s)) \end{aligned}$$

$$\begin{aligned} & ((p \wedge q) \vee r) \vee s = \\ & ((p \wedge q) \vee (r \vee s) = \\ & ((p \vee r \vee s) \wedge (q \vee r \vee s)) \end{aligned}$$

CNF alternative notation in ASCII

$(\neg A \vee \neg B \vee E) \wedge (\neg E \vee A) \wedge (\neg E \vee B) \wedge (\neg C \vee F) \wedge (\neg F \vee C) \wedge (I)$

$(A' + B' + E)(E' + A)(E' + B)$
 $(C' + F)(F' + C)$
 (I)

DIMACS Format:
-1 -2 5 0
-5 1 0
-5 2 0
-3 6 0
-6 3 0
9 0

Davis Putnam (DPLL)

- SAT is NP Complete (proved by Cook in 1971, as the first NP complete problem)
- This means that any fast SAT solver could also be used to solve lots of other exponential problems
- This is widely believed to be impossible: nobody believes NP=P (but still an open problem)
- So should we give up all hope?
- Remarkably, even though they are exponential in the worst case (Cook was right in 1971....), modern SAT solver are very fast most of the time

Davis Putnam algorithm

- It was not developed by Davis & Putnam, but by Davis, Logemann & Loveland
- The original algorithm by Davis & Putnam used potentially exponential space
- The improved algorithm by DLL is now known as DP...
- It dates from the '60s, is (of course) exponential in the worst case, but performs remarkably well (if we add a few heuristics)

DPLL version 1

Davis – Putnam – Loveland – Logemann 1962

```
Let pa be a partial assignment  $\alpha$  a formula
dpll_1(pa) {
    if (pa makes  $\alpha$  false) return false;
    if (pa makes  $\alpha$  true) return true;

    choose P in  $\alpha$  ;
    if (dpll_1(pa  $\cup$  {P=false}))
        return true;
    else return dpll_1(pa  $\cup$  {P=true});
}
```

Returns true if α is satisfiable, false otherwise

We need a few new terms

- A “**variable**” is one of the simple statements (the letters) that make up the problem statement
- A **clause** is one of the disjuncts in the big disjunction in the input statement (which is in clausal normal form)
- A **literal** is a variable or its negation (both P and $\neg P$ are literals)
- A **unit clause** is a clause of length 1 (= only one literal)
- A literal P is **pure** if $\neg P$ appears nowhere (and vice versa)

$(\neg A \vee B \vee E) \wedge (\neg E \vee A) \wedge (\neg E \vee B) \wedge (\neg C \vee F) \wedge (\neg F \vee C) \wedge (I)$

<i>clause</i>	<i>literal</i>	<i>pure</i>	<i>literal</i>	<i>unit clause</i>
---------------	----------------	-------------	----------------	--------------------

- An **empty clause** is (trivially) false (disjunction: at least one of these must be true)
- An **empty set of clauses** is (trivially) true (conjunction: all of these have to be true)

23

Davis-Putnam procedure

It consists of two stages:

simplify
split

1. It **simplifies** the formula, and checks a few easy cases.
If such an easy case exists
it can assign the right value to one of the truth values
and **repeats the process**
2. If there is no such easy case,
it **picks one of the predicates** and **picks a truth value to it**,
and then **repeats the process**
3. The process **stops**
 1. when all clauses have been satisfied (success!), or
 2. when one becomes inconsistent (failure)

backtrack

If there is failure, that could be because we picked the wrong truth value for a letter in step 2. So we have to “**backtrack**”:

- 4a. for the letter from step 2 we pick the other truth value,
and try again.
- 4b. In step 2 we pick another letter

24

Part 1: Simplification rules

1.1 Tautology:

if a clause contains $\neg P \vee P$ it can be removed
(because the clause will be satisfied in all cases)

1.2 if there is a **pure literal** (occurring only positive or negative),
it can be set to true (or false)
(because this choice doesn't affect any other clause)

1.3 if there is a **unit clause** (consisting of only one literal),
that literal can be set to true
(because this is the only way the clause can be satisfied)
(of course, if the literal is of the form $\neg P$, that makes P false)

Part 1: Simplification rules (b)

If literal L_1 is true, then clause $(L_1 \vee L_2 \vee \dots)$ is true

If clause C_1 is true, then $C_1 \wedge C_2 \wedge C_3 \wedge \dots$ has the same
value as $C_2 \wedge C_3 \wedge \dots$

Therefore: Okay to delete clauses containing true literals!

If literal L_1 is false, then clause $(L_1 \vee L_2 \vee L_3 \vee \dots)$ has
the same value as $(L_2 \vee L_3 \vee \dots)$

Therefore: Okay to shorten clauses containing false literals!

If literal L_1 is false, then clause (L_1) is false

Therefore: the empty clause means false!

DPLL version 2

Davis – Putnam – Loveland – Logemann

```
dpll_2(a, literal) {
    remove clauses from a containing literal
    shorten clauses from a containing ¬literal

    if (a contains no clauses) return true;
    if (a contains empty clause return false;

    choose P in a;
    if (dpll_2(a, ¬P)) return true;
    return dpll_2(a, P);
}
```

Partial assignment corresponding to a node is the set of chosen literals on
the path from the root to the node

Further Improvements

Formula $(L) \wedge C_2 \wedge C_3 \wedge \dots$ is only true when literal L is true

Therefore: Branch immediately on unit literals!

If literal L does not appear negated in formula F , then setting L true preserves satisfiability of F

Therefore: Branch immediately on pure literals!

Davis-Putnam algorithm

► Input: A set Σ of clauses, i.e. sets of Literals

► Output: **Yes**, if Σ is satisfiable, **no** otherwise.

```
procedure DP( $\Sigma$ )
(Sat)      if  $\Sigma = \emptyset$  then "sat"
(Empty)    if  $\Sigma$  contains the empty clause then "unsat"
(Taut)     if  $\Sigma$  contains  $p \vee \neg p$  then  $DP(\Sigma \setminus \{p \vee \neg p\})$ 
(Unit Pr)  if  $\Sigma$  has unit clause  $\{l\}$  then  $DP(\Sigma \{l = \text{true}\})$ 
(Pure)     if  $\Sigma$  has pure literal  $l$  then  $DP(\Sigma \{l = \text{true}\})$ 
(Split)    if  $DP(\Sigma \{l = \text{true}\})$  is satisfiable then "sat"
            else  $DP(\Sigma \{l = \text{false}\})$ 
```

Example: The diplomatic puzzle

We have to find a satisfying assignment for

$$(P \vee \neg Q) \wedge (Q \vee R) \wedge (\neg R \vee \neg P)$$

P	Q	R	P or $\neg Q$	Q or R	$\neg R$ or $\neg P$	Σ
T	T	T	T	T	F	F
T	T	F	T	T	T	T
T	F	T	T	T	F	F
T	F	F	T	F	T	F
F	T	T	F	T	T	F
F	T	F	F	T	T	F
F	F	T	T	T	T	T
F	F	F	T	F	T	F

Example: The diplomatic puzzle

We have to find a satisfying assignment for

$$(P \vee \neg Q) \wedge (Q \vee R) \wedge (\neg R \vee \neg P)$$

- No easy cases, so **split**. Let's say $P=1$
- **Simplify**: $(Q \vee R) \wedge (\neg R)$
- **Pure literal** Q , so $Q=1$
- **Simplify**: $(\neg R)$
- **Pure literal** (and unit clause), so $\neg R=1$, and $R=0$

Solution: $P=1, Q=1, R=0$

You can easily check that this does indeed satisfy the input formula.

Therefore: invite **Peru** and **Quatar**, but not **Rumania**

Example: The diplomatic puzzle

We could have made another split choice in the first step:

$$(P \vee \neg Q) \wedge (Q \vee R) \wedge (\neg R \vee \neg P)$$

- No easy cases, so **split**. This time let's say $R=1$
- Simplify: $(P \vee \neg Q) \wedge (\neg P)$
- Pure literal $\neg Q$, so $\neg Q=1$, so $Q=0$
- Simplify: $(\neg P)$
- Unit clause (and pure literal), $\neg P=1$, so $P=0$
- Solution: $P=0, Q=0, R=1$
- You can easily check that this also satisfies the input formula: only invite **Rumania**
- We know from the truth table method that these are the only two satisfying assignments (= the only two solutions to the puzzle)

3.4

Another example (continued)

- We had $((Q, R), (\neg R, Q))$
- Now split with $Q=1$
- Simplify: all clauses satisfied, success;

Thus: $P=1, Q=1$ satisfies the original formula.

Notice that this is independent of the value of R .

So there are two full solutions:

$(P=1, Q=1, R=1)$ and $(P=1, Q=1, R=0)$

Another example:

Find a satisfying truth assignment for

$$((P \vee \neg R) \wedge (P \vee \neg Q, R) \wedge (Q \vee R \vee \neg P) \wedge (\neg R \vee \neg P \vee Q))$$

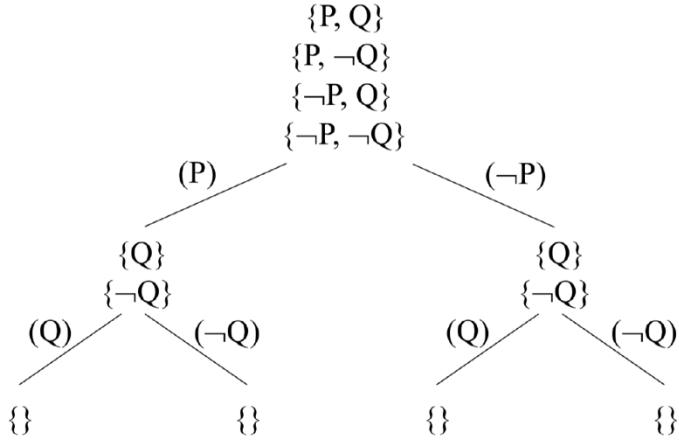
We can also write this formula as:

$$((P, \neg R), (P, \neg Q, R), (Q, R, \neg P), (\neg R, \neg P, Q))$$

- No easy case, split: let's choose $P=1$
- Simplify: $((Q, R), (\neg R, Q))$
- With naïve DPLL: let's choose $Q=0$
- Simplify: $((R), (\neg R))$
- Unit clause: $R=1$
- Simplify: false! We must backtrack to our most recent choice, which was $Q=0$, and change it to $Q=1$



Davis-Putnam as Trees



Exercise: write DP as a tree on
 $((P \vee \neg R) \wedge (P \vee \neg Q, R) \wedge (Q \vee R \vee \neg P) \wedge (\neg R \vee \neg P \vee Q))$

Variations

```

procedure DP( $\Sigma$ )
(Sat)   if  $\Sigma = \emptyset$  then "sat"
(Empty)  if  $\Sigma$  contains the empty clause then "unsat"
(Taut)   if  $\Sigma$  contains  $p \vee \neg p$  then  $DP(\Sigma \setminus \{p \vee \neg p\})$ 
(Unit Pr) if  $\Sigma$  has unit clause  $\{l\}$  then  $DP(\Sigma \{l = true\})$ 
(Pure)   if  $\Sigma$  has pure literal  $l$  then  $DP(\Sigma \{l = true\})$ 
(Split)  if  $DP(\Sigma \{l = true\})$  is satisfiable then "sat",
          else  $DP(\Sigma \{l = false\})$ 
  
```

- The **Tautology rule** only needs to be checked once at the start
- The **Pure rule** is often deleted because the cost of using it might be more than the benefits (Question: why is it so expensive?)
- The **Unit rule** is not strictly essential (if you do a clever choice for split), but it is essential for good performance
- The **Unit rule** by itself is complete for Horn clauses!

38

Heuristics for splitting

The secret sauce: how to split?

Notation:

$CP(v)$ is nr of occurrences of v
 $CN(v)$ is nr of occurrences of $\neg v$

- RAND: pick a random literal
- DLCS (Dynamic Largest Combined Sum):
 - Pick v with largest $CP(v) + CN(v)$ (= most frequent v)
 - If $CP(v) > CN(v)$ then $v=1$ else $v=0$
- DLIS (Dynamic Largest Individual Sum):
 - Pick v with largest $CP(v)$ or largest $CN(v)$
 - If $CP(v) > CN(v)$ then $v=1$ else $v=0$

Jeroslow-Wang Heuristics

- For a given literal I I compute:
$$J(I) := \sum_{l \in \omega, \omega \in \varphi} 2^{-|\omega|}$$
 Sum over 2 to the power of length of clause $\{\{A,B\}, \{-B,D\}, \{-A,D,E\}\}$
- One sided Jeroslow-Wang (JW-OS)
 - select the literal with the highest value of J
- Two-sided Jeroslow-Wang (JW-TS)
 - select the variable with the highest value of $(J(x) + J(x'))$
 - if $J(x) \geq J(x')$ set x true, else set x false
- Intuition:
 - Weight occurrences in small clauses higher

MOM's Heuristic

$\{\{A,B\}, \{-B,D\}, \{-A,D,E\}\}$

- Maximum Occurrence's in Clauses of Minimum Size
- Select the literal that maximizes the function:

$$[f^*(x) + f^*(x')] * 2^k + f^*(x) * f^*(x')$$
 - $f^*(I)$: Number of occurrences of I in the smallest non-satisfied clauses
 - k is a tuning parameter
- Intuition: Preference is given to clauses:
 - with a large number of occurrences of either x or x' in them
 - and also variables that have a large number of clauses of both phases of x in them
 - focus on the currently smallest size clauses

Sudokus in propositional logic

Checking for (in)consistency can be used to solve problems:

Many problems can be formulated as a set of constraints on the solution

These constraints can be stated in PL

Finding a solution

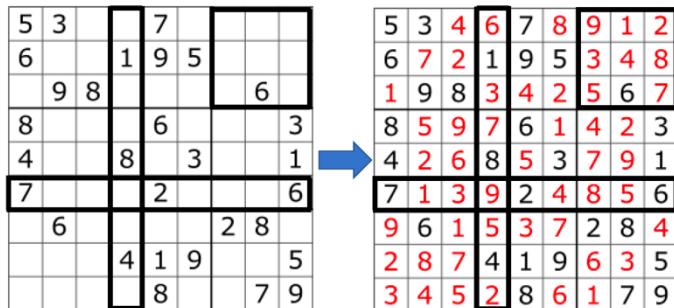
= asking if the set of constraints is satisfiable

= finding a satisfying truth assignment

Thus: Solving the problem = finding satisfying truth assignment

(that satisfies all the constraints)

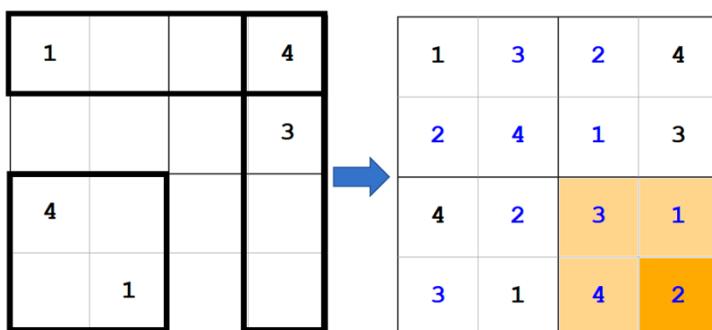
Example: solving Sudokus using PL



Formulate the problem as Constraints:

- all squares must have exactly one number from 1-9
- no number can appear twice in a row, column or square

Example: solving Sudokus using PL



Formulate the problem as Constraints:

all squares must have precisely one number from 1-4

no number can appear twice in a row, column or square

How to write the constraints in PL?

We need to describe “the world”:
one propositional variable
for each value in each position on the board

position_❶_❷_is_❸,
position_❷_❸_is_❹

in short: 144, 421,

This puzzle is described as
 $111 \wedge 144 \wedge 243 \wedge 314 \wedge 421$

in total we need
 $16 \times 4 = 64$ variables

NB: this is propositionalisation

	❶	❷	❸	❹
❶	1			4
❷				3
❸	4			
❹		1		47

How to write the constraints in PL?

We need to describe the constraints:

At least one number in each square:

“position_❶_❷ is a 1 or a 2 or 3 or a 4”
 $111 \vee 112 \vee 113 \vee 114$ (16 statements, one for each cell)

At most one number in each square: ($4 \times 16 = 64$ statements)

“if position_❶_❷ is a 1 it is not a 2 and it is not a 3 and ...”

$111 \rightarrow \neg 112 \wedge \neg 113 \wedge \neg 114$
 $112 \rightarrow \neg 111 \wedge \neg 113 \wedge \neg 114$
 $113 \rightarrow \neg 111 \wedge \neg 112 \wedge \neg 114$
 $114 \rightarrow \neg 111 \wedge \neg 112 \wedge \neg 113$

	❶	❷	❸	❹
❶	1			4
❷				3
❸	4			
❹		1		48

Logic, Hardness, GSAT, MaxSat

Knowledge Representation as Logic Engineering

- Previously , you might have learned that Logics exist , how they are defined , what their theoretical properties are, etc.
- Knowledge Representation is
- The field of using the right logic for the right AI task
- Evaluating logics w.r.t. a task or in general
- Analytically (e.g. soundness , completeness , decidability , complexity , etc
- Empirically (practical complexity , practical completeness
- Adapt existing logics for a task (we do that in the course)
- Develop new logics if needed

KR in a nutshell: (syntax, semantics, calculus)

- ▶ A possible practical problem: Politics or parties in the news
- ▶ Syntax: $\mathcal{PPP} = (PN)^*$ e.g. $PNPP \in \mathcal{PPP}$
- ▶ Semantics: (intended meaning)
 - Interpretation:
 - $(Pstr)^\mathcal{I} = (str)^\mathcal{I} + 1$
 - $(Nstr)^\mathcal{I} = (str)^\mathcal{I} - 2$
 - $(\emptyset)^\mathcal{I} = 0$
 - Truth (of equality of terms): $\models str_1 = str_2$ iff, $str_1^\mathcal{I} = str_2^\mathcal{I}$
- ▶ Calculus: $\vdash str_1 = str_2$
 - iff, $\#N(str_1) = \#N(str_2)$ and $\#P(str_1) = \#P(str_2)$

KR in a nutshell 2 (Theory about Logics)

- ▶ Properties of a logic:
 - Soundness: $\vdash str_1 = str_2$ implies $\models str_1 = str_2$
 - Completeness: $\models str_1 = str_2$ implies $\vdash str_1 = str_2$
- ▶ Our algorithm is **sound** but **incomplete**: why?
 - $\models PNP = P$, but not $\vdash PNP = P$
 - There are several sound and complete calculi for \mathcal{PPP}
- ▶ Decision problem, and decidability:
 - A terminating, sound and complete calculus (for \mathcal{PPP}) **decides** a decision problem of the logic (whether $\models str_1 = str_2$).
 - In **decidable** logics, such algorithms exist, otherwise a logic is called **undecidable**.
 - Finally, we are interested in the complexity of the logic, i.e. how difficult it is to solve the satisfiability problem.

Reasoning Tasks and Properties of Logics

Propositional Logic (reminder)

- Syntax
 - Propositions P, Q
 - Sentences S1, S2, ... based on complex operators v, -
- Semantics
 - Assignments I
 - $I(P) = \text{True or False}$
 - Assignment
 - $I(S1 \vee S2) = \text{True iff } I(S1) = \text{True or } I(S2) = \text{True}$
 - $I(\neg S) = \text{True iff } I(S) = \text{False}$
- Calculus and their properties
 - Check with Truth Tables
 - Check with DPLL

Defining Syntax of Propositional Logic

- Context-free grammars written in Bachus-Naur form (BNF) have 4 components
 - A set of terminal symbols (Words, Propositions, "True", "False", etc)
 - A set of non-terminal symbols (<sentence>, <op>, <formulas>)
 - A start symbol (*Sentence, Expr, Program*)
 - A set of rewrite rules of the form LHS ::= RHS
- Propositional Logic in BNF
 - $\langle \text{sentence} \rangle ::= p \mid T \mid F$ where $p \in \text{Propositions}$,
and T and F be used for True and False
 - $\langle \text{sentence} \rangle ::= \neg \langle \text{sentence} \rangle$
 - $\langle \text{sentence} \rangle ::= (\langle \text{sentence} \rangle \langle \text{op} \rangle \langle \text{sentence} \rangle)$
 - $\langle \text{op} \rangle ::= \wedge \mid \vee \mid \Rightarrow \mid \Leftarrow$ alternatively $\langle \text{op} \rangle ::= \& \mid \vee \mid \rightarrow \mid \Leftarrow$

Reminder: Semantics

Entailment, Validity and satisfiability

- A sentence α is **entailed by** a KB K iff α holds/is true in all models of KB. We write $\text{KB} \models \alpha$
- A sentence is **satisfiable** if it is true in *some* model
—e.g., $A \vee B, C$
- A sentence is **unsatisfiable** if it is true in *no* models

Semantics versus Calculus

- A sentence S is **satisfiable** if it is true in *some* model (we write $\models S$)
- DPLL(S) returns SAT or UNSAT (we write $\vdash S$)

Intention is that these two things deliver the same things!

DPLL complexity

- DPLL determines **satisfiability** of formulas in CNF, i.e.
 - dpll(KB)=sat iff Knowledge Base KB is satisfiable
- ▶ The worst case complexity of the DP algorithm $O(1,618^n)$.
- ▶ This is an improvement!... Notice that, for example,

$$\begin{array}{rcl} 2^{50} & = & 1.125.899.906.842.624 \\ 1,618^{50} & = & 28.114.208.661 \end{array}$$

DEFINITION: Properties of inference

- If an algorithm \vdash proves **only satisfiable** sentences S it is **sound**, i.e. if
 - $\vdash S$ also $\vDash S$
- It is **complete** if it derives **all satisfiable** sentence, i.e.
 - $\vDash S$ implies $\vdash S$
- A calculus **terminates** if it finds a model in finite time
- A logic is **decidable** if there is a sound and complete calculus that terminates

Properties of DPLL based SAT solving

Our calculus

1. $\vdash S$: DPLL(S) is **sound**
2. $\vdash S$: DPLL(S) is **complete**
3. $\vdash S$: DPLL(S) **terminates**

Which implies that $\text{KB} \vdash \alpha$ is a decision procedure for SAT solving of propositional logic.

Proof(sketch) for DPLL based Entailment

Our calculus

1. $\vdash S : \text{DPLL}(S)$ is **sound**
2. $\vdash S : \text{DPLL}(S)$ is **complete**
3. $\vdash S : \text{DPLL}(S)$ **terminates**

Proof sketch: Show that:

1. If $\models S$ then $\models S$. $\vdash S$ means that DPLL returns SAT, i.e. it found a partial assignment. This assignment is a model for S .
2. If $\models S$, S is satisfiable, i.e. has a model M . This model M will be an assignment found by DPLL, as DPLL tries all combinations.
3. $\vdash S$ assignes, in the worst case, truth values True and False to all propositional variables in S . There are only finitely many variables.

Proof of Entailment by Refutation

We define a procedure **KB** $\vdash a$ {

1. Translate KB into CNF to get $\text{cnf}(\text{KB})$
2. Translate $\neg a$ into CNF to get $\text{cnf}(\neg a)$
3. Add $\text{cnf}(\neg a)$ to $\text{cnf}(\text{KB})$
4. Apply SAT solver (such as DPLL) until
 - either model is found (satisfiable) or
 - search exhausted (unsatisfiable)
5. If satisfiable, return "FALSE", otherwise "TRUE".

}

This works in Propositional Logic, not in all logics

Semantics versus Calculus

- **KB** $\models \alpha$: A sentence α is **entailed by** a KB K iff α holds/is true in all models of KB.
- **KB** $\vdash \alpha$: $\text{DPLL}(\text{cnf}(\text{KB}) + \text{cnf}(\{\neg a\}))$ returns UNSAT

Intention is that these two things deliver the same things!

Properties DPLL based Entailment

Our calculus

- $\text{KB} \vdash \alpha : \text{DPLL}(\text{cnf}(\text{KB}) + \text{cnf}(\{\neg a\}))$ is **sound**
- $\text{KB} \vdash \alpha : \text{DPLL}(\text{cnf}(\text{KB}) + \text{cnf}(\{\neg a\}))$ is **complete**
- $\text{KB} \vdash \alpha : \text{DPLL}(\text{cnf}(\text{KB}) + \text{cnf}(\{\neg a\}))$ **terminates**

Which implies that $\text{KB} \vdash \alpha$ is a decision procedure (or decides) entailment of propositional logic $\text{KB} \models a$

.

Propositionalizing

Propositional logic is a weak language

- Propositional letters only describe complete “states” of the world, we cannot talk about “individuals” (e.g., Mary, 3)
- Can’t directly talk about properties of individuals or relations between individuals (e.g., “Bill is tall”)
- Generalizations, patterns, regularities can’t easily be represented (e.g., “all triangles have 3 sides”)
- First Order Logic (abbreviated FOL or FOPC) is expressive enough to concisely represent this kind of information

FOL adds relations, variables, and quantifiers, e.g.,

- “Every elephant is $\forall x$ (elephant(x) gray(x))
- “There is a white $\exists x$ (alligator(X) ^ white(X))

Propositionalising First Order Logic over finite domains

$\forall x \text{ Smokes}(x) \Rightarrow \text{Cancer}(x)$

$\forall x, y \text{ Friends}(x, y) \Rightarrow (\text{Smokes}(x) \Leftrightarrow \text{Smokes}(y))$

2 individuals Alice & Bob

smokesA, smokesB, cancerA, cancerB
friendsAB, friendsBA, friendsAA, friendsBB

smokesA \Rightarrow cancerA,
smokesB \Rightarrow cancerB
friendsAB \Rightarrow (smokesA \Leftrightarrow smokesB)
friendsBA \Rightarrow (smokesB \Leftrightarrow smokesA)
friendsAA \Rightarrow (smokesA \Leftrightarrow smokesA)
friendsBB \Rightarrow (smokesB \Leftrightarrow smokesB)

Problem solving with Propositionalisation

In order to solve a problem using Propositionalisation

1. formulate problem as a complex formula F in first order logic
2. propositionalise F to get a Propositional statement S
3. find a “satisfying assignment” of truth values (in other words: a model of S)
4. This model gives us the solution to the original problem.

Remember:

if S contains n symbols, we need to check $\sim 1,681^n$ nodes

Complexity: if S is the result of propositionalising a FOL sentence, then n is exponential in the number of individuals i in the domain

--> DPLL may have to check doubly exponentially many nodes: $1,681^{2n}$

n=10 gives a truth table of size 10230

Hard Satisfiability Problems

Structured vs. Random Problems

- So far, we've been dealing with SAT problems that encode other problems
- Most not as hard as # of variables & clauses suggests
- Small crossword grid + medium sized dictionary may turn into a big formula ... but still a small puzzle at some level
- Unit propagation does a lot of work for you
- Clause learning picks up on the structure of the encoding
- But some random SAT problems really are hard!
- Fix the length (often 3)
- Vary the ratio between number of clauses l and number of variables n

The magic ratio 4.26

Complexity peak is very stable ...

across problem sizes

across solver types

systematic (last lecture)

stochastic (this lecture)

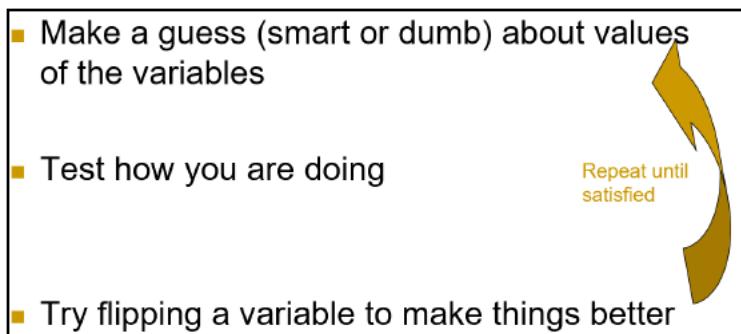
That's called a "phase transition"

- Problems $< 0^{\circ}\text{C}$ are like ice; $> 0^{\circ}\text{C}$ are like water
- Similar "phase transitions" for other NP hard problems
- job shop scheduling
- traveling salesperson
- exam timetables
- Boolean circuit synthesis
- Latin squares (alias sports scheduling)

Incomplete SAT methods (GSAT)

Local search for SAT

- Make a guess (smart or dumb) about values of the variables
 - Test how you are doing
 - Try flipping a variable to make things better
 - Algorithms differ on which variable to flip



- Flip a randomly chosen variable?
 - No, blundering around blindly takes exponential time
 - Ought to pick a variable that improves ...
 - Increase the # of satisfied clauses as much as possible
 - Break ties randomly
 - Note: Flipping a var will repair some clauses and break others
 - This is the “GSAT” (Greedy SAT) algorithm (almost)
 - Flip most-improving variable. Guaranteed to work?
 - What if first guess is A=1, B=1, C=1?
 - 2 clauses satisfied
 - Flip A to 0 → 3 clauses satisfied
 - Flip B to 0 → all 4 clauses satisfied (pick one)
 - Flip C to 0 → 3 clauses satisfied
 - Flip most-improving variable. Guaranteed to work?
 - But what if first guess is A=0, B=1, C=1?
 - 3 clauses satisfied
 - Flip A to 1 → 3 clauses satisfied
 - Flip B to 0 → 3 clauses satisfied
 - Flip C to 0 → 3 clauses satisfied
 - Pick one anyway ... (picking A wins on next step)

- Flip most-improving variable. Guaranteed to work?
- Yes for 2-SAT: probability $\rightarrow 1$ within $O(n^2)$ steps
- No in general: can & usually does get locally stuck
- Therefore, GSAT just restarts periodically
 - New random initial guess

The GSAT Procedure

► This algorithm is due to Selman, Levesque and Mitchell

- Adds restarts to the simple "greedy" algorithm, and also allows sideways flips (not necessarily increasing the "cost function.")

```

procedure GSAT(Sigma)
  for i := 1 to MAX-TRIES ; These are the restarts
    T := random(Sigma) ; random assignment
    for j := 1 to MAX-FLIPS ; To ensure termination
      if T satisfies Sigma then return T
      else T := T with variable flipped to maximize
            number of satisfied clauses
      ; It doesn't matter if the number does
      ; not increase. This are the sideways flips
    end
  end

```

Believe it or not, GSAT outperformed everything else when it was invented in 1992.

Decidability? Completeness?

- What does that mean for decidability?
- What does that mean for completeness?
- Remember: Completeness means that all correct solutions are found
- Soundness: If our algorithm returns a variable assignment with input KB, this variable assignment is a model of KB
- Completeness: If an assignment is a model, then our algorithm will find it.

Who cares about decidability?

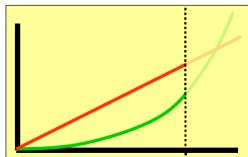
- Decidability \approx completeness guarantee to find an answer, or tell you it doesn't exist, given enough run time & memory
- Sources of incompleteness
- incompleteness of the input data
- insufficient run time to wait for the answer
- ④ Completeness is unachievable in practice anyway, regardless of the completeness of the algorithm

Who cares about undecidability?

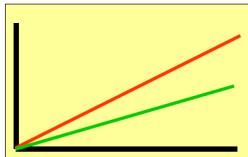
- Undecidability ≠ always guaranteed not to find an answer
- Undecidability = not always guaranteed to find an answer
- Undecidability may be harmless in many cases; in all cases that matter

Who cares about complexity?

- **worst-case**: may be exponentially rare
- **Asymptotic $O(n)$ vs $O(n^2)$**



- **ignores constants**



Clause Learning, Planning and Applications

Clause Learning and Dependency directed backtracking

Chronological Backtracking

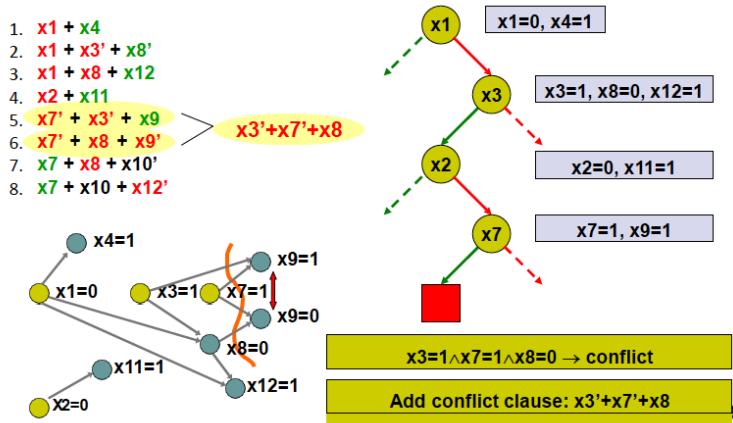
- ☒ Backtracking to the highest decision level that has not been tried with both values
- ☒ Originally proposed in the DLL paper in 1962
- ☒ OK for randomly generated instances, bad for instances generated in practical applications
- ☒ We can do better than that

Conflict Driven Learning and Non-chronological Backtracking

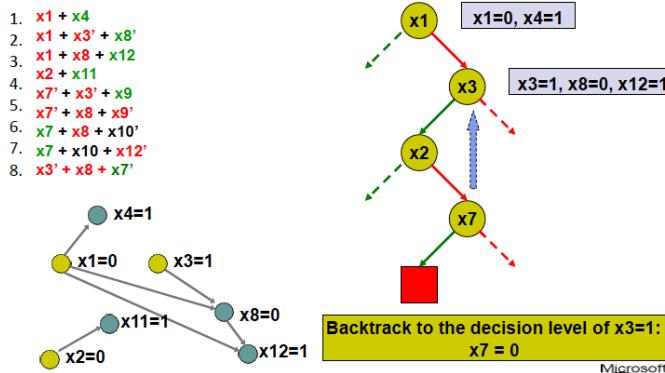
1. $x_1 + x_4$
2. $x_1 + x_3' + x_8'$
3. $x_1 + x_8 + x_{12}'$
4. $x_2 + x_{11}$
5. $x_7' + x_3' + x_9$
6. $x_7' + x_8 + x_9'$
7. $x_7 + x_8 + x_{10}'$
8. $x_7 + x_{10} + x_{12}'$

Notation: + is "or" (\vee), a ' is "not", \neg

So: $x_1 + x_3' + x_8' = x_1 \vee \neg x_3 \vee \neg x_8$



DLL with Non-Chronological Backtracking and Learning



Summarizing

We have now gained two advantages:

1. We have “learned” a new clause that “summarizes” a potential conflict. Adding this clause to the clause base will help us to avoid making the same mistake in other parts of the search space.
2. We can make a bigger “jump back” in the search tree, replacing chronological backtracking (= just undoing the most recent choice), with dependency directed backtracking (= undoing the most recent choice that was involved in the conflict, ie the most recent choice that mattered).

Chapter 7: Propositional Satisfiability Techniques

- 1 Propositional satisfiability: given a boolean formula

» e.g., $(P \vee Q) \wedge (\neg Q \vee R \vee S) \wedge (\neg R \vee \neg P)$,

does there exist a model

» i.e., an assignment of truth values to the propositions that makes the formula true?

1 approach:

Try translating classical planning problems into satisfiability problems, and solving them that way

State Transition System

$$\Sigma = (S, A, \gamma)$$

1 states $S = \{ \dots, s_i, \dots \}$

1 actions $A = \{ \dots, a_i, \dots \}$

1 State-transition function

$$\gamma: S \times A \rightarrow 2^S$$

$$u S = \{s_0, \dots, s_5\}$$

$$u A = \{\text{move1, move2, put, take, load, unload}\}$$

u γ : see the arrows

Planning Problem

$$P = (\Sigma, s_0, s_G)$$

- Description of Σ
- Initial state or set of states

Initial state = s_0

- Objective

Goal state, set of goal states, set of tasks, “trajectory” of states, objective function, ...

Goal state = s_5

Plans

u **Classical plan**: a sequence of actions

$$\pi = \langle a_0, a_1, \dots, a_{n-1} \rangle$$

$\langle \text{take, move1, load, move2} \rangle$

Policy: partial function from S into A

$\{(s_0, \text{take}),$

$(s_1, \text{move1}),$

$(s_3, \text{load}),$

$(s_4, \text{move2})\}$

Solution length = 4

Overall Approach

1 A *bounded planning problem* is a pair (P, n) :

u P is a planning problem; n is a positive integer

u Any solution for P of length n is a solution for (P, n)

1 Planning algorithm: Do iterative deepening

u for $n = 0, 1, 2, \dots,$

» encode (P, n) as a satisfiability problem Φ

» if Φ is satisfiable, then

- From the set of truth values that satisfies Φ , a solution plan can be constructed, so return it and exit

Notation

- 1 For satisfiability problems we use propositional logic
- 1 Need to encode ground atoms into propositions
 - » Atom: $\text{at}(r1, \text{loc}1)$
 - » Proposition: $\text{at-}r1\text{-loc}1$
- u But we won't bother to do a syntactic rewrite
Just use $\text{at}(r1, \text{loc}1)$ itself as the proposition
- 1 Also, we'll write plans starting at a_0
 - u $\pi = \langle a_0, a_1, \dots, a_{n-1} \rangle$

Fluents

- 1 If $\pi = \langle a_0, a_1, \dots, a_{n-1} \rangle$ is a solution for (P, n) , it generates these states:
 $s_0, s_1 = \gamma(s_0, a_0), s_2 = \gamma(s_1, a_1), \dots, s_n = \gamma(s_{n-1}, a_{n-1})$
- 1 *Fluent*: proposition saying a particular atom is true in a particular state
 - u $\text{at}(r1, \text{loc}1, i)$ is a fluent that's true iff $\text{at}(r1, \text{loc}1)$ is in s_i
 - u We'll use l_i to denote the fluent for literal l in state s_i
 - » e.g., if $l = \text{at}(r1, \text{loc}1)$
then $l_i = \text{at}(r1, \text{loc}1, i)$
 - u a_i is a fluent saying that a is the i 'th step of π
 - » e.g., if $a = \text{move}(r1, \text{loc}2, \text{loc}1)$
then $a_i = \text{move}(r1, \text{loc}2, \text{loc}1, i)$

Encoding Planning Problems

- 1 Encode (P, n) as a formula Φ such that
 $\pi = \langle a_0, a_1, \dots, a_{n-1} \rangle$ is a solution for (P, n) if and only if
 Φ can be satisfied in a way that makes the fluents a_0, \dots, a_{n-1} true
 - 1 Let
 - u $A = \{\text{all actions in the planning domain}\}$
 - u $S = \{\text{all states in the planning domain}\}$
 - u $L = \{\text{all literals in the language}\}$
 - 1 Φ is the conjunction of many formulas ...

Formulas in Φ

- 1 Formula describing the initial state:
 $\Lambda\{l_0 \mid l \in s_0\} \wedge \Lambda\{\neg l_0 \mid l \in L - s_0\}$
- 1 Formula describing the goal:
 $\Lambda\{l_n \mid l \in g^+\} \wedge \Lambda\{\neg l_n \mid l \in g^-\}$
- 1 For every action a in A , formulas describing what changes a would make if it were the i 'th step of the plan:
 - u $a_i \Rightarrow \Lambda\{p_i \mid p \in \text{Precond}(a)\} \wedge \Lambda\{e_{i+1} \mid e \in \text{Effects}(a)\}$
- 1 Complete exclusion axiom:

- u For all actions a and b , formulas saying they can't occur at the same time
 $\neg a_i \vee \neg b_i$
 - u this guarantees there can be only one action at a time
- 1 Is this enough?

Frame Axioms

- 1 *Frame axioms:*
 - u Formulas describing what *doesn't* change between steps i and $i+1$
- 1 Several ways to write these
- 1 One way: *explanatory frame axioms*
 - u One axiom for every literal l
 - u Says that if l changes between s_i and s_{i+1} , then the action at step i must be responsible:
$$(\neg l_i \wedge l_{i+1} \Rightarrow \forall_{a \in A} \{a_i \mid l \in \text{effects}^+(a)\})$$

$$\wedge (l_i \wedge \neg l_{i+1} \Rightarrow \forall_{a \in A} \{a_i \mid l \in \text{effects}^-(a)\})$$

Example

- 1 Planning domain:
 - u one robot r1
 - u two adjacent locations l1, l2
 - u one operator (move the robot)
- 1 Encode (P, n) where $n = 1$
 - u Initial state: {at(r1,l1)}
 - Encoding: at(r1,l1,0) \wedge \neg at(r1,l2,0)
 - u Goal: {at(r1,l2)}
 - Encoding: at(r1,l2,1) \wedge \neg at(r1,l1,1)
 - u Operator: see next slide
- 1 Operator: move(r,l,l')
 - precond: at(r,l)
 - effects: at(r,l'), \neg at(r,l)
- Encoding:
 - move(r1,l1,l2,0) \Rightarrow at(r1,l1,0) \wedge at(r1,l2,1) \wedge \neg at(r1,l1,1)
 - move(r1,l2,l1,0) \Rightarrow at(r1,l2,0) \wedge at(r1,l1,1) \wedge \neg at(r1,l2,1)
 - move(r1,l1,l1,0) \Rightarrow at(r1,l1,0) \wedge at(r1,l1,1) \wedge \neg at(r1,l1,1) {contradiction}
 - move(r1,l2,l2,0) \Rightarrow at(r1,l2,0) \wedge at(r1,l2,1) \wedge \neg at(r1,l2,1) {easy to detect}
 - move(l1,r1,l2,0) \Rightarrow ... {
 - move(l2,l1,r1,0) \Rightarrow ... non-essential
 - move(l1,l2,r1,0) \Rightarrow ... }
 - move(l2,l1,r1,0) \Rightarrow ... }
- 1 How to avoid generating the last four actions?
 - u Assign data types to the constant symbols: order-sorted logic
- 1 Locations: l1, l2
- 1 Robots: r1

1 Operator: $\text{move}(r : \text{robot}, l : \text{location}, l' : \text{location})$
 precond: $\text{at}(r,l)$
 effects: $\text{at}(r,l'), \neg\text{at}(r,l)$

Encoding:

$\text{move}(r1,l1,l2,0) \Rightarrow \text{at}(r1,l1,0) \wedge \text{at}(r1,l2,1) \wedge \neg\text{at}(r1,l1,1)$
 $\text{move}(r1,l2,l1,0) \Rightarrow \text{at}(r1,l2,0) \wedge \text{at}(r1,l1,1) \wedge \neg\text{at}(r1,l2,1)$

1 Complete-exclusion axiom:
 $\neg\text{move}(r1,l1,l2,0) \vee \neg\text{move}(r1,l2,l1,0)$

1 Explanatory frame axioms:

$\neg\text{at}(r1,l1,0) \wedge \text{at}(r1,l1,1) \Rightarrow \text{move}(r1,l2,l1,0)$
 $\neg\text{at}(r1,l2,0) \wedge \text{at}(r1,l2,1) \Rightarrow \text{move}(r1,l1,l2,0)$
 $\text{at}(r1,l1,0) \wedge \neg\text{at}(r1,l1,1) \Rightarrow \text{move}(r1,l1,l2,0)$
 $\text{at}(r1,l2,0) \wedge \neg\text{at}(r1,l2,1) \Rightarrow \text{move}(r1,l2,l1,0)$

Extracting a Plan

- 1 Suppose we find an assignment of truth values that satisfies Φ .
 - u This means P has a solution of length n (*as used in the encoding*)
- 1 For $i=1,\dots,n$, there will be exactly one action a such that $a_i = \text{true}$
 - u This is the i 'th action of the plan.
- 1 Example (from the previous slides):
 - u Φ can be satisfied with $\text{move}(r1,l1,l2,0) = \text{true}$
 - u Thus $\langle \text{move}(r1,l1,l2,0) \rangle$ is a solution for $(P,0)$
 - » It's the only solution - no other way to satisfy Φ

Planning

- 1 How to find an assignment of truth values that satisfies Φ ?
 - u Use a satisfiability algorithm:
 - u Davis-Putname
 - u Local Search
 - u GSAT
 - u ...

Discussion

- 1 Recall the overall approach:
 - u for $n = 0, 1, 2, \dots,$
 - » encode (P,n) as a satisfiability problem Φ
 - » if Φ is satisfiable, then
 - From the set of truth values that satisfies Φ , extract a solution plan and return it
 - How well does this work?
 - u By itself, not very practical (takes too much memory and time)
 - u But it can be combined with other techniques
 - » e.g., planning graphs

SatPlan

- 1 SatPlan combines planning-graph expansion and satisfiability checking, roughly as follows:

- u for $k = 0, 1, 2, \dots$
 - » Create a planning graph that contains k levels
 - » Encode the planning graph as a satisfiability problem
 - » Try to solve it using a SAT solver
 - If the SAT solver finds a solution within some time limit,
 - Remove some unnecessary actions
 - Return the solution
 - Memory requirement still is combinatorially large
 - u but less than what's needed by a direct translation into satisfiability
- 1 BlackBox (predecessor to SatPlan) was one of the best planners in the 1998 planning competition
- 1 SatPlan was one of the best planners in the 2004 and 2006 planning competitions

Exponential Complexity Growth: The Challenge of Complex Domains

SAT Encoding (automatically generated from problem specification)

Applications of SAT (1)

- Model Checking (hardware/software verification)
- Bounded model checking (BMC)
- Enhancement techniques: Abstraction and refinement
- Major groups at Intel, IBM, Microsoft, and universities such as CMU, Cornell, and Princeton
- **SAT has become a dominant back-end technology**

Key questions:

- Safety property: can system ever reach state S?
- Liveness property: will system always reach T after it gets to S?

Main idea:

- Encode “reachability” in a finite transition graph
- CNF encoding: in simple terms, “similar” to planning

Classical Planning

- Parallel step optimal plans
- E.g., SATPLAN-06 fastest in this category at IPC-06 planning competition

Main idea:

- Create planning graph, unfolded up to T steps
- Encode legal transitions as CNF
 - **Variables**: action vars, fact vars for each time step t
 - **Clauses** encode action preconditions, effects, mutual exclusion, etc.
- e.g., (actionA_(t+1) AND preA1_t AND preA2_t AND preA3_t),
 (NOT actionA_t OR NOT actionB_t), ...

Combinatorial Design

- Complex mathematical structures with desirable properties
- Latin squares

- Quasi groups
- Ramsey numbers
- Steiner systems
- ...
- Highly useful for
- design of experiments
- coding theory, cryptography
- drug design and testing
- crop rotation schedules

Has been applied to solve sub-problems in many domains!

- Test pattern generation
- Scheduling
- Optimal Control
- Protocol Design (e.g., for networks)
- Multi-agent systems
- E-Commerce (E-auctions and electronic trading agents), etc.

Propositional Logic, MaxSat + Exam Preparation

The MAXSAT Problem

- The well-known SAT problem is to determine if a boolean formula in Conjunctive Normal Form (CNF) has a satisfying truth assignment
- A CNF formula is a conjunction of clauses
- A clause is a disjunction of literals (variables or their negations)
- \perp denotes the empty clause (falsified by every truth assignment)
- A satisfying assignment assigns true to at least one literal of every clause.
- MAXSAT is an optimization extension of SAT that asks what is the maximum number of clauses that can be simultaneously satisfied

Example

$$F = (\neg x) \vee (x \wedge y) \vee (\neg y) \vee (z \wedge w)$$

F is unsatisfiable, so no truth assignment can satisfy all 4 clauses.

Truth assignment $u = \{x, \neg y, z, \neg w\}$ satisfies all clauses except $(\neg x)$.)

$U(\text{pie})$ is a solution to the MAXSAT problem F.

Some clauses may be more important to satisfy than others

- This can be modeled by associating a positive **cost** with each clause C that will be incurred if C is falsified

If it is mandatory to satisfy C, its cost is / and C is called hard

- Otherwise, C is called soft

Example

$$F = (\neg x, /(\infty)) \vee (x \wedge 2y, 4) \vee (\neg y, 1) \vee (z \wedge 2w, /(\infty))$$

In F, $(\neg x, /)$ is a hard clause, and $(x \wedge 2y, 4)$ is soft with cost 4.

A truth assignment u has cost equal to the sum of the costs of the clauses it falsifies

- Goal: find an optimal feasible truth assignment, i.e., a truth assignment of minimum finite cost $\text{mincost}(F)$

Example

$F = (\neg x, \infty) \vee (x, 2y, 4) \vee (\neg y, 1) \vee (z, 2w, \infty)$

$u = \{\neg x, y, z, \neg w\}$ satisfies all clauses except $(\neg y, 1)$

u is optimal: $\text{mincost}(F) = \text{cost}(u) = 1$.

We use $\text{cost}(C)$ to denote the cost of clause C .

- A solution must satisfy all hard clauses (else its cost will be infinite).
- A solution also satisfies a maximum total cost of soft clauses.
- Casting as minimization problem more closely corresponds to how most MAXSAT solvers work.
- Many solutions might exist—typically we are only interested in finding one, sometimes only interested in finding out the cost of a solution.

MAXSAT (ms) (standard MAXSAT): no hard clauses and all clause have weight 1.

- Solution maximizes the number of satisfied clauses.
- **Weighted MAXSAT (wms)**: no hard clauses.
- **Partial MAXSAT (pms)**: have hard clauses but all soft clauses have weight 1.
- **Weighted Partial MAXSAT (wpms)**: the version we have defined here (subsumes all other versions).
- Standard MAXSAT is most interesting for theory: it already has sufficient structure for theoretical insights.
- Other versions mostly an artifact of the limitations of earlier MAXSAT solvers.

Exam Preparation (I tried to answer them, emphasis on tried...)

1)

In a proof by contradiction, such as DPLL, I can prove that a formula F is entailed by a knowledge base KB by showing that

A the knowledge base KB and the formula F are together unsatisfiable.

B the knowledge base is unsatisfiable, which implies that the formula F must be entailed.

C the knowledge base KB and the negation of the formula are unsatisfiable.

D the formula F is unsatisfiable, which implies that it must be entailed by the KB .

2)

Logic Engineering

Consider a language L defined as follows

- Syntax

{A,B,C} L

If F1, F2, F3 L then "(F1*F2*F3)" L

•Semantics: let I be an interpretation

$I(A)=1, I(B)=2, I(C)=3$

$I(F1*F2*F3) = I(F3).I(F1).I(F2)$

What is $I(A*B*(C*A*B))$?

3)

Semantics: the meaning of expressions

Fill in the truth value of the following formulas under the assumption that A,B,C and D are all false. Fill in the truth values of the formulas (use the words "True, False").

$(A \& B) v -C v -D$ is : **True**

Conjunction, or, disjunction

$(A \& B) v (-C \& D)$ is **True**

Conjunction, or, conjunction

$(A \rightarrow B) v (C \rightarrow D)$ is **True**

Implies, or, implies _____

$(-A \rightarrow B) v (-C \rightarrow D)$ is **False**

4)

Semantics 2

Which of the following is true? The propositional statement

$(P \vee Q) \rightarrow (P \& Q)$

To be either false or true (truth assignment)

contingency: neither a tautology (true and valid) nor a contradiction (false)

A is satisfiable, but not valid

B is valid

C is a contradiction

D neither valid, satisfiable nor a contradiction

5)

Automated Reasoning as Search

Automated reasoning can often be seen as search procedure, that tries to find an order of rules to formulas in a knowledge base.

Which of the following statements about the DPLL (Davis Putnam Logemann Loveland) Algorithm and Search is correct?

A DPLL recursively searches through all possible variable assignments for a model, i.e. an interpretation that satisfies all the clauses.

B DPLL iteratively searches through the set of all clauses for one that is satisfied by a given interpretation.

C DPLL exhaustively applies a set of transformation rules to produce a contradiction

D DPLL cannot be seen as a search procedure, it only randomly assigns values to propositions until it finds an assignment that satisfies all the clauses

6)

DPLL

Consider the following knowledge base (set of clauses):

Definition 2.1.1. A tautology is a proposition that is always true.

Example 2.1.1. $p \vee \neg p$

Definition 2.1.2. A contradiction is a proposition that is always false.

Example 2.1.2. $p \wedge \neg p$

Definition 2.1.3. A contingency is a proposition that is neither a tautology nor a contradiction.

Example 2.1.3. $p \vee q \rightarrow \neg r$

1: $\neg p \vee q \vee r$
2: $p \vee \neg q \vee \neg r$
3: $\neg p \vee \neg q \vee \neg r$
4: $p \vee q$
5: $p \vee r$
6: $\neg p \vee \neg q \vee r$

and show whether the knowledge base is satisfiable or not using DPLL. Use the following order: unit, pure, split alphabetic with a positive value.

Assign $p = \text{true}$

2: $p \vee \neg q \vee \neg r$

4: $p \vee q$

5: $p \vee r$

Assign q=true

1: $\neg p \vee q \vee r$

4: $p \vee q$

Find unit clause 6: $\neg p \vee \neg q \vee r$, so r = true

3: $\neg p \vee \neg q \vee \neg r$

7)

Properties of GSAT

GSAT is complete w.r.t Propositional Logic satisfiability. True or false?

GSAT: Local greedy search/algorithm of picking the next variable to flip which increases the number of satisfied clauses the most (ties are broken randomly, note that flips will break also some clauses).

A True

B False

Probabilistic Graphical Models 1

Foundations

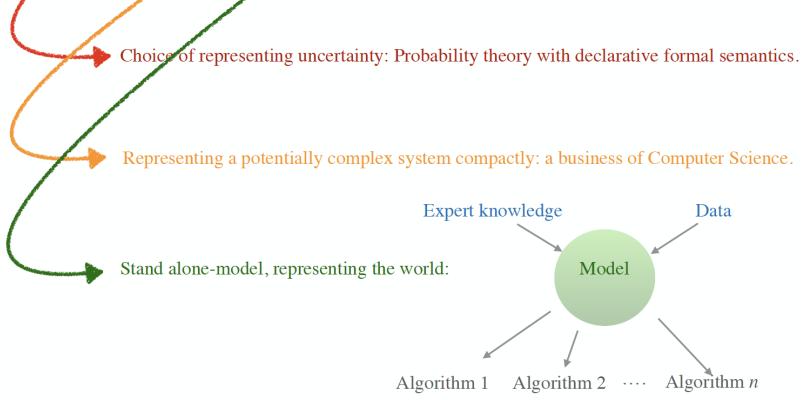
Introduction

Probabilistic Graphical Models is a major field in AI at the crossroads of KR, Statistics and Learning.

- It has vast amount of applications:
- Decision Making
- Image Segmentation
- Medical Diagnoses
- Speech Recognition
- Natural Language Processing

MOTIVATION

- Probabilistic Graphical Models is a major field in AI at the crossroads of KR, Statistics and Learning.



DEGREES OF BELIEFS

- (Classical) Knowledge base: a coarse (binary!) classification of the world:
 - Possible
 - Not Possible
- We can do better(!): a finer classification of sentences through a finer classification of worlds.
 - A *degree of belief* or *probability* in $[0, 1]$ assigned to each world $\omega : \Pr(\omega)$
- Then, the probability/belief of/in, for a given sentence $\alpha : \Pr(\alpha) \stackrel{\text{def}}{=} \sum_{\omega \models \alpha} \Pr(\omega)$ (B1)
which is the sum of probabilities assigned to worlds that α is true.

DEGREES OF BELIEF

Example:

world	Earthquake	Burglary	Alarm	Pr(.)
ω_1	true	true	true	.0190
ω_2	true	true	false	.0010
ω_3	true	false	true	.0560
ω_4	true	false	false	.0240
ω_5	false	true	true	.1620
ω_6	false	true	false	.0180
ω_7	false	false	true	.0072
ω_8	false	false	false	.7128

A state of belief, is also called a joint probability distribution.

$$\begin{aligned}
 \Pr(\text{Earthquake}) &= \Pr(\omega_1) + \Pr(\omega_2) + \Pr(\omega_3) + \Pr(\omega_4) = .1 \\
 \Pr(\text{Burglary}) &= .2 \\
 \Pr(\neg\text{Burglary}) &= .8 \\
 \Pr(\text{Alarm}) &= .2442
 \end{aligned}$$

PROPERTIES OF BELIEF

- $0 \leq \Pr(\alpha) \leq 1$ for any sentence α .
- $\Pr(\alpha) = 0$ when α is inconsistent.
- $\Pr(\alpha) = 1$ when α is valid.
- $\Pr(\alpha) + \Pr(\neg\alpha) = 1$.
- $\Pr(\alpha \vee \beta) = \Pr(\alpha) + \Pr(\beta) - \Pr(\alpha \wedge \beta)$.

PROPERTIES OF BELIEF (PARTITIONED WORLDS)

$$\Pr(\alpha) + \Pr(\neg\alpha) = 1.$$

Example:

world	Earthquake	Burglary	Alarm	$\Pr(\cdot)$
ω_1	true	true	true	.0190
ω_2	true	true	false	.0010
ω_3	true	false	true	.0560
ω_4	true	false	false	.0240
ω_5	false	true	true	.1620
ω_6	false	true	false	.0180
ω_7	false	false	true	.0072
ω_8	false	false	false	.7128

Observation: The

set of all worlds.

$$\Pr(\text{Burglary}) = \Pr(\omega_1) + \Pr(\omega_2) + \Pr(\omega_5) + \Pr(\omega_6) = .2$$

$$\Pr(\neg\text{Burglary}) = \Pr(\omega_3) + \Pr(\omega_4) + \Pr(\omega_7) + \Pr(\omega_8) = .8$$



PROPERTIES OF BELIEF (PARTITIONED WORLDS)

$$\Pr(\alpha \vee \beta) = \Pr(\alpha) + \Pr(\beta) - \Pr(\alpha \wedge \beta).$$

$$\alpha \quad \curvearrowright \quad \beta$$

Example:

world	Earthquake	Burglary	Alarm	$\Pr(\cdot)$
ω_1	true	true	true	.0190
ω_2	true	true	false	.0010
ω_3	true	false	true	.0560
ω_4	true	false	false	.0240
ω_5	false	true	true	.1620
ω_6	false	true	false	.0180
ω_7	false	false	true	.0072
ω_8	false	false	false	.7128

Observation: Th

ose those satisfying $\alpha \wedge \neg\beta$,

$$\Pr(\text{Earthquake}) = \Pr(\omega_1) + \Pr(\omega_2) + \Pr(\omega_3) + \Pr(\omega_4) = .1$$

$$\Pr(\text{Burglary}) = \Pr(\omega_1) + \Pr(\omega_2) + \Pr(\omega_5) + \Pr(\omega_6) = .2$$

$$\Pr(\text{Earthquake} \wedge \text{Burglary}) = \Pr(\omega_1) + \Pr(\omega_2) = .02$$

$$\Pr(\text{Earthquake} \vee \text{Burglary}) = .1 + .2 - .02 = .28$$



PROPERTIES OF BELIEF (PARTITIONED WORLDS)

$$\Pr(\alpha \vee \beta) = \Pr(\alpha) + \Pr(\beta) - \Pr(\alpha \wedge \beta).$$

OBSERVATIONS:

- $\Pr(\alpha \vee \beta) = \Pr(\alpha) + \Pr(\beta)$ when α and β are mutually exclusive.
- In this case, there is no world that satisfies both α and β . Hence, $\alpha \wedge \beta$ is inconsistent and $\Pr(\alpha \wedge \beta) = 0$.

Belief Dynamics

UPDATING BELIEFS

- Recall that:

<i>world</i>	Earthquake	Burglary	Alarm	$\Pr(\cdot)$
ω_1	true	true	true	.0190
ω_2	true	true	false	.0010
ω_3	true	false	true	.0560
ω_4	true	false	false	.0240
ω_5	false	true	true	.1620
ω_6	false	true	false	.0180
ω_7	false	false	true	.0072
ω_8	false	false	false	.7128

- Now assume that we know **alarm is triggered** i.e., **Alarm = true** (also called **evidence**).
- Yet, this knowledge is **not compatible** with the belief state above (i.e., 0.2442) alarm being true.

- Given evidence β , our goal is to update the state of belief $\Pr(\cdot)$ to $\Pr(\cdot|\beta)$

- Hence we expect $\Pr(\cdot|\beta)$ to assign a belief of 1 to β : $\Pr(\beta|\beta) = 1$.

- This implies $\Pr(\neg\beta|\beta) = 0$ and, hence,

every world ω that satisfies $\neg\beta$ must be assigned the belief 0:

$$\Pr(\omega|\beta) = 0 \quad \text{for all } \omega \models \neg\beta.$$

- To define $\Pr(\cdot|\beta)$, we need to define it in every world ω that satisfies β . But how?

- We already know that the sum of all such beliefs must be 1: $\sum_{\omega \models \beta} \Pr(\omega|\beta) = 1$. (C1)

- But this leaves us with many options for $\Pr(\omega|\beta)$ when world ω satisfies β .

What should we do?

The Idea: Minimal Change

- Impossible worlds (0-probability worlds) stay 0:

$$\Pr(\omega|\beta) = 0 \quad \text{for all } \omega \text{ where } \Pr(\omega) = 0. \quad (\text{C2})$$

- Relative probability of positive probability worlds, also stay the same:

$$\frac{\Pr(\omega)}{\Pr(\omega')} = \frac{\Pr(\omega|\beta)}{\Pr(\omega'|\beta)} \quad \text{for all } \omega, \omega' \models \beta, \Pr(\omega) > 0, \Pr(\omega') > 0. \quad (\text{C3})$$

Having those constraints...

- We already know that the sum of all such beliefs must be 1: $\sum_{\omega \models \beta} \Pr(\omega|\beta) = 1.$ (C1)

- Impossible worlds (0- probability worlds) stay 0:

$$\Pr(\omega|\beta) = 0 \quad \text{for all } \omega \text{ where } \Pr(\omega) = 0. \quad (\text{C2})$$

- Relative probability of positive probability worlds, also stay the same:

$$\frac{\Pr(\omega)}{\Pr(\omega')} = \frac{\Pr(\omega|\beta)}{\Pr(\omega'|\beta)} \quad \text{for all } \omega, \omega' \models \beta, \Pr(\omega) > 0, \Pr(\omega') > 0. \quad (\text{C3})$$

...leaves us with only one option:

$$\Pr(\omega|\beta) = \frac{\Pr(\omega)}{\Pr(\beta)} \quad \text{for all } \omega \models \beta. \quad (\text{B2})$$



which is a result of normalising the old belief w.r.t. β

Now, our new state of belief is completely defined:

$$\Pr(\omega|\beta) \stackrel{\text{def}}{=} \begin{cases} 0, & \text{if } \omega \models \neg\beta \\ \frac{\Pr(\omega)}{\Pr(\beta)}, & \text{if } \omega \models \beta. \end{cases}$$

The new state of belief $\Pr(.|\beta)$ is obtained via *conditioning* the old state \Pr on evidence β .

Consider the earlier example with the evidence $\text{Alarm} = \text{true}$.

<i>world</i>	<i>Earthquake</i>	<i>Burglary</i>	<i>Alarm</i>	$\Pr(.)$	$\Pr(. \text{Alarm})$
ω_1	true	true	true	.0190	.0190/.2442
ω_2	true	true	false	.0010	0
ω_3	true	false	true	.0560	.0560/.2442
ω_4	true	false	false	.0240	0
ω_5	false	true	true	.1620	.1620/.2442
ω_6	false	true	false	.0180	0
ω_7	false	false	true	.0072	.0072/.2442
ω_8	false	false	false	.7128	0

$$\Pr(\text{Burglary}) = .2 \quad \Pr(\text{Earthquake}) = .1$$

$$\Pr(\text{Burglary}|\text{Alarm}) \approx .741 \uparrow \quad \Pr(\text{Earthquake}|\text{Alarm}) \approx .307 \uparrow$$



$$\begin{aligned}
\Pr(\alpha|\beta) &= \sum_{\omega \models \alpha} \Pr(\omega|\beta) \quad \text{by (B1)} \\
&= \sum_{\omega \models \alpha, \omega \models \beta} \Pr(\omega|\beta) + \sum_{\omega \models \alpha, \omega \models \neg\beta} \Pr(\omega|\beta) \quad \text{since } \omega \text{ satisfies } \beta \text{ or } \neg\beta \text{ but not both} \\
&= \sum_{\omega \models \alpha, \omega \models \beta} \Pr(\omega|\beta) \quad \text{by (B2)} \\
&= \sum_{\omega \models \alpha \wedge \beta} \Pr(\omega|\beta) \quad \text{by properties of } \models \\
&= \sum_{\omega \models \alpha \wedge \beta} \Pr(\omega)/\Pr(\beta) \quad \text{by (B2)} \\
&= \frac{1}{\Pr(\beta)} \sum_{\omega \models \alpha \wedge \beta} \Pr(\omega) \\
&= \frac{\Pr(\alpha \wedge \beta)}{\Pr(\beta)} \quad \text{by (B1)}
\end{aligned}$$

$$\Pr(\alpha) \stackrel{\text{def}}{=} \sum_{\omega \models \alpha} \Pr(\omega) \quad (\text{B1})$$

$$\Pr(\omega|\beta) \stackrel{\text{def}}{=} \begin{cases} 0, & \text{if } \omega \models \neg\beta \\ \frac{\Pr(\omega)}{\Pr(\beta)}, & \text{if } \omega \models \beta. \end{cases} \quad (\text{B2})$$

Hence,

$$\Pr(\alpha|\beta) = \frac{\Pr(\alpha \wedge \beta)}{\Pr(\beta)} \quad (\text{Bayes Conditioning}) \quad (\text{B3})$$

Remark: Note that defined only when $\Pr(\beta) \neq 0$.

SUMMARY

- Worlds that contradict the evidence β will have zero probability.
- Worlds that have zero probability continue to have zero probability.
- Worlds that are consistent with evidence β and have positive probability will maintain their relative beliefs.

SCENARIO I

Evidence: Earthquake = true

$$\begin{aligned} \Pr(\text{Burglary}) &= .2 \\ \Pr(\text{Burglary}|\text{Earthquake}) &= .2 \\ \Pr(\text{Alarm}) &= .2442 \\ \Pr(\text{Alarm}|\text{Earthquake}) &\approx .75 \uparrow \end{aligned}$$

Observation: Belief in Burglary does not change while Alarm increases.

SCENARIO II

Evidence: Burglary = true

$$\begin{aligned} \Pr(\text{Alarm}) &= .2442 \\ \Pr(\text{Alarm}|\text{Burglary}) &\approx .905 \uparrow \\ \Pr(\text{Earthquake}) &= .1 \\ \Pr(\text{Earthquake}|\text{Burglary}) &= .1 \end{aligned}$$

Observation: Belief in Alarm increases while Earthquake stays the same.

SCENARIO I

We know: $\Pr(\text{Burglary}) \uparrow$ if $\Pr(\text{Alarm}) \uparrow$ and $\text{Alarm} = \text{true}$.

Question: How does our belief (on Burglary) change with the new evidence of an [Earthquake](#)?

Answer:

$$\begin{aligned} \Pr(\text{Burglary}|\text{Alarm}) &\approx .741 \\ \Pr(\text{Burglary}|\text{Alarm} \wedge \text{Earthquake}) &\approx .253 \downarrow \end{aligned}$$

Observation: Our belief in Burglary decreases, when we have an explanation for the Alarm.

SCENARIO II

We know: $\Pr(\text{Burglary}) \uparrow$ if $\Pr(\text{Alarm}) \uparrow$ and $\text{Alarm} = \text{true}$.

Question: How does our belief (on Burglary) change with the confirmation of [no-Earthquake](#)?

Answer:

$$\begin{aligned} \Pr(\text{Burglary}|\text{Alarm}) &\approx .741 \\ \Pr(\text{Burglary}|\text{Alarm} \wedge \neg\text{Earthquake}) &\approx .957 \uparrow \end{aligned}$$

Observation: The new evidence of [no-Earthquake](#) strengthens our belief on Burglary (as an explanation for the Alarm).

Recall the original state of belief:

world	Earthquake	Burglary	Alarm	$\Pr(\cdot)$
ω_1	true	true	true	.0190
ω_2	true	true	false	.0010
ω_3	true	false	true	.0560
ω_4	true	false	false	.0240
ω_5	false	true	true	.1620
ω_6	false	true	false	.0180
ω_7	false	false	true	.0072
ω_8	false	false	false	.7128

Observation:

$$\begin{aligned} \Pr(\text{Earthquake}) &= .1 \\ \Pr(\text{Earthquake}|\text{Burglary}) &= .1 \end{aligned}$$

Conclusion:

According to Pr, evidence of [Burglary](#) does not change the belief in [Earthquake](#), hence two events are *independent*.

Independence

More generally, given a distribution \Pr , the events α and β are **independent** in \Pr iff

$$\Pr(\alpha|\beta) = \Pr(\alpha) \quad \text{or} \quad \Pr(\beta) = 0.$$

Observe that Burglary is also independent of Earthquake, according to \Pr :

$$\begin{aligned}\Pr(\text{Burglary}) &= .2 \\ \Pr(\text{Burglary}|\text{Earthquake}) &= .2\end{aligned}$$

Joint probability of independent events

Given a \Pr , the event α is independent from the event β w.r.t \Pr iff $\Pr(\alpha \wedge \beta) = \Pr(\alpha)\Pr(\beta)$.

Remark: Do not mix it up with logical disjointness.



Independence is a **dynamic notion**: two events which are independent can be dependent after a new **evidence**.

Example:

$\Pr(\text{Burglary} \text{Alarm})$	$\approx .741$
$\Pr(\text{Burglary} \text{Alarm} \wedge \text{Earthquake})$	$\approx .253$

Observation: Note that initially independent Burglary and Earthquake became dependent after the evidence Alarm. Why?

Explanation: These are two competing explanations for Alarm confirming one decreases the belief in the other.

<i>world</i>	Temp	Sensor1	Sensor2	$\Pr(\cdot)$
ω_1	normal	normal	normal	.576
ω_2	normal	normal	extreme	.144
ω_3	normal	extreme	normal	.064
ω_4	normal	extreme	extreme	.016
ω_5	extreme	normal	normal	.008
ω_6	extreme	normal	extreme	.032
ω_7	extreme	extreme	normal	.032
ω_8	extreme	extreme	extreme	.128

We have two different sensors which have noise and different reliabilities.

It follows,

$$\Pr(\text{Temp} = \text{normal}) = .80$$

$$\Pr(\text{Sensor1} = \text{normal}) = .76$$

$$\Pr(\text{Sensor2} = \text{normal}) = .68$$

Initial Setting:

$$\Pr(\text{Temp} = \text{normal}) = .80$$

$$\Pr(\text{Sensor1} = \text{normal}) = .76$$

$$\Pr(\text{Sensor2} = \text{normal}) = .68$$

<i>world</i>	Temp	Sensor1	Sensor2	$\Pr(\cdot)$
ω_1	normal	normal	normal	.576
ω_2	normal	normal	extreme	.144
ω_3	normal	extreme	normal	.064
ω_4	normal	extreme	extreme	.016
ω_5	extreme	normal	normal	.008
ω_6	extreme	normal	extreme	.032
ω_7	extreme	extreme	normal	.032
ω_8	extreme	extreme	extreme	.128

Follow up:

Suppose, we checked the first sensor and it reads **normal**. Our belief in the second sensor reading **normal** as well, would increase (as expected):

$$\Pr(\text{Sensor2} = \text{normal} | \text{Sensor1} = \text{normal}) \approx .768 \uparrow$$

Hence, they are dependent!

Initial Setting:

$$\Pr(\text{Temp} = \text{normal}) = .80$$

$$\Pr(\text{Sensor1} = \text{normal}) = .76$$

$$\Pr(\text{Sensor2} = \text{normal}) = .68$$

$$\Pr(\text{Sensor2} = \text{normal} | \text{Sensor1} = \text{normal}) \approx .768 \uparrow$$

<i>world</i>	Temp	Sensor1	Sensor2	$\Pr(\cdot)$
ω_1	normal	normal	normal	.576
ω_2	normal	normal	extreme	.144
ω_3	normal	extreme	normal	.064
ω_4	normal	extreme	extreme	.016
ω_5	extreme	normal	normal	.008
ω_6	extreme	normal	extreme	.032
ω_7	extreme	extreme	normal	.032
ω_8	extreme	extreme	extreme	.128

They are dependent.

Now, assume you observed a new evidence: **Temp = normal**

$$\Pr(\text{Sensor2} = \text{normal} | \text{Temp} = \text{normal}) = .80$$

$$\Pr(\text{Sensor2} = \text{normal} | \text{Temp} = \text{normal}, \text{Sensor1} = \text{normal}) = .80$$



Now, they have become independent!

This urges us to look into the more general notion of **conditional independence**.

Conditional Independence

- Given a state of belief \Pr , event α is *conditionally independent* from the event β given γ iff

$$\Pr(\alpha|\beta \wedge \gamma) = \Pr(\alpha|\gamma) \quad \text{or} \quad \Pr(\beta \wedge \gamma) = 0. \quad (\text{B4})$$

- Realise that this is **symmetric** (i.e., β is *conditionally independent* from the event α given γ), hence, alternatively:

$$\Pr(\alpha \wedge \beta|\gamma) = \Pr(\alpha|\gamma)\Pr(\beta|\gamma) \quad \text{or} \quad \Pr(\gamma) = 0. \quad (\text{B4}')$$

Bayes' Rule

Recall: $\Pr(\alpha|\beta) = \frac{\Pr(\alpha \wedge \beta)}{\Pr(\beta)} \quad (\text{Bayes Conditioning}) \quad (\text{B3})$

By repeatedly applying *Bayes conditioning* (B3), we get the following:

Chain rule

$$\Pr(\alpha_1 \wedge \alpha_2 \wedge \dots \wedge \alpha_n) = \Pr(\alpha_1|\alpha_2 \wedge \dots \wedge \alpha_n)\Pr(\alpha_2|\alpha_3 \wedge \dots \wedge \alpha_n) \dots \Pr(\alpha_n). \quad (\text{B5})$$

...which will be very important for *Bayesian Networks* later on.

A set of events β_1, \dots, β_n is called:

- logically disjoint* (or *mutually exclusive*) iff $\text{Mods}(\beta_j) \cap \text{Mods}(\beta_k) = \emptyset$ for $j \neq k$ where $j, k \in \{1, \dots, n\}$.
- exhaustive* iff $\bigcup_{i=1}^n \text{Mods}(\beta_i) = \Omega$, where Ω is the set of all worlds.

Marginalisation (also called *case analysis* or *law of total probability*)

Assume a set of events β_1, \dots, β_n which are mutually exclusive and exhaustive, then

$$\Pr(\alpha) = \sum_{i=1}^n \Pr(\alpha \wedge \beta_i) \quad (\text{B6})$$

$$\Pr(\alpha) = \sum_{i=1}^n \Pr(\alpha|\beta_i)\Pr(\beta_i) \quad (\text{B6}')$$

Why does it work?

Examples: Consider the following two simple examples of case analysis:

- $\Pr(\alpha) = \Pr(\alpha \wedge \beta) + \Pr(\alpha \wedge \neg\beta)$
- $\Pr(\alpha) = \Pr(\alpha|\beta)\Pr(\beta) + \Pr(\alpha|\neg\beta)\Pr(\neg\beta).$

Observe:

These expressions hold because β and $\neg\beta$ are mutually exclusive and exhaustive.

Value:

This is useful when, as in many situations, computing the belief in cases is easier than the whole α .

Bayes' rule or Bayes' theorem

$$\Pr(\alpha|\beta) = \frac{\Pr(\beta|\alpha)\Pr(\alpha)}{\Pr(\beta)}$$



Classical Use:

- When an event α is perceived to be a cause of an event β e.g., α is a disease, β is a symptom.
- The belief in an effect given its cause, $\Pr(\beta|\alpha)$, is usually more readily available than the belief in a cause given one of its effects, $\Pr(\alpha|\beta)$.
- Bayes' theorem helps us to compute the latter from the former.

Exercise: False Positive

Problem Setting: Suppose that we have a patient who was just tested for a particular disease and the test came out positive. We know that one in every thousand people has this disease. We also know that the test is not completely reliable: it has a false positive rate of 2% and a false negative rate of 5%.

Question: What should be our belief in the patient having the disease given that the test came out positive?

Hint: Let the propositional variable D stand for “the patient has the disease” and the propositional variable T stand for “the test came out positive,” our goal is then to compute $\Pr(D|T)$.

Problem Statement

Suppose that we have a patient who was just tested for a particular disease and the test came out positive. We know that one in every thousand people has this disease. We also know that the test is not reliable: it has a false positive rate of 2% and a false negative rate of 5%. **What should be our belief in the patient having the disease given that the test came out positive?** In particular, let the propositional variable D stand for “the patient has the disease” and the propositional variable T stand for “the test came out positive,” our goal is then to compute $\Pr(D|T)$.

Our prior belief in the patient having the disease before we run any tests: $\Pr(D) = \frac{1}{1,000}$

The false positive rate of the test is $\Pr(T|\neg D) = \frac{2}{100}$ hence, $\Pr(\neg T|\neg D) = \frac{98}{100}$

Similarly, the false negative rate is $\Pr(\neg T|D) = \frac{5}{100}$ hence, $\Pr(T|D) = \frac{95}{100}$

Our prior belief in the patient having the disease before we run any tests: $\Pr(D) = \frac{1}{1,000}$

The false positive rate of the test is $\Pr(T|\neg D) = \frac{2}{100}$ hence, $\Pr(\neg T|\neg D) = \frac{98}{100}$

Similarly, the false negative rate is $\Pr(\neg T|D) = \frac{5}{100}$ hence, $\Pr(T|D) = \frac{95}{100}$

Using Bayes rule, we get $\Pr(D|T) = \frac{\frac{95}{100} \times \frac{1}{1,000}}{\Pr(T)}$

$\Pr(T)$ is not readily available, but can be obtained by case analysis:

$$\begin{aligned}\Pr(T) &= \Pr(T|D)\Pr(D) + \Pr(T|\neg D)\Pr(\neg D) \\ &= \frac{95}{100} \times \frac{1}{1,000} + \frac{2}{100} \times \frac{999}{1,000} = \frac{2,093}{100,000}\end{aligned}$$

which yields,



$$\Pr(D|T) = \frac{95}{2,093} \approx 4.5\%$$

2nd Solution: Constructing the Belief State

Because we have only two events of interest, T and D, leading to only four worlds, this solution is feasible.

world	D	T	
ω_1	true	true	has disease, test positive
ω_2	true	false	has disease, test negative
ω_3	false	true	has no disease, test positive
ω_4	false	false	has no disease, test negative

which gives rise to ...

$$\Pr(\omega_1) = \Pr(T \wedge D) = \Pr(T|D)\Pr(D)$$

$$\Pr(\omega_2) = \Pr(\neg T \wedge D) = \Pr(\neg T|D)\Pr(D)$$

$$\Pr(\omega_3) = \Pr(T \wedge \neg D) = \Pr(T|\neg D)\Pr(\neg D)$$

$$\Pr(\omega_4) = \Pr(\neg T \wedge \neg D) = \Pr(\neg T|\neg D)\Pr(\neg D).$$

$$\begin{aligned}
 \Pr(\omega_1) &= \Pr(T \wedge D) = \Pr(T|D)\Pr(D) \\
 \Pr(\omega_2) &= \Pr(\neg T \wedge D) = \Pr(\neg T|D)\Pr(D) \\
 \Pr(\omega_3) &= \Pr(T \wedge \neg D) = \Pr(T|\neg D)\Pr(\neg D) \\
 \Pr(\omega_4) &= \Pr(\neg T \wedge \neg D) = \Pr(\neg T|\neg D)\Pr(\neg D).
 \end{aligned}$$

..whose values are readily available in the problem setting, which yields:

world	D	T	Pr(.)		
ω_1	true	true	95/100	\times	$1/1,000 = .00095$
ω_2	true	false	5/100	\times	$1/1,000 = .00005$
ω_3	false	true	2/100	\times	$999/1,000 = .01998$
ω_4	false	false	98/100	\times	$999/1,000 = .97902$

$$\frac{\Pr(\omega_1)}{\Pr(\omega_1) + \Pr(\omega_3)} \approx 4.5\%$$

VL

Probabilistic Graphical Models 2

Introduction

We have seen that joint probability distribution is **exponential** in the size of **variables**.

i.e., $\Pr(X_1, \dots, X_n)$ needs a state of belief/distribution, specified for 2^n cases/possible worlds.

n = 3 variables

world	Earthquake	Burglary	Alarm	Pr(.)
ω_1	true	true	true	.0190
ω_2	true	true	false	.0010
ω_3	true	false	true	.0560
ω_4	true	false	false	.0240
ω_5	false	true	true	.1620
ω_6	false	true	false	.0180
ω_7	false	false	true	.0072
ω_8	false	false	false	.7128

2ⁿ = 8 cases/possible worlds

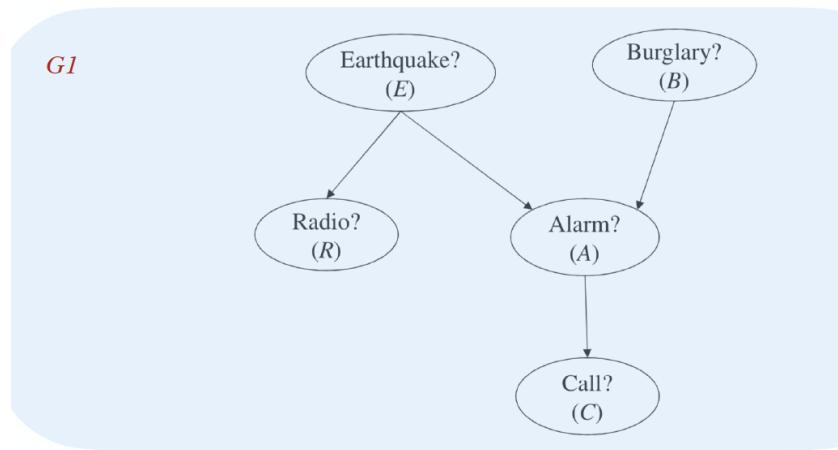
Even if this is addressed, one still needs to ensure that the synthesized distribution is *faithful* i.e., matches the beliefs held about a given situation.

Representing an expert knowledge in some domain, one needs to ensure some correspondence between the independencies held by the distribution and those believed by the expert.

Ex: $\Pr(\text{Earthquake} \mid \text{Burglary}) = \Pr(\text{Earthquake})$ or $\Pr(\text{Aids} \mid \text{Covid-19}) = \Pr(\text{Aids})$.

The *Bayesian network (BN)* is a modelling tool to specify probability distribution which can, in principle, address these challenges. A BN relies on the basic insight that independence forms a significant aspect of beliefs and that it can be elicited relatively easily using the language of graphs.

A Graphical Tool for capturing Independence

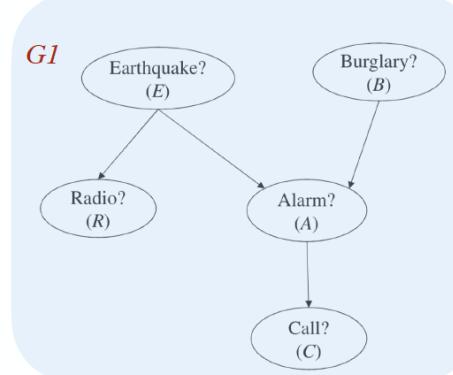


Assume a DAG (*directed acyclic graph*) in which:

- nodes represent propositional variables.
- edges represent *direct causal influences*.

Assuming that the links are causal, what can be said on the dependencies in *G1*:

- Alarm triggering (*A*) is a direct cause of a call from the neighbour (*C*).
- If we get a radio report (*R*) (that an earthquake (*E*) took place in our neighbourhood), then our belief in *A* could increase, which in turn increase our belief in *C*.
- Yet, if we know already that $\neg A$, then our belief in *C* stays unchanged.
i.e., *C* is independent of *R* given $\neg A$ |



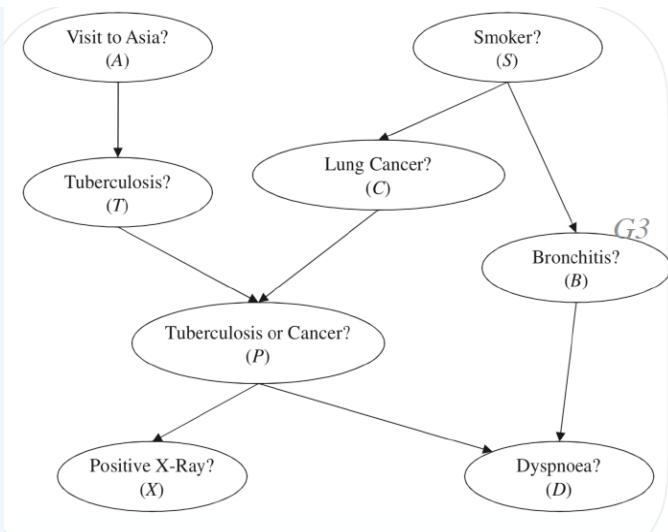
Example:

In a limited medical domain,

we could find a visit to Asia relevant to our belief in the x-ray test coming out positive.

we would find the visit irrelevant if we know for sure that the patient does not have tuberculosis

That is, X is dependent on A but is independent of A given $\neg T$.



First introducing some notions:

Given a DAG G, and a variable V in G,

- Parents(V) are the parent nodes of V i.e., a set of variables N which has a directed edge to V.

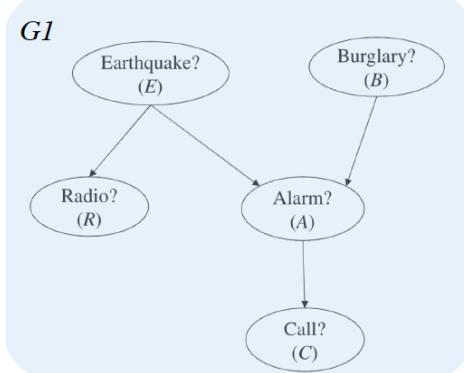
Ex.: Parents(A) = {E, B}

- Descendants(V) are the descendant nodes of V i.e., a set of variables N s.t. there is a directed path from V.

Ex.: Descendants(B) = {A, C}

- Non_Descendants(V) are the nodes that are neither in Descendants(V) nor Parents(V) other than V.

Ex.: Non_Descendants(B) = {E, R}



Then, each DAG G corresponds to a compact representation of the following independence statements:

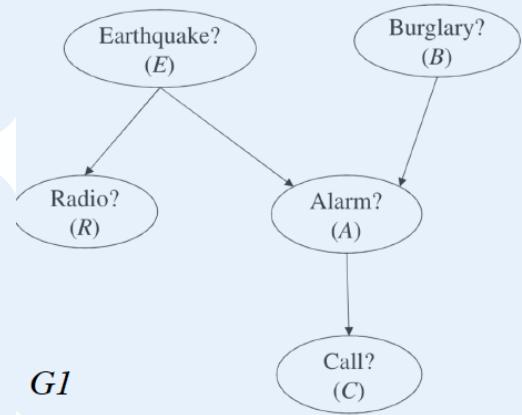
(also called) **Markovian assumptions** or **Markov(G)**

$I(V, \text{Parents}(V), \text{Non_Descendants}(V))$ for all variables V in DAG G .

which says every variable is conditionally independent of its nondescendants given its parents.

Example

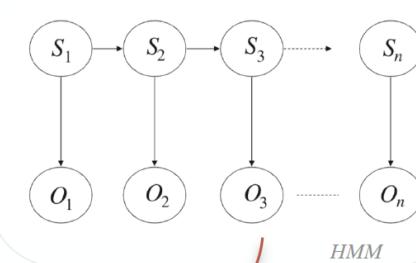
Markov(G_1):
 $I(C, A, \{B, E, R\})$
 $I(R, E, \{A, B, C\})$
 $I(A, \{B, E\}, R)$
 $I(B, \emptyset, \{E, R\})$
 $I(E, \emptyset, B)$



Markov(G):

$I(V, \text{Parents}(V), \text{Non_Descendants}(V))$ for all variables V in DAG G .

which says every variable is conditionally independent of its nondescendants given its parents.



Example

This is the famous **Hidden Markov Model (HMM)** *

variables S_1, S_2, \dots, S_n represent the state of a dynamic system and

variables O_1, O_2, \dots, O_n represent sensors that measure the system state at time points $1, 2, \dots, n$, respectively.

For every S_i , $I(S_i, \{S_{i-1}\}, \{S_1, \dots, S_{i-2}, O_1, \dots, O_{i-1}\}) \in \text{Markov}(\text{HMM})$.

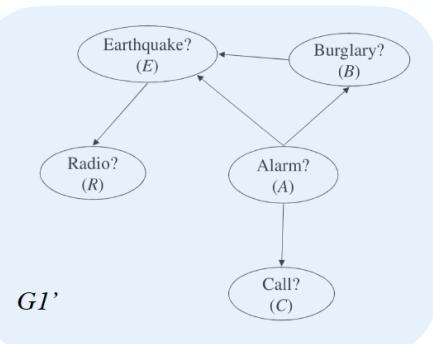
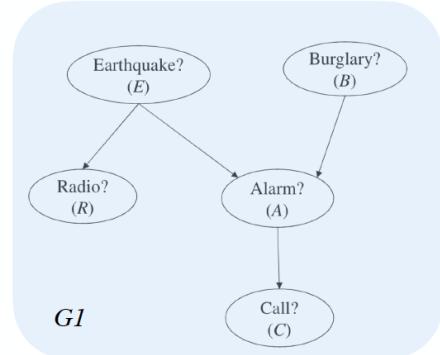


*: Vast amount of applications in *reinforcement learning*, *temporal pattern recognition* and etc.

Remark 1:

If one views the DAG as a causal structure, then

- Parents(V) are the *direct causes* of V ,
- Descendants(V) are the *effects* of V .



Remark 2:

Although we used causality to motivate the interpretation, the formal interpretation makes no reference to the notion of *causality*.

$$\text{Markov}(G1) \implies \text{Markov}(G1')$$

Bayesian Networks

Parameterising the Independence Structure

Suppose, we would like to construct a probability distribution \Pr that captures state of our belief in a given domain.

1st step: Construct a DAG

To construct a DAG G whose independence structure is consistent with our belief state about that domain.

This, indeed, constrains the possible choices, yet does not uniquely defines it.

2nd step: Add conditional probability tables

For every variable X in the DAG, and its parents U , For every value x of variable X , and every instantiation \mathbf{u} of U , define $\Pr(x \mid \mathbf{u})$.

The second step corresponds to constructing conditional probability table (CPT). This ensures that \Pr is unique.

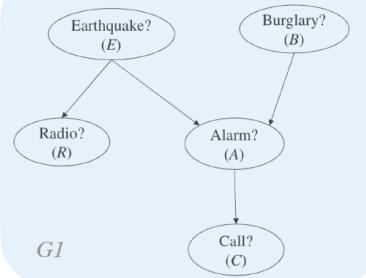
Conditional Probability Tables

Example:

For $G1$, we need to provide the following conditional probabilities:

$$\Pr(c|a), \Pr(r|e), \Pr(alb,e), \Pr(e), \Pr(b)$$

where a, b, c, e , and r are values of variables A, B, C, E, R , respectively.



Example:

A	C	$\Pr(c a)$
true	true	.80
true	false	.20
false	true	.001
false	false	.999

Such a table is called **CPT** (for C).

Observe: $\Pr(c|a) + \Pr(\bar{c}|a) = 1$ and $\Pr(c|\bar{a}) + \Pr(\bar{c}|\bar{a}) = 1$.

Two entries in the CPT are redundant i.e., can be inferred from one another, hence can be omitted.



Bayesian Networks

Finally, we can formally define a Bayesian network [Pearl, 1985]:

Given a set Z of variables, a **Bayesian network** (BN) for Z is a pair (G, Θ) where:

- G is a DAG over Z , called the network *structure*.
- Θ is a set of CPTs, one for each element in Z , called the network *parameterisation*.

Notation

- $\Theta_{X|U}$ denotes the CPT for variable X and its parents U .
- $\theta_{x|u}$, the network *parameter*, denotes the value assigned by $\Theta_{X|U}$ to the $\Pr(x|u)$.

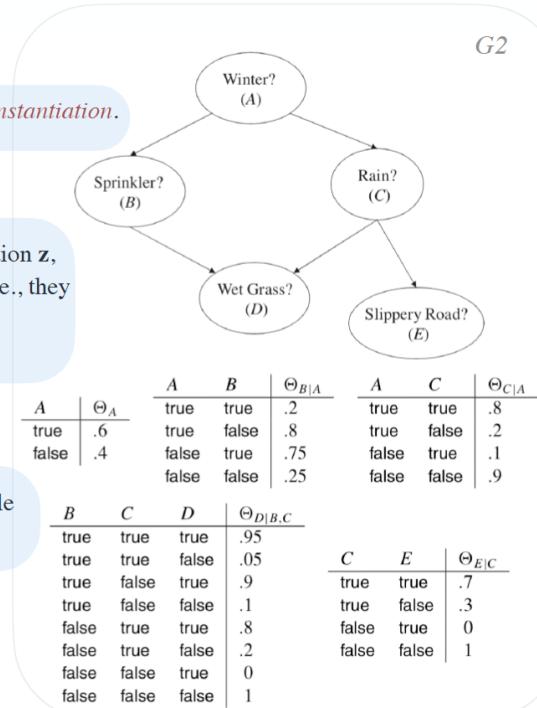
Remark: $\sum_x \theta_{x|u} = 1$ for every parent instantiation u .

Auxilliary notions:

An instantiation of all network variables will be called a *network instantiation*.

A network parameter $\theta_{x|u}$ is *compatible* with a network instantiation z , denoted $\theta_{x|u} \sim z$, iff the instantiations $x|u$ and z are compatible (i.e., they agree on the values they assign to their common variables).

$\theta_a, \theta_{b|a}, \theta_{\bar{c}|a}, \theta_{d|\bar{b}, \bar{c}}$ and $\theta_{\bar{e}|\bar{c}}$ are the network parameters compatible with network instantiation $a, b, \bar{c}, d, \bar{e}$ in $G2$.



Chain Rule for Bayesian networks:

$$\Pr(\mathbf{z}) \stackrel{\text{def}}{=} \prod_{\theta_{x|u} \sim \mathbf{z}} \theta_{x|u}.$$

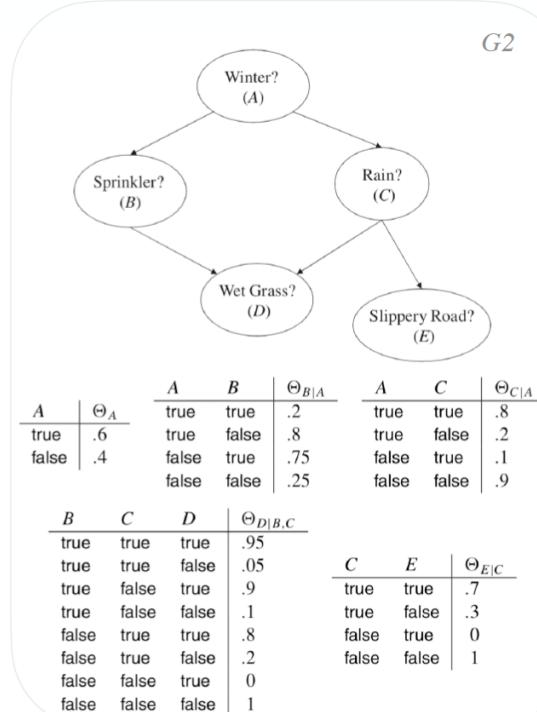
Intuitively:

The probability assigned to a network instantiation \mathbf{z} is simply the product of all network parameters compatible with \mathbf{z} .

Examples:

$$\begin{aligned} \Pr(a, b, \bar{c}, d, \bar{e}) &= \theta_a \theta_{b|a} \theta_{\bar{c}|a} \theta_{d|\bar{b}, \bar{c}} \theta_{\bar{e}|\bar{c}} \\ &= (.6)(.2)(.2)(.9)(1) \\ &= .0216 \end{aligned}$$

$$\begin{aligned} \Pr(\bar{a}, \bar{b}, \bar{c}, \bar{d}, \bar{e}) &= \theta_{\bar{a}} \theta_{\bar{b}|\bar{a}} \theta_{\bar{c}|\bar{a}} \theta_{\bar{d}|\bar{b}, \bar{c}} \theta_{\bar{e}|\bar{c}} \\ &= (.4)(.25)(.9)(1)(1) \\ &= .09 \end{aligned}$$



Remark! Every Bayesian network implicitly represents a unique \Pr .

The distribution Pr specified by a Bayesian network (G, Θ) is guaranteed to satisfy every independence assumption in $\text{Markov}(G)$:

$$I_{\text{Pr}}(X, \text{Parents}(X), \text{Non Descendants}(X))$$

for every variable X in the network.

Recall:

$I_{\text{Pr}}(\mathbf{X}, \mathbf{Z}, \mathbf{Y})$ means distribution Pr finds variables \mathbf{X} independent of variables \mathbf{Y} given variables \mathbf{Z} :

$$\text{Pr}(\mathbf{x}|\mathbf{z}, \mathbf{y}) = \text{Pr}(\mathbf{x}|\mathbf{z}) \text{ or } \text{Pr}(\mathbf{y}, \mathbf{z}) = 0,$$

for all instantiations $\mathbf{x}, \mathbf{y}, \mathbf{z}$ of variables $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$, respectively.

Properties of Probabilistic Independence

In principle, a Bayesian network G might imply more independence statement than the ones in $\text{Markov}(G)$.

Observe: According to G_2 , D and E are independent given A and C ,

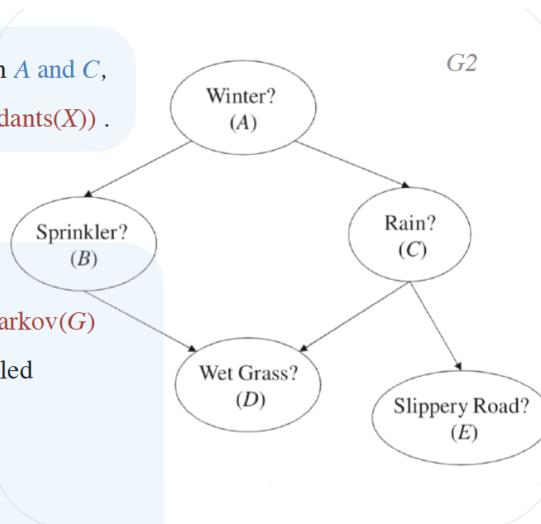
Yet they are not included in $I_{\text{Pr}}(X, \text{Parents}(X), \text{Non Descendants}(X))$.

This one and the additional ones follow from the ones in $\text{Markov}(G)$

Using a set of properties for probabilistic independence, called

Graphoid axioms:

- Symmetry
- Decomposition
- Weak Union
- Contraction



Properties of Probabilistic Independence: Symmetry

The first and the simplest property of probabilistic independence we consider is symmetry:

Symmetry:

$$I_{\Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y}) \text{ if and only if } I_{\Pr}(\mathbf{Y}, \mathbf{Z}, \mathbf{X}).$$

Intuition (Symmetry):

If learning \mathbf{y} does not influence our belief in \mathbf{x} , then learning \mathbf{x} does not influence our belief in \mathbf{y} , either.

Example:

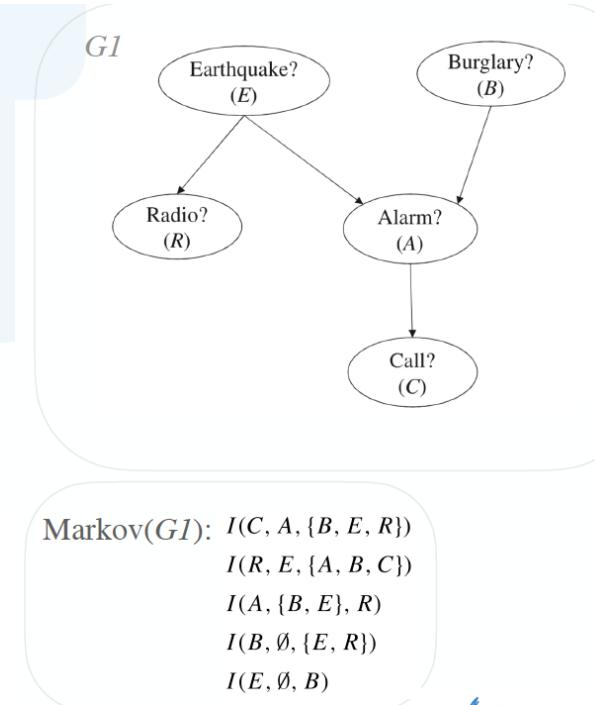
We know that $I_{\Pr}(A, \{B, E\}, R) \in \text{Markov}(G1)$.

Using [symmetry](#), we can conclude that

$I_{\Pr}(R, \{B, E\}, A)$ which is **not** in $\text{Markov}(G1)$.

Little reminder: “only if” means “ \implies ”

“if” means “ \iff ”



$\text{Markov}(G1):$

- $I(C, A, \{B, E, R\})$
- $I(R, E, \{A, B, C\})$
- $I(A, \{B, E\}, R)$
- $I(B, \emptyset, \{E, R\})$
- $I(E, \emptyset, B)$

Decomposition

Decomposition:

$$\implies$$

$I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y} \cup \mathbf{W})$ only if $I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y})$ and $I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{W})$.

Intuition (Decomposition):

If learning \mathbf{yw} does not influence our belief in \mathbf{x} , then learning \mathbf{y} alone or \mathbf{w} alone, will not influence our belief in \mathbf{x} .

Remark

Note that the opposite, say, *composition* is not true in general:

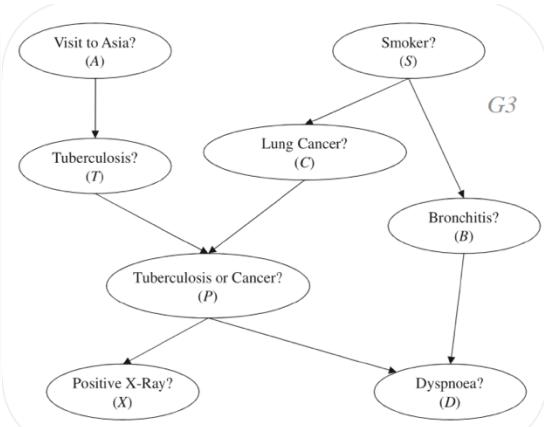
$I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y})$ and $I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{W})$ only if $I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y} \cup \mathbf{W})$

Example:

We know that $I_{Pr}(B, S, \{A, C, P, T, X\}) \in \text{Markov}(G3)$.

Using *decomposition*, we can conclude that

$I_{Pr}(B, S, C)$: Once we know that the person is a smoker, our belief in developing bronchitis is not influenced by information about developing cancer.



More generally, **decomposition**, let us to state

$$I_{\Pr}(X, \text{Parents}(X), \mathbf{W}) \text{ for every } \mathbf{W} \subseteq \text{Non Descendants}(X)$$

which can be seen as a strengthening of **Markov(G)**:

$$I_{\Pr}(X, \text{Parents}(X), \text{Non Descendants}(X))$$

which is a special case when \mathbf{W} contains all nondescendants of X .

Another application:

Decomposition can simplify the initial chain rule.

Example:

Assume, we want to calculate $\Pr(r, c, a, e, b)$.

Then, by $\Pr(\alpha_1 \wedge \alpha_2 \wedge \dots \wedge \alpha_n) = \Pr(\alpha_1|\alpha_2 \wedge \dots \wedge \alpha_n)\Pr(\alpha_2|\alpha_3 \wedge \dots \wedge \alpha_n)\dots\Pr(\alpha_n)$. (B5)

$\Pr(r, c, a, e, b) = \Pr(r|c, a, e, b)\Pr(c|a, e, b)\Pr(a|e, b)\Pr(e|b)\Pr(b)$.

By, $I_{\Pr}(X, \text{Parents}(X), \mathbf{W})$ for every $\mathbf{W} \subseteq \text{Non Descendants}(X)$, we get:

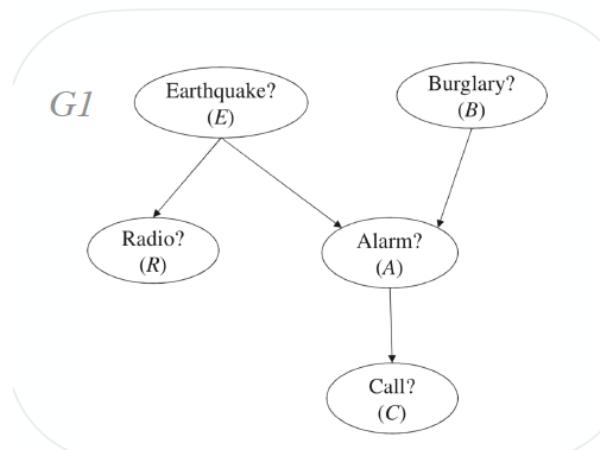
$$\Pr(r, c, a, e, b) = \Pr(r|e)\Pr(c|a)\Pr(a|e, b)\Pr(e|b)\Pr(b)$$

$$= \theta_{r|e} \theta_{c|a} \theta_{a|e,b} \theta_e \theta_b$$

$$\Pr(r|c, a, e, b) = \Pr(r|e)$$

$$\Pr(c|a, e, b) = \Pr(c|a)$$

$$\Pr(e|b) = \Pr(e).$$



Weak Union:

Weak Union

$I_{Pr}(X, Z, Y \cup W)$ only if $I_{Pr}(X, Z \cup Y, W)$.

Intuition (Weak Union):

If the information yw is not relevant to our belief in x , then the partial information y will not make the rest of the information w , relevant.

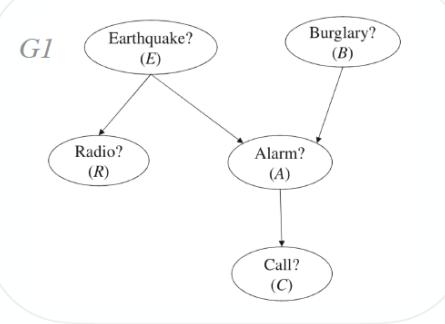
Example:

Assume a Θ on $G1$ and also assume that Pr is generated by $(G1, \Theta)$.

We have $I_{Pr}(C, A, \{B, E, R\}) \in \text{Markov}(G1)$

By [weak-union](#), we can conclude

$I_{Pr}(C, \{A, E, B\}, R)$ which was initially not in $\text{Markov}(G1)$



More generally:

$I_{Pr}(X, \text{Parents}(X) \cup W, \text{Non Descendants}(X) \setminus W)$, for any $W \subseteq \text{Non Descendants}(X)$.

Contraction

Contraction:

$I_{\text{Pr}}(\mathbf{X}, \mathbf{Z}, \mathbf{Y})$ and $I_{\text{Pr}}(\mathbf{X}, \mathbf{Z} \cup \mathbf{Y}, \mathbf{W})$ only if $I_{\text{Pr}}(\mathbf{X}, \mathbf{Z}, \mathbf{Y} \cup \mathbf{W})$.

Intuition (Contraction):

If after learning the irrelevant information \mathbf{y} the information \mathbf{w} is found to be irrelevant to our belief in \mathbf{x} , then the combined information \mathbf{yw} must have been irrelevant from the beginning.

Remark

Realize that contraction is a weaker version of composition:

$I_{\text{Pr}}(\mathbf{X}, \mathbf{Z}, \mathbf{Y})$ and $I_{\text{Pr}}(\mathbf{X}, \mathbf{Z}, \mathbf{W})$ only if $\cancel{I_{\text{Pr}}(\mathbf{X}, \mathbf{Z}, \mathbf{Y} \cup \mathbf{W})}$

d-Separation

We have seen that deriving new independencies from $\text{Markov}(G)$ is a non-trivial task.

Good news is that there is an easy graphical test called d-separation which captures the inferential power of graphoid axioms.

This means $\text{dsep}_G(\mathbf{X}, \mathbf{Z}, \mathbf{Y})$ (which reads \mathbf{X} and \mathbf{Y} is d-separated by \mathbf{Z}) implies $I_{\text{Pr}}(\mathbf{X}, \mathbf{Z}, \mathbf{Y})$ for every probability distribution Pr induced by G .

The intuition is similar to a water flow along a pipe (a path in the network) controlled by a set of valves:

Two variables are independent if all the paths between them are blocked by a closed valve (to be defined).

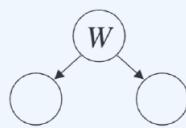
We will need only three **types of valves**:

- A *sequential valve* ($\rightarrow W \rightarrow$) arises when W is a parent of one of its neighbors and a child of the other.
- A *divergent valve* ($\leftarrow W \rightarrow$) arises when W is a parent of both neighbors.
- A *convergent valve* ($\rightarrow W \leftarrow$) arises when W is a child of both neighbors.

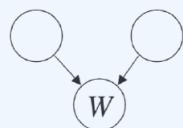
Types of Valves (visually)



Sequential



Divergent



Convergent

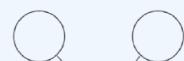
Types of Valves



Sequential

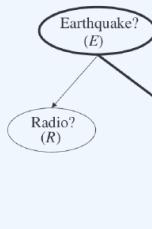


Divergent

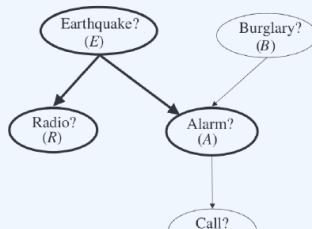


Convergent

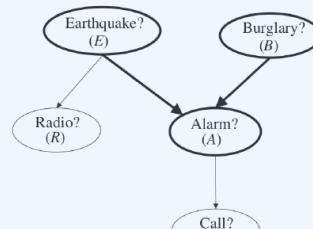
Examples:



Sequential valve



Divergent valve

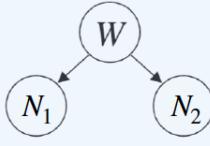


Convergent valve

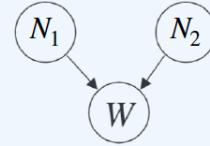
Types of Valves:



Sequential



Divergent



Convergent

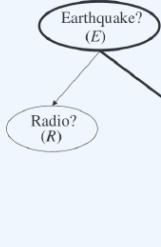
A Causal Interpretation to Help Intuition:

A *sequential valve* $N_1 \rightarrow W \rightarrow N_2$ declares variable W as an intermediary between a cause N_1 and its effect N_2 .

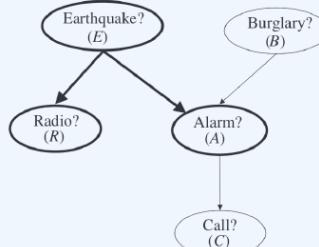
A *divergent valve* $N_1 \leftarrow W \rightarrow N_2$ declares variable W as a common cause of two effects N_1 and N_2 .

A *convergent valve* $N_1 \rightarrow W \leftarrow N_2$ declares variable W as a common effect of two causes N_1 and N_2 .

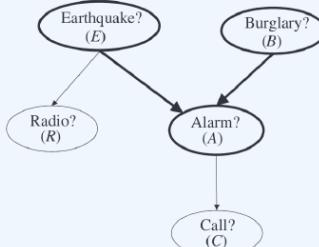
Examples:



Sequential valve



Divergent valve



Convergent valve

Closedness:

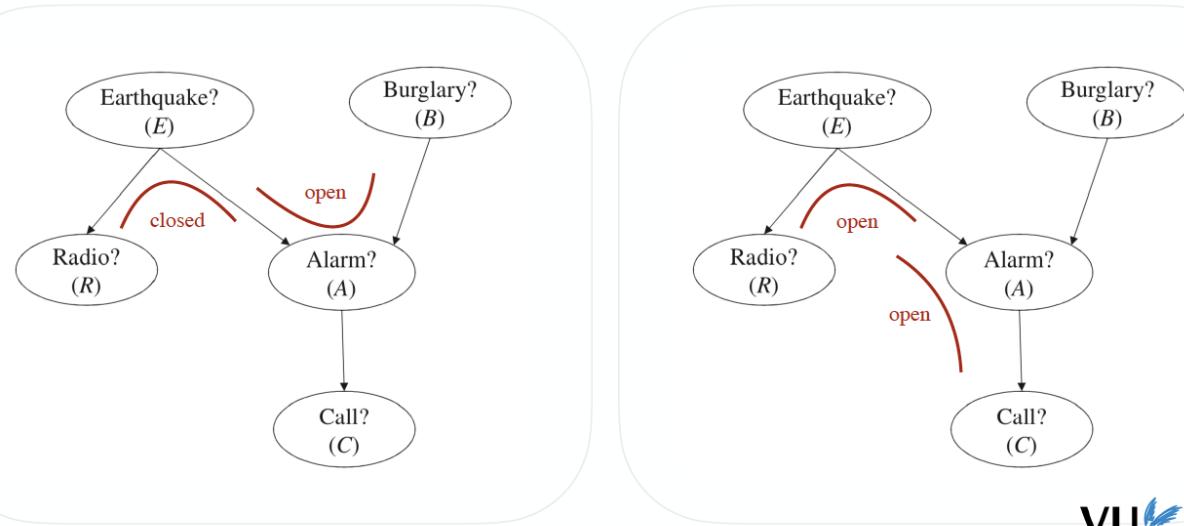
A *sequential valve* ($\rightarrow W \rightarrow$) is *closed* iff variable W appears in Z .

A *divergent valve* ($\leftarrow W \rightarrow$) is *closed* iff variable W appears in Z .

A *convergent valve* ($\rightarrow W \leftarrow$) is *closed* iff neither variable W nor any of its descendants appears in Z .

d-Separation:

Given disjoint sets \mathbf{X} , \mathbf{Y} , and \mathbf{Z} of nodes in a DAG G , \mathbf{X} and \mathbf{Y} are *d-separated* by \mathbf{Z} , denoted $dsep_G(\mathbf{X}, \mathbf{Z}, \mathbf{Y})$, iff every path between a node in \mathbf{X} and a node in \mathbf{Y} is blocked by \mathbf{Z} (i.e., has at least one closed valve $W \in \mathbf{Z}$).

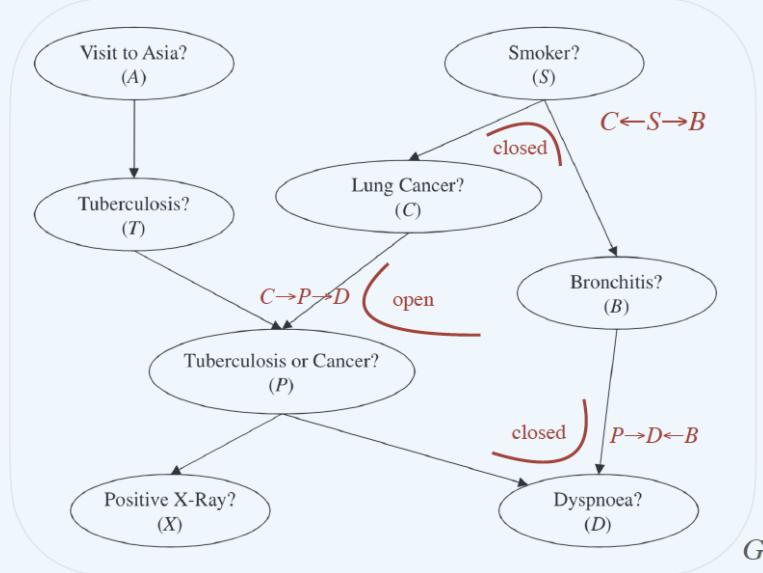


Are R and B d-separated by E, C ? Yes!

Are R and C d-separated? No!

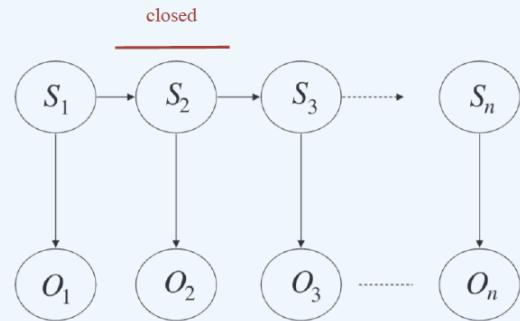


More example:



Does $dsep_G(B, S, C)$ hold? Yes!

More example:



$\text{dsep}_G(S_1, S_2, \{S_3, S_4\})?$ Yes!

Any path between S_1 and $\{S_3, S_4\}$ must have the valve $S_1 \rightarrow S_2 \rightarrow S_3$ on it, which is closed given S_2 .

Pruning Method

Paths between sets of nodes can be exponentially many. The following method guarantees it can be decided in linear time/space in the size of the graph.

Pruning method:

Given a DAG G and disjoint sets of nodes $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$ in G , a pruned G' w.r.t. \mathbf{Z} is obtained by:

Deleting every leaf node $W \notin \mathbf{X} \cup \mathbf{Y} \cup \mathbf{Z}$,

Deleting all edges outgoing from nodes in \mathbf{Z} ,

iteratively until both rules do not apply anymore.

Theorem:

\mathbf{X} and \mathbf{Y} are **d-separated** by \mathbf{Z} in G iff \mathbf{X} and \mathbf{Y} are disconnected in the pruned G' w.r.t. \mathbf{Z} .

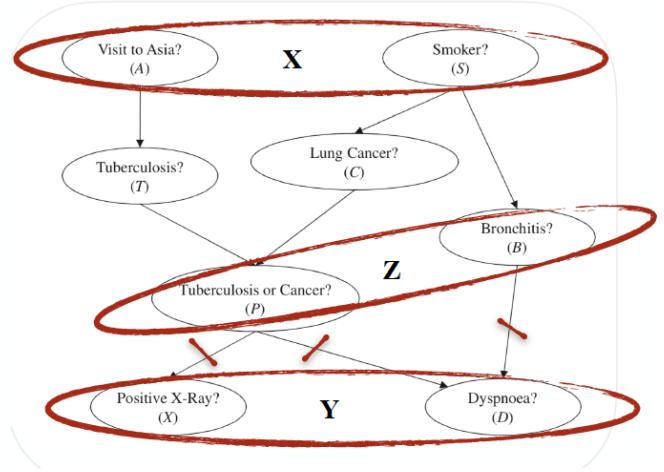
Recall the Pruning rules (iteratively):

Delete every leaf node $W \notin \mathbf{X} \cup \mathbf{Y} \cup \mathbf{Z}$.

Delete all edges outgoing from nodes in \mathbf{Z} .

Example:

Is $\mathbf{X} = \{A, S\}$ d-separated from $\mathbf{Y} = \{D, X\}$ by $\mathbf{Z} = \{B, P\}$?



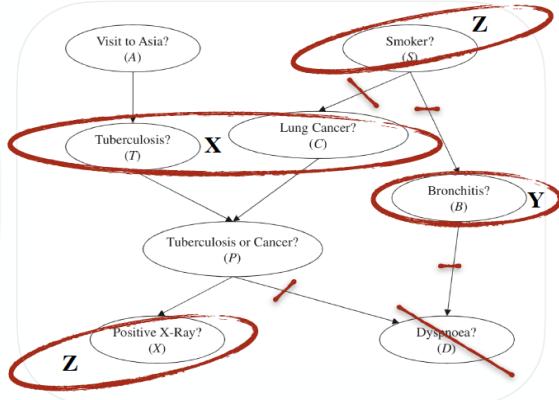
Recall the Pruning rules:

Delete every leaf node $W \notin \mathbf{X} \cup \mathbf{Y} \cup \mathbf{Z}$.

Delete all edges outgoing from nodes in \mathbf{Z} .

Example:

Is $\mathbf{X} = \{T, C\}$ d-separated from $\mathbf{Y} = \{B\}$ by $\mathbf{Z} = \{S, X\}$?



Yes!



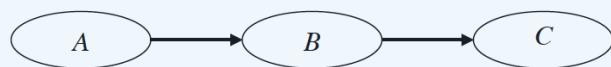
Theorem

Let Pr be a probability distribution induced by a BN (G, Θ) , then $\text{dsep}_G(\mathbf{X}, \mathbf{Z}, \mathbf{Y}) \implies I_{\text{Pr}}(\mathbf{X}, \mathbf{Z}, \mathbf{Y})$.

Theorem

$$\text{dsep}_G(\mathbf{X}, \mathbf{Z}, \mathbf{Y}) \not\implies I_{\text{Pr}}(\mathbf{X}, \mathbf{Z}, \mathbf{Y}).$$

Proof by contradiction:



Assume a CPT for variable B is with $\theta_{bla} = \theta_{b|\bar{a}}$. Then, although there is an unblocked edge between A and B , B is independent from A , hence C is independent from A .

However, C is not d-separated from A .

Probabilistic Graphical Models 3

Variable Elimination:
A tool for Inference in Bayesian Networks

The Idea of Variable Elimination

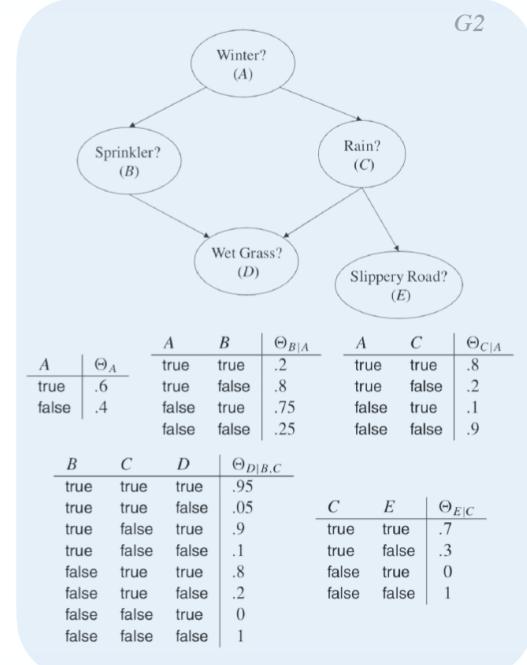
THE IDEA OF VARIABLE ELIMINATION

Consider the Bayesian network $G2$ from the previous lecture:

And suppose we are interested in computing the marginal $\Pr(D, E)$, that will look like:

D	E	$\Pr(D, E)$
true	true	.30443
true	false	.39507
false	true	.05957
false	false	.24093

How can we do that?



Suppose we are interested in computing the marginal $\Pr(D, E)$

How can we do that? One way to do that:

First, compute the joint probability distribution (recall the chain rule):

And then summing out each variable (i.e., A, B, C) one by one:

For instance, for A the following rows

A	B	C	D	E	$\Pr(.)$
true	true	true	true	true	.06384
false	true	true	true	true	.01995

Would merge into:

B	C	D	E	$\Pr(.)$
true	true	true	true	.08379 = .06384 + .01995

A	B	C	D	E	$\Pr(.)$
true	true	true	true	true	.06384
true	true	true	true	false	.02736
true	true	true	false	true	.00336
true	true	true	false	false	.00144
true	true	false	true	true	0
true	true	false	true	false	.02160
true	true	false	false	true	0
true	true	false	false	false	.00240
true	false	true	true	true	.21504
true	false	true	true	false	.09216
true	false	true	false	true	.05376
true	false	true	false	false	.02304
true	false	false	true	true	0
true	false	false	true	false	0
true	false	false	false	true	0
true	false	false	false	false	.09600
false	true	true	true	true	.01995
false	true	true	true	false	.00855
false	true	true	false	true	.00105
false	true	true	false	false	.00045
false	true	false	true	true	0
false	true	false	true	false	.24300
false	true	false	false	true	0
false	true	false	false	false	.02700
false	false	true	true	true	.00560
false	false	true	true	false	.00240
false	false	true	false	true	.00140
false	false	true	false	false	.00060
false	false	false	true	true	0
false	false	false	true	false	0
false	false	false	false	true	0
false	false	false	false	false	.0900

It will work: Resulting new distribution \Pr' we have is as good as \Pr i.e., $\Pr'(\alpha) = \Pr(\alpha)$ for any α that does not include A .

But it will necessarily require an exponential effort in the number of variables in the BN.

The **value** of **variable elimination** method is that it can *sometimes* avoid such complexity (i.e., having to construct the joint distribution explicitly by brute-force.)

An auxiliary notion: Factor

Given a set of variables X , a **factor** f is a (total) function* mapping each instantiation to a non-negative number i.e., $f: X \rightarrow \mathbb{R}_+$.

Hence, a factor does not necessarily represent a probability distribution.

Some notation:

Denote the variables that f is defined for, by $\text{vars}(f)$.

I.e., $\text{vars}(f) = X_1, \dots, X_n$ where $f(X_1, \dots, X_n) \in \mathbb{R}_+$.

We will also allow f to be defined over an empty set of variables.

Such factors called *trivial* as they assign a single number to the trivial instantiation T .

Example:

B	C	D	f_1
true	true	true	.95
true	true	false	.05
true	false	true	.9
true	false	false	.1
false	true	true	.8
false	true	false	.2
false	false	true	0
false	false	false	1

$$f_1(b, c, d) = \Pr(d|b, c)$$

D	E	f_2
true	true	.448
true	false	.192
false	true	.112
false	false	.248

$$f_2(d, e) = \Pr(d, e)$$

We will have two key operations on factors:
summing-out and multiplication.

We will have **two key operations** on factors:

summing-out and *multiplication*!

Summing-out (or Marginalising)

Given a factor f over variables \mathbf{X} , **summing-out** variable $X \in \mathbf{X}$, from factor f results in another factor $\sum_X f$ over variables $\mathbf{X} \setminus \{X\}$: and defined as

$$\left(\sum_X f \right)(y) \stackrel{\text{def}}{=} \sum_x f(x, y)$$

Example (cont'd):

The result of summing-out D from f_1 .

Example (cont'd):

If we sum-out all variables, we end up with the trivial factor:

$$\begin{array}{c|c} & \sum_B \sum_C \sum_D f_1 \\ \hline \top & 4 \end{array}$$

Example:

B	C	D	f_1
true	true	true	.95
true	true	false	.05
true	false	true	.9
true	false	false	.1
false	true	true	.8
false	true	false	.2
false	false	true	0
false	false	false	1

$$f_1(b, c, d) = \Pr(db, c)$$

B	C	$\sum_D f_1$
true	true	1
true	false	1
false	true	1
false	false	1

Remark: Summing-out is commutative.

$$\sum_Y \sum_X f = \sum_X \sum_Y f$$

Summing-out

Given a factor f over variables \mathbf{X} , **summing-out** variable $X \in \mathbf{X}$, from factor f results in another factor $\sum_X f$ over variables $\mathbf{X} \setminus \{X\}$: and defined as

$$\left(\sum_X f \right)(y) \stackrel{\text{def}}{=} \sum_x f(x, y)$$

Remark: Summing-out is commutative.

$$\sum_Y \sum_X f = \sum_X \sum_Y f$$

Complexity: $O(\exp(|\mathbf{X}|))$ of time & space.

Notation: Summing-out multiple variables:

$$\sum_{\mathbf{X}} f$$

where \mathbf{X} is a set of variables.

If \mathbf{Y} are other variables of f , then $\sum_{\mathbf{X}} f$

is also called **projection** of f on \mathbf{Y} .

Algorithm 1 SumOutVars($f(\mathbf{X})$, \mathbf{Z})

input:

$f(\mathbf{X})$: factor over variables \mathbf{X}

\mathbf{Z} : a subset of variables \mathbf{X}

output: a factor corresponding to $\sum_{\mathbf{Z}} f$

main:

- 1: $\mathbf{Y} \leftarrow \mathbf{X} - \mathbf{Z}$
 - 2: $f' \leftarrow$ a factor over variables \mathbf{Y} where $f'(y) = 0$ for all y
 - 3: **for** each instantiation \mathbf{y} **do**
 - 4: **for** each instantiation \mathbf{z} **do**
 - 5: $f'(\mathbf{y}) \leftarrow f'(\mathbf{y}) + f(\mathbf{yz})$
 - 6: **end for**
 - 7: **end for**
 - 8: **return** f'
-

Example:

Multiplication of Factors

Given factors f_1 and f_2 over variables \mathbf{X} and \mathbf{Y} , respectively. Multiplying $f_1 f_2$ is another factor over variables $\mathbf{Z} = \mathbf{X} \cup \mathbf{Y}$ s.t.

$$(f_1 f_2)(\mathbf{z}) \stackrel{\text{def}}{=} f_1(\mathbf{x}) f_2(\mathbf{y}),$$

where \mathbf{x} and \mathbf{y} are compatible with \mathbf{z} ; that is, $\mathbf{x} \sim \mathbf{z}$ and $\mathbf{y} \sim \mathbf{z}$.

Example (cont'd):

B	C	D	E	$f_1(B, C, D) f_2(D, E)$
true	true	true	true	.4256 = (.95)(.448)
true	true	true	false	.1824 = (.95)(.192)
true	true	false	true	.0056 = (.05)(.112)
:	:	:	:	:
false	false	false	false	.2480 = (1)(.248)

B	C	D	f_1
true	true	true	.95
true	true	false	.05
true	false	true	.9
true	false	false	.1
false	true	true	.8
false	true	false	.2
false	false	true	0
false	false	false	1

$$f_1(b, c, d) = \Pr(d|b, c)$$

D	E	f_2
true	true	.448
true	false	.192
false	true	.112
false	false	.248

$$f_2(d, e) = \Pr(d, e)$$

Remark: Factor multiplication is commutative and associative.



Multiplication of Factors

Given factors f_1 and f_2 over variables \mathbf{X} and \mathbf{Y} , respectively. Multiplying $f_1 f_2$ is another factor over variables $\mathbf{Z} = \mathbf{X} \cup \mathbf{Y}$ s.t.

$$(f_1 f_2)(\mathbf{z}) \stackrel{\text{def}}{=} f_1(\mathbf{x}) f_2(\mathbf{y}),$$

where \mathbf{x} and \mathbf{y} are compatible with \mathbf{z} ; that is, $\mathbf{x} \sim \mathbf{z}$ and $\mathbf{y} \sim \mathbf{z}$.

Remark: Factor multiplication is commutative and associative.

Algorithm 2 `MultiplyFactors($f_1(\mathbf{X}_1), \dots, f_m(\mathbf{X}_m)$)`

input: $f_1(\mathbf{X}_1), \dots, f_m(\mathbf{X}_m)$: factors

output: a factor corresponding to the product $\prod_{i=1}^m f_i$

main:

- 1: $\mathbf{Z} \leftarrow \bigcup_{i=1}^m \mathbf{X}_i$
 - 2: $f \leftarrow$ a factor over variables \mathbf{Z} where $f(\mathbf{z}) = 1$ for all \mathbf{z}
 - 3: **for** each instantiation \mathbf{z} **do**
 - 4: **for** $i = 1$ to m **do**
 - 5: $\mathbf{x}_i \leftarrow$ instantiation of variables \mathbf{X}_i consistent with \mathbf{z}
 - 6: $f(\mathbf{z}) \leftarrow f(\mathbf{z}) f_i(\mathbf{x}_i)$
 - 7: **end for**
 - 8: **end for**
 - 9: **return** f
-

Complexity: $O(m \exp(|\mathbf{Z}|))$ of time & space.

The Process of Elimination

Consider again the Bayesian network $G2$:

And suppose we are interested in computing the joint probability distribution.

We can do this in **two ways**:

- Either use **chain rule** for BN for each (compatible) instantiation e.g., a, b, c, d, e

$$\Pr(a, b, c, d, e) = \theta_{e|c} \theta_{d|bc} \theta_{c|a} \theta_{b|a} \theta_a.$$

- Or interpreting each CPT as a factor and multiplying them getting the following factor:

$$\Theta_{E|C} \Theta_{D|BC} \Theta_{C|A} \Theta_{B|A} \Theta_A$$

which is indeed the joint probability distribution induced by $G2$.

Remark

This shows one of the key applications of factor multiplication: Express the joint distribution of any BN as a product of its CPTs.

Now, suppose we are interested in computing the marginal distribution on variables D and E .

As a combination of summing-out and multiplication, and we get:

$$\Pr(D, E) = \sum_{A,B,C} \Theta_{E|C} \Theta_{D|BC} \Theta_{C|A} \Theta_{B|A} \Theta_A.$$

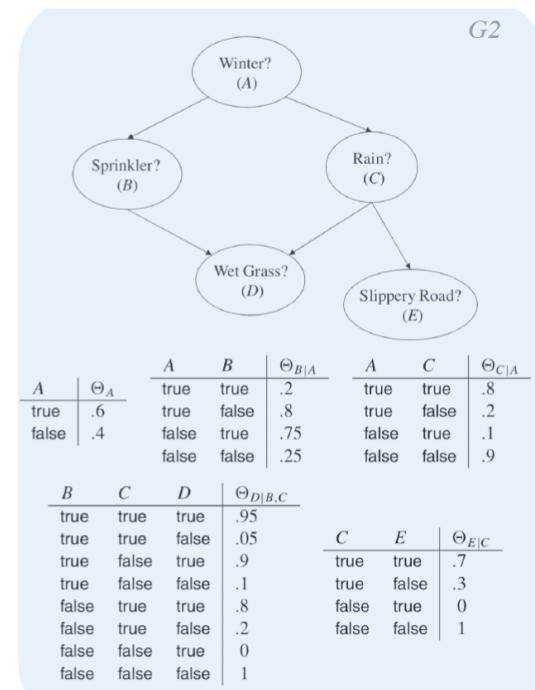
In fact, this can be done for any set of marginals.

Exponential complexity (in the size of the variables due to multiplying CPTs) is still on the table.

The following result shows that such multiplication is not necessary in general.

Theorem 3.1 If f_1 and f_2 are factors and if variable X appears only in f_2 , then

$$\sum_X f_1 f_2 = f_1 \sum_X f_2.$$



Why is this cool?

It implies that, if f_1, \dots, f_n are CPTs of a BN, to sum-out X , it may not be necessary to multiply these factors, first.

That is, if X only appears in f_n , we don't need to multiply any of these factors before one sums-out X since

$$\sum_X f_1 \dots f_n = f_1 \dots f_{n-1} \sum_X f_n.$$

Or, if X only appears in f_{n-1} and f_n , then one only needs to multiply these two factors before summing-out X since

$$\sum_X f_1 \dots f_n = f_1 \dots f_{n-2} \sum_X f_{n-1} f_n.$$

In general, to sum-out X from the product f_1, \dots, f_n , we need to just multiply the factors f_k that include X before summing out X from the resulting factor $\prod_k f_k$.

Remember the idea: To sum-out X from the product f_1, \dots, f_n , we need to just multiply the factors f_k that include X before summing out X from the resulting factor $\prod_k f_k$.

Example: Consider the following BN $G4$

Let's say we want to compute $\Pr(C)$ (*the prior marginal on C*), by first eliminating A and then B .

		A		B		C		G4
				$\Theta_{B A}$				
		A	B	true	false	true	false	$\Theta_{C B}$
				.9	.1			
		true	true			.3		
		true	false			.7		
		false	true			.5		
		false	false			.5		
				.2	.8			

Note that there are two factors that mention variable A : Θ_A and $\Theta_{B|A}$. We must multiply them first, then sum out variable A .

Eliminating A

Multiplying Θ_A and $\Theta_{B|A}$, we get:

A	B	$\Theta_A \Theta_{B A}$
true	true	.54
true	false	.06
false	true	.08
false	false	.32

Summing out variable A , we get:

B	$\sum_A \Theta_A \Theta_{B A}$
true	.62 = .54 + .08
false	.38 = .06 + .32

Eliminating B

Since both factors have B , we multiply them and get:

B	C	$\Theta_{C B} \sum_A \Theta_A \Theta_{B A}$
true	true	.186
true	false	.434
false	true	.190
false	false	.190

Summing out:

C	$\sum_B \Theta_{C B} \sum_A \Theta_A \Theta_{B A}$
true	.376
false	.624

$\Pr(C)$

We now have factors the $\sum_A \Theta_A \Theta_{B|A}$ and $\Theta_{C|B}$, and we want to eliminate B .

Computing Prior Marginals

COMPUTING PRIOR MARGINALS

Observe: All the crucial work is on Line 3 and 4.

Recall that in the previous examples, we have eliminated first A and then B , yielded single variable factors:

$$\sum_B \Theta_{C|B} \underbrace{\sum_A \Theta_A \Theta_{B|A}}_1 .$$

store CPTs

eliminate variable $\pi(i)$

$\Pr(\mathbf{Q})$

Algorithm 3 VE_PR1($\mathcal{N}, \mathbf{Q}, \pi$)

input:

\mathcal{N} : Bayesian network

\mathbf{Q} : variables in network \mathcal{N}

π : ordering of network variables not in \mathbf{Q}

output: the prior marginal $\Pr(\mathbf{Q})$

main:

```

1:  $\mathcal{S} \leftarrow$  CPTs of network  $\mathcal{N}$ 
2: for  $i = 1$  to length of order  $\pi$  do
3:    $f \leftarrow \prod_k f_k$ , where  $f_k$  belongs to  $\mathcal{S}$  and mentions variable  $\pi(i)$ 
4:    $f_i \leftarrow \sum_{\pi(i)} f$ 
5:   replace all factors  $f_k$  in  $\mathcal{S}$  by factor  $f_i$ 
6: end for
7: return  $\prod_{f \in \mathcal{S}} f$ 

```

Assume instead we eliminated B first and A afterwards.

This yields:

$$\sum_A \Theta_A \underbrace{\sum_B \Theta_{B|A} \Theta_{C|B}}_2 .$$

Now, inner-part having two variables, costing significantly more computation (both in terms of space and time).

Bottom-line: Order matters!

Choosing An Elimination Order

Computing Prior Marginals

Remember: All the crucial work is on Line 3 and 4.

Theorem 3.2: Given a BN, if the largest built factor on Line 4 has w variables, then the complexity of Lines 3-5 is $O(n \exp(w))$ where n is the number of variables in the BN.

Quality of order

The number w is known as the *width* of used order π and is taken as a measure of the *order quality*.

Idea: Choose an order that has the smallest possible width.

Remark: Note that total time complexity of the algorithm:

$O(n \exp(w) + n \exp(|\mathbf{Q}|))$.

Algorithm 3 VE_PR1($\mathcal{N}, \mathbf{Q}, \pi$)

input:

\mathcal{N} : Bayesian network

\mathbf{Q} : variables in network \mathcal{N}

π : ordering of network variables not in \mathbf{Q}

output: the prior marginal $\Pr(\mathbf{Q})$

main:

```

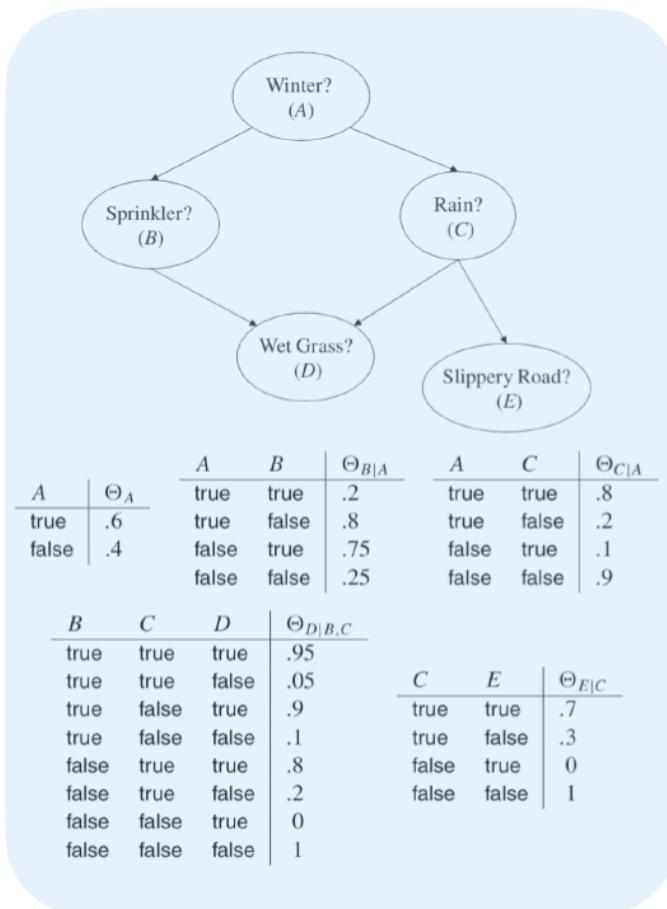
1:  $\mathcal{S} \leftarrow$  CPTs of network  $\mathcal{N}$ 
2: for  $i = 1$  to length of order  $\pi$  do
3:    $f \leftarrow \prod_k f_k$ , where  $f_k$  belongs to  $\mathcal{S}$  and mentions variable  $\pi(i)$ 
4:    $f_i \leftarrow \sum_{\pi(i)} f$ 
5:   replace all factors  $f_k$  in  $\mathcal{S}$  by factor  $f_i$ 
6: end for
7: return  $\prod_{f \in \mathcal{S}} f$ 

```

Choosing an Order of Elimination

Observe: A straightforward idea is to track the degrees of each order.

i	$\pi(i)$	\mathcal{S}	f_i	w
		$\Theta_A \Theta_{B A} \Theta_{C A} \Theta_{D BC} \Theta_{E C}$		
1	B	$\Theta_A \Theta_{C A} \Theta_{E C} f_1(A, C, D)$	$f_1 = \sum_B \Theta_{B A} \Theta_{D BC}$	3
2	C	$\Theta_A f_2(A, D, E)$	$f_2 = \sum_C \Theta_{C A} \Theta_{E C} f_1(A, C, D)$	3
3	A	$f_3(D, E)$	$f_3 = \sum_A \Theta_A f_2(A, D, E)$	2
4	D	$f_4(E)$	$f_4 = \sum_D f_3(D, E)$	1



Interaction Graphs

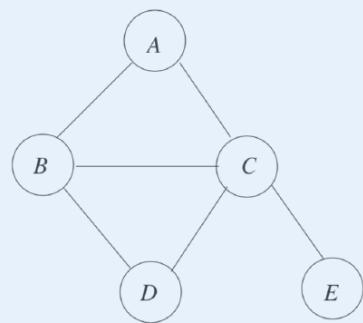
It is simpler to visualize such width with interaction graphs which expresses the interactions between CPTs of a BN.

The interaction graph G of the given factors , is an undirected graph constructed as follows:
 f_1, \dots, f_n

- The nodes of G are the variables that appear in factors $.f_1, \dots, f_n$
- Edges connect the variables which appear in the same factor.

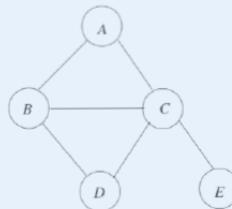
Example:

$$\mathcal{S}_1 : \Theta_A \quad \Theta_{B|A} \quad \Theta_{C|A} \quad \Theta_{D|BC} \quad \Theta_{E|C}$$

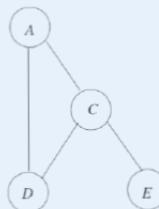


Example:

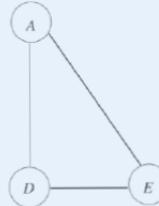
$$\mathcal{S}_1 : \Theta_A \quad \Theta_{B|A} \quad \Theta_{C|A} \quad \Theta_{D|BC} \quad \Theta_{E|C}$$



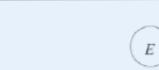
$$\mathcal{S}_2 : \Theta_A \quad \Theta_{C|A} \quad \Theta_{E|C} \quad f_1(A, C, D)$$



$$\mathcal{S}_3 : \Theta_A \quad f_2(A, D, E)$$



$$\mathcal{S}_4 : f_3(D, E)$$



$$\mathcal{S}_5 : f_4(E)$$



Observations:

- Eliminating a variable $\pi(i)$ from S leads to constructing a factor over the neighbours of $\pi(i)$.

- For $S \xrightarrow{\text{eliminating } \pi(i)} S'$ which corresponds to $G \xrightarrow{\text{eliminating } \pi(i)} G'$

G' can be obtained by G as follows:

1. Add an edge to G between every pair of neighbours of variable $\pi(i)$ that are not already connected by an edge.
 2. Delete variable $\pi(i)$ from G .
- Multiplying all factors containing $\pi(i)$ in S .
- Summing out variable $\pi(i)$ from the resulting factor.

Also observe that each set of variables per factor forms a *clique*.

Algorithm 4 OrderWidth(\mathcal{N}, π)

input:

\mathcal{N} : Bayesian network
 π : ordering of the variables in network \mathcal{N}

output: the width of elimination order π

main:

- 1: $G \leftarrow$ interaction graph of the CPTs in network \mathcal{N}
 - 2: $w \leftarrow 0$
 - 3: **for** $i = 1$ to length of elimination order π **do**
 - 4: $w \leftarrow \max(w, d)$, where d is the number of $\pi(i)$'s neighbors in G
 - 5: add an edge between every pair of non-adjacent neighbors of $\pi(i)$ in G
 - 6: delete variable $\pi(i)$ from G
 - 7: **end for**
 - 8: **return** w
-

A simple heuristic: Always choose the node with the smallest degree.

Algorithm 5 MinDegreeOrder(\mathcal{N}, \mathbf{X})

input:

\mathcal{N} : Bayesian network
 \mathbf{X} : variables in network \mathcal{N}

output: an ordering π of variables \mathbf{X}

main:

- 1: $G \leftarrow$ interaction graph of the CPTs in network \mathcal{N}
 - 2: **for** $i = 1$ to number of variables in \mathbf{X} **do**
 - 3: $\pi(i) \leftarrow$ a variable in \mathbf{X} with smallest number of neighbors in G
 - 4: add an edge between every pair of non-adjacent neighbors of $\pi(i)$ in G
 - 5: delete variable $\pi(i)$ from G and from \mathbf{X}
 - 6: **end for**
 - 7: **return** π
-

Another heuristic: Always choose the node whose elimination adds smallest number of edges.

Algorithm 6 MinFillOrder(\mathcal{N}, \mathbf{X})

input:

\mathcal{N} : Bayesian network
 \mathbf{X} : variables in network \mathcal{N}

output: an ordering π of variables \mathbf{X}

main:

- 1: $G \leftarrow$ interaction graph of the CPTs in network \mathcal{N}
 - 2: **for** $i = 1$ to number of variables in \mathbf{X} **do**
 - 3: $\pi(i) \leftarrow$ a variable in \mathbf{X} that adds the smallest number of edges on Line 4
 - 4: add an edge between every pair of non-adjacent neighbors of $\pi(i)$
 - 5: delete variable $\pi(i)$ from G and from \mathbf{X}
 - 6: **end for**
 - 7: **return** π
-

Computing Posterior Marginals

Now assume that given a BN , \mathcal{N} , a set of variables \mathbf{Q} , and an instantiation e , we want to compute the posterior marginal $\Pr(\mathbf{Q}|e)$ for variables \mathbf{Q} .

Example:

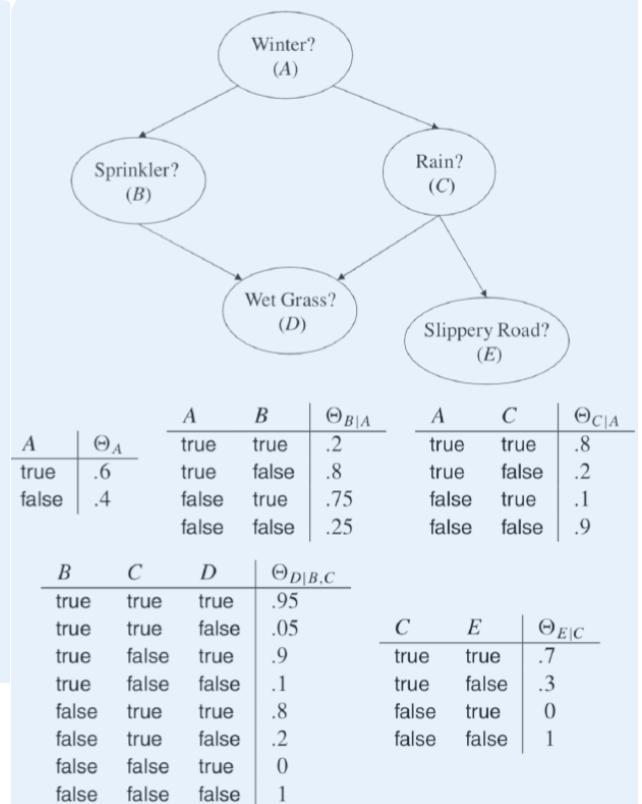
Now considering G_2 , assume that we have $\mathbf{Q} = \{D, E\}$ and evidence for e : $A = \text{true}$, $B = \text{false}$.

We would like to compute the following factor:

D	E	$\Pr(\mathbf{Q} e)$
true	true	.448
true	false	.192
false	true	.112
false	false	.248

For instance, 3rd row corresponds to:

$$\Pr(D = \text{false}, E = \text{true} | A = \text{true}, B = \text{false}) = .112$$



Remark: Note that prior marginal is a special case of posterior marginal where e is empty.

Example (cont'd):

It might be more straightforward to compute a variation on posterior marginals called *joint marginals*, $\text{Pr}(\mathbf{Q}, \mathbf{e})$.

That is, instead of computing, say, $\text{Pr}(\mathbf{q}|\mathbf{e})$, we compute $\text{Pr}(\mathbf{q}, \mathbf{e})$

D	E	$\text{Pr}(\mathbf{Q}, \mathbf{e})$
true	true	.21504
true	false	.09216
false	true	.05376
false	false	.11904

Goal: Computing $\text{Pr}(\mathbf{Q}, \mathbf{e})$ and normalising by \mathbf{e} .

In doing so, we can start by zeroing out the rows that are inconsistent with \mathbf{e} .

Factor given evidence: Given a factor f over \mathbf{X} , and an evidence \mathbf{e} , reduced factor $f^{\mathbf{e}}$ over \mathbf{X} is defined as follows:

$$f^{\mathbf{e}}(\mathbf{x}) \stackrel{\text{def}}{=} \begin{cases} f(\mathbf{x}), & \text{if } \mathbf{x} \sim \mathbf{e} \\ 0, & \text{otherwise.} \end{cases}$$

Example:

For instance, 3rd row corresponds to:

$$\text{Pr}(D=\text{false}, E=\text{true}, A=\text{true}, B=\text{false}) = .05376$$

Note that if we keep on summing up, we end up with the probability of evidence \mathbf{e} : $A=\text{true}, B=\text{false}$ (since case analysis).

D	E	f	D	E	$f^{\mathbf{e}}$
true	true	.448	true	true	.448
true	false	.192	true	false	0
false	true	.112	false	true	.112
false	false	.248	false	false	0

This means, we can compute the posterior marginal $\text{Pr}(\mathbf{Q}|\mathbf{e})$ simply by normalising $\text{Pr}(\mathbf{Q}, \mathbf{e})$ by \mathbf{e} which is obtained for free.

We can omit zeroed rows:

D	E	$f^{\mathbf{e}}$
true	true	.448
false	true	.112

Example (con'td):

$\mathbf{Q} = \{D, E\}$ and $\mathbf{e} : A = \text{true}, B = \text{false}$.

$$\text{Pr}(\mathbf{Q}, \mathbf{e}) = \sum_{A, B, C} (\Theta_{E|C} \Theta_{D|BC} \Theta_{C|A} \Theta_{B|A} \Theta_A)^{\mathbf{e}}.$$

Realise that we have to multiply all CPTs before we start eliminating variables, which is undesirable.

The following result helps.

Theorem 3.3: If f_1 and f_2 are two factors and \mathbf{e} is an instantiation, then $(f_1 f_2)^{\mathbf{e}} = f_1^{\mathbf{e}} f_2^{\mathbf{e}}$.

Hence, by the theorem it reduces to:

$$\text{Pr}(\mathbf{Q} = \{D, E\}, \mathbf{e}) = \sum_{A, B, C} \Theta_{E|C}^{\mathbf{e}} \Theta_{D|BC}^{\mathbf{e}} \Theta_{C|A}^{\mathbf{e}} \Theta_{B|A}^{\mathbf{e}} \Theta_A^{\mathbf{e}}$$

hence, *variable elimination algorithm* can be used.

Example: Consider G4 once again.

Suppose $Q = \{C\}$, $e : A=\text{true}$,

Goal: $\Pr(Q, e)$ by eliminating A first, B second.



G4

		Θ_A
		.9
A	true	.6
	false	.4

		$\Theta_{B A}$			$\Theta_{C B}$
		.9			.3
B	true	.1	true	false	.7
	false	.8	false	true	.5

Reducing the CPTs given evidence e :

		Θ_A^e
		.6
A	true	.6
	false	.4

		$\Theta_{B A}^e$
		.9
A	true	.9
	false	.1

		$\Theta_{C B}^e$
		.3
B	true	.3
	false	.7

The formula we need to evaluate then becomes

$$\begin{aligned}\Pr(Q, e) &= \sum_B \sum_A \Theta_A^e \Theta_{B|A}^e \Theta_{C|B}^e \\ &= \sum_B \Theta_{C|B}^e \sum_A \Theta_A^e \Theta_{B|A}^e.\end{aligned}$$

Example (cont'd): Consider G4 once again.

Suppose $Q = \{C\}$, $e : A=\text{true}$,

Goal: $\Pr(Q | e)$ by eliminating A first, B second.

$$\Pr(Q, e) = \sum_B \Theta_{C|B}^e \sum_A \Theta_A^e \Theta_{B|A}^e.$$

All intermediate factors needed to evaluate this formula are shown here:

		$\Theta_A^e \Theta_{B A}^e$			$\sum_A \Theta_A^e \Theta_{B A}^e$
		.54			.54
A	true	.54	B	true	.54
	false	.06		false	.06

		$\Theta_{C B}^e \sum_A \Theta_A^e \Theta_{B A}^e$			$\sum_B \Theta_{C B}^e \sum_A \Theta_A^e \Theta_{B A}^e$
		.162			.192
B	true	.162	C	true	.192
	false	.378		false	.408

Therefore,

$$\Pr(C=\text{true}, A=\text{true}) = .192$$

$$\Pr(C=\text{false}, A=\text{true}) = .408$$

$$\Pr(A=\text{true}) = .600$$

To compute, $\Pr(C | A=\text{true})$, we normalise by $\Pr(A=\text{true})$.

		$\Pr(C A=\text{true})$
		.32
C	true	.32
	false	.68

All we explained so far (except normalisation by e):

Algorithm 7 VE_PR2($\mathcal{N}, \mathbf{Q}, \mathbf{e}, \pi$)

input:

\mathcal{N} : Bayesian network
 \mathbf{Q} : variables in network \mathcal{N}
 \mathbf{e} : instantiation of some variables in network \mathcal{N}
 π : an ordering of network variables not in \mathbf{Q}

output: the joint marginal $\Pr(\mathbf{Q}, \mathbf{e})$

main:

```

1:  $\mathcal{S} \leftarrow \{f^{\mathbf{e}} : f \text{ is a CPT of network } \mathcal{N}\}$ 
2: for  $i = 1$  to length of order  $\pi$  do
3:    $f \leftarrow \prod_k f_k$ , where  $f_k$  belongs to  $\mathcal{S}$  and mentions variable  $\pi(i)$ 
4:    $f_i \leftarrow \sum_{\pi(i)} f$ 
5:   replace all factors  $f_k$  in  $\mathcal{S}$  by factor  $f_i$ 
6: end for
7: return  $\prod_{f \in \mathcal{S}} f$ 

```

The only difference to
to the previous algorithm
(i.e., VE_PR1).

Canvas Quiz with Answers

1.

LAS8 In a proof by contradiction, such as DPLL or tableau, I can prove that a formula F is entailed by a knowledge base KB by showing that

A: the knowledge base KB and the negation of the formula are unsatisfiable.

the formula F is unsatisfiable, which implies that it must be entailed by the KB.

the knowledge base is unsatisfiable, which implies that the formula F must be entailed.

the knowledge base KB and the formula F are together unsatisfiable

2.

LAS5 A calculus is called complete w.r.t. the semantics of a logic if, and only, if

all the formulas it proves are semantically entailed.

A: it can prove all the semantically entailed formulas.

it proves all the correct formulas in finite time.

all the formulas it proves are tautologies.

3.

LAS4 A formula is in clause normal form if it is a A: conjunction of A: disjunctions of A:
literals

CNF

4.

AS6 The semantics (meaning) of a formula in Propositional Logic is determined as:

multiple choice

a numeric value

set membership

A: a truth value

5.

Consider the sentence

$$((\neg A \vee \neg B) \rightarrow (A \rightarrow \neg B)) \wedge (A \vee B)$$

Which of the following statements are true?

The sentence is neither valid, nor satisfiable, nor a contradiction

A: The sentence is satisfiable, but not valid

The sentence is valid, but not satisfiable

The sentence is not valid, and thus also not satisfiable

The sentence is not valid, and thus a contradiction

6.

LAC2 In the following truth table X1, X2 and X3 stand for possible truth values:

A	B	$B \rightarrow A$	$A \rightarrow (B \rightarrow A)$
True	True	X1(T)	T
True	False	T	X2 (T)
False	True	F	X3 (T)
False	False	T	T

Which of the following statements is correct (multiple answers possible):

$A \rightarrow (B \rightarrow A)$ is valid because there is a model that interprets it as true

X1 = True

X1 = False

$A \rightarrow (B \rightarrow A)$ is not valid but satisfiable

A $\rightarrow (B \rightarrow A)$ is valid because all possible models interpret it as true

X2 = True

X2 = False

7.

LAC1 Match formulas that are logically equivalent.

$A \rightarrow (B \ \& \ C)$

A: $(\neg A \vee B) \ \& \ (\neg A \vee C)$

$A \vee (B \ \& \ C)$

A: $(A \vee B) \ \& \ (A \vee C)$

$A \rightarrow \neg(B \ \& \ C)$

A: $\neg A \vee \neg B \vee \neg C$

$A \ \& \ \neg(B \ \& \ C)$

A: $A \ \& \ (\neg B \vee \neg C)$

$A \ \& \ (B \rightarrow C)$

A: $(A \ \& \ \neg B) \vee (A \ \& \ C)$

$A \rightarrow (B \ \& \ C)$ is equivalent to $\neg A \vee (B \ \& \ C)$ by definition of \rightarrow , which is equivalent to $(\neg A \vee B) \ \& \ (\neg A \vee C)$ by rule 2 of slide 23

$A \vee (B \ \& \ C)$ is $(A \vee B) \ \& \ (A \vee C)$ by rule 2 on slide 23.

$A \rightarrow \neg(B \ \& \ C)$ is equivalent to $\neg A \vee \neg(B \ \& \ C)$ by definition of \rightarrow . $\neg(B \ \& \ C)$ is equivalent to $\neg B \vee \neg C$ by rule 1.

$A \vee \neg(B \ \& \ C)$ applies rule 1 again.

8.

Weighted partial MAXSAT formulas

both hard and soft clauses are allowed. Subsumes all previously mentioned variations.

Consider the following statements about geese

1. all geese are white X
2. geese often have two legs Y
3. It is very likely that a goose is either white or has two legs or both $(X \vee Y) \vee (X \wedge Y)$

4. if a goose does not have wings, it cannot fly. ($\neg W \vee Z$)

and the following variables:

W stands for goose has wings

X stands for goose is white

Y stands for goose has two legs

Z stands for goose can not fly

Which of the following statements is a faithful representation of the knowledge described above?

$$F = (X, \infty) \wedge (X \vee Y, 0.4) \wedge (Y, -5) \wedge (Z \vee W, \infty)$$

$$F = (X, \infty) \wedge (\neg(X \vee Y), 0.4) \wedge (Y, 5) \wedge (\neg(Z \vee W), \infty)$$

$$A: F = (\neg X, \infty) \wedge (X \vee Y, 0.4) \wedge (\neg Y, 5) \wedge (Z \vee W, \infty)$$

$$F = X \wedge (X \vee Y, 0.4) \wedge (\neg Y, 5) \wedge (Z \vee W)$$

9.

Use DPLL procedure to prove or disprove satisfiability of the formula

$$(X \vee Y \vee Z) \wedge (X \vee \neg Y) \wedge (Y \vee \neg Z) \wedge (Z \vee \neg X) \wedge (\neg X \vee \neg Y \vee \neg Z)$$

Label each step which part of the algorithm you have used

1. Simplification

Tautology: remove tautologies like $P \vee \neg P$ from knowledge base (once in the beginning)

$$(X \vee Y \vee Z)$$

2. Split

3. pick a predicate and assume a truth value
4. – Heuristics of which literal to pick next can improve the efficiency of DPLL a lot
5. – DLCS (Dynamic Largest Combined Sum): Pick v with the largest count of positive and negative occurrences: $CP(v) + CN(v)$. If $CP(v) > CN(v)$, choose $v = 1$, else $v = 0$
6. – DLIS (Dynamic Largest Individual Sum): Pick v with either largest CP or CN. Same truth
8. assignment as for DLCS.
9. – Jeroslow-Wang: weight of literal depends on the length of clauses it occurs in. Thereby, we
10. prefer small clauses. The score of a literal is $J(v) =$
11. P
12. $c_2 C_v$
13. $2^{-|c_j|}$, and we pick the highest value.
14. One-sided JW looks at v and $\neg v$ independently, whereas the two-sided approach looks at sum
15. of v and $\neg v$, and just picks the truth value based on $J(v) - J(\neg v)$ $v = 1$, else $v = 0$.
16. – MOMs (Maximum Occurrences in Clauses of Minimum Size): similar to JW, but we only look at

17. the smallest clauses in the knowledge base. The number of occurrences in those is indicated by
18. the function f . We choose the literal that maximizes $[f(v) + f(\neg v)] - 2k + f(v) - f(\neg v)$. k is
19. a tuning parameter for the trade-off between the balanced distribution of v and $\neg v$ and their individual ones.

- Pure literals: set predicates that solely occur in their positive or negative form to the corresponding truth value
- Unit clauses: set literals for which the knowledge base contains a unit clause to true (or the predicate to false respectively)

10.

Give a pseudocode description of GSAT

Note that this algorithm tends to get stuck in local minimum (flipping a single variable does not increase score). Thus, we perform random restarts to start new. That's also why GSAT spends most time on plateaus where score is not improved

```

procedure GSAT(Sigma)
  for i := 1 to MAX-TRIES ; These are the restarts
    T := random(Sigma) ; random assignment
    for j := 1 to MAX-FLIPS ; To ensure termination
      if T satisfies Sigma then return T
      else T := T with variable flipped to maximize
            number of satisfied clauses
      ; It doesn't matter if the number does
      ; not increase. This are the sideways flips
    end
  end

```

11.

What are advantages of a short restart and a long restart?

More than 1 answer might be correct.

Advantage long restart: if we've randomly jumped to a part of the search space without solutions, we will waste too much time on local walks

Advantage short restart: if we've randomly jumped to a part of the search space close to a solution, we might still miss it because we didn't do enough local exploration

A: Advantage short restart: if we've randomly jumped to a part of the search space without solutions, we won't waste too much time on local walks

A: Advantage long restart: if we've randomly jumped to a part of the search space close to a solution, we are likely to actually find the solution because we'll do enough local exploration

12.

In Description Logics, a model is an interpretation function that assigns a value true or false to each concept name.

- ontologies

-

True or false?

True

False

13.

In Description Logics a concept is called unsatisfiable if there is no model.

True or false?

True

False

14.

Description Logics are truly more expressive than Propositional Logics (i.e. everything that can be said in PL can be said in all Description Logics).

True or false?

True

False

15.

Suppose that the concept names to use are Kid and Ice-Cream, and the role name is likes.

Translate the following sentence from DL ALC to English

- [∃ \(Links to an external site.\)](#) likes.Ice-Cream \sqsubseteq Kid

A: Everybody who likes ice-cream is a kid

there is a kid that likes ice-cream

every kid likes ice-cream

everybody who only likes ice-cream is a kid

16.

Use a tableau to prove whether $(\forall r.A \sqcap \forall r.B) \sqsubseteq \forall r.(A \sqcap B)$

Explain the individual steps.

17.

1. Every Bayesian network implicitly represents a unique probability distributions.
2. True or false?

A: True

False

18.

Given a Bayesian Network G, if two random variables X and Y are independent (w.r.t. to the probability distribution induced by G) given a third random variable Z, then it follows that X and Y are d-separated by Z.

True or false?

True

A: False

19.

A divergent valve ($\leftarrow W \rightarrow$) is closed iff neither variable W nor any of its descendants appears in Z .

True or false?

True

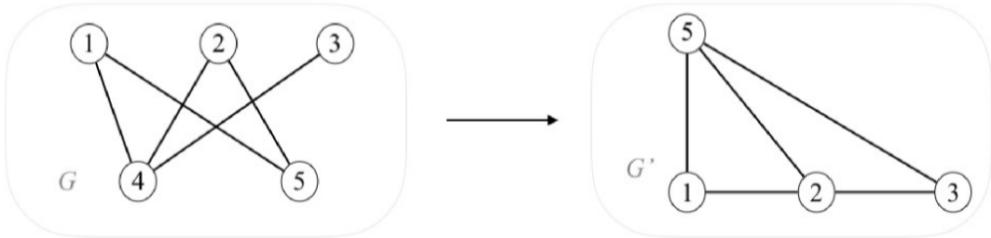
False

20.

Q5 - Check out the statements below whether they are true or not, and label them accordingly with either “true” or “false”.

- a) If $\alpha \models \beta$ and $\Pr(\alpha)=0$, then $\Pr(\beta)=0$.
- b) If $\alpha \models \beta \models \gamma$, then $\Pr(\alpha \mid \beta) \geq \Pr(\alpha \mid \gamma)$.

- c) In the interaction graph G below, if we eliminate 4, we get G'



- d) Assume G' above, according to the order elimination heuristic of “choosing the node with smallest degree”, one possible elimination order is: 1, 5, 3 ,2.

Group of answer choices

A: c) false

a) true

c) true

b) false

A: b) true

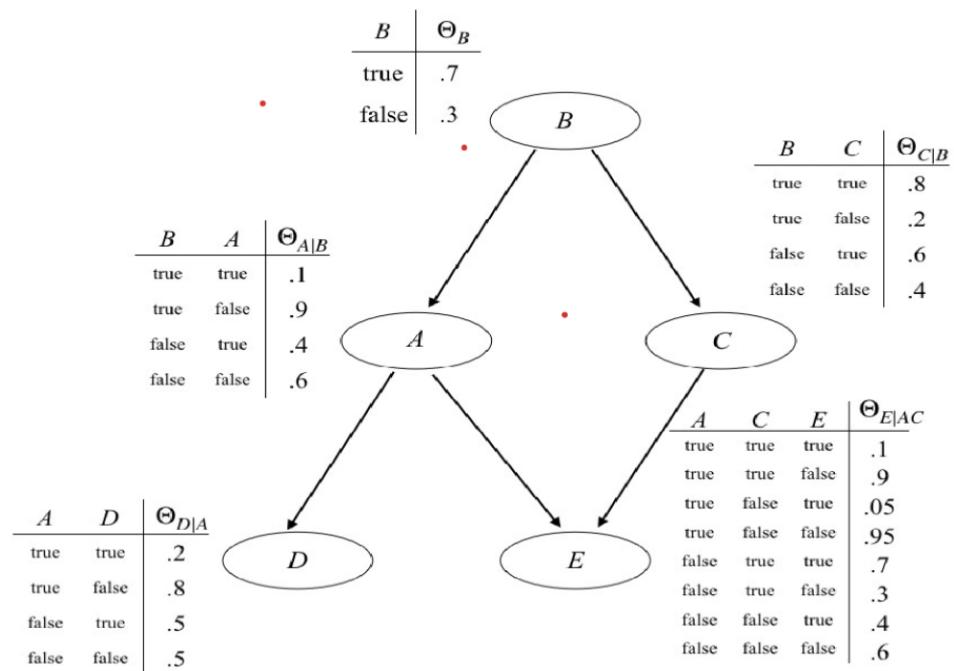
A: d) true

A: a) false

d) false

21.

Q6 - Consider the Bayesian network in the figure, and answer the questions below accordingly.



Are the following statements true or false?

$\text{IPr}(D, \emptyset, C)$ is True

A: $\text{IPr}(C, \{B, E\}, A)$ is False

A: $\text{IPr}(E, \{A\}, B)$ is False

$\text{IPr}(E, \{A\}, B)$ is True

A: $\text{IPr}(D, \emptyset, C)$ is False

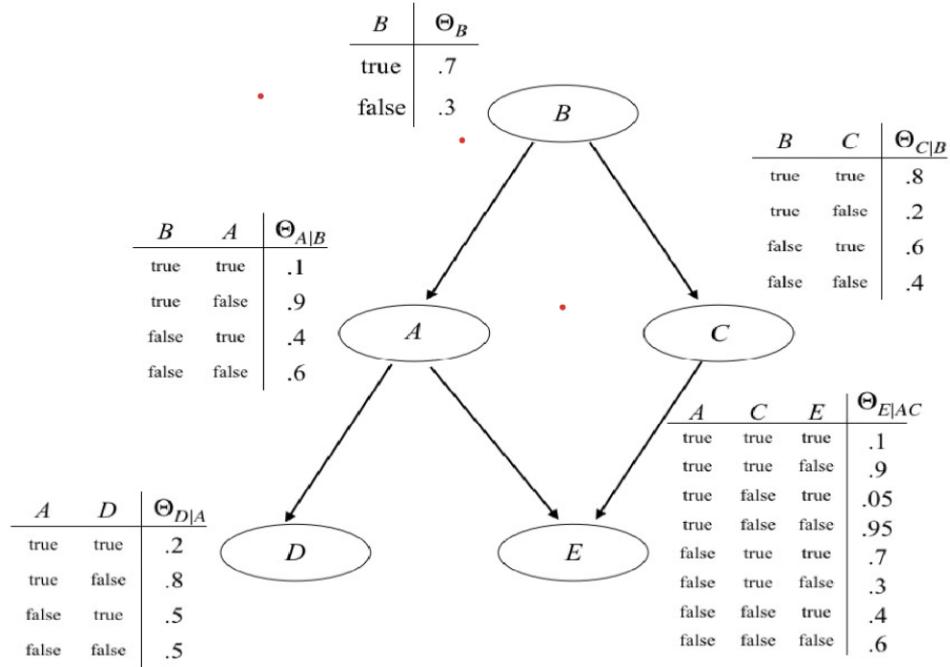
$\text{IPr}(C, \{B, E\}, A)$ is True

$\text{IPr}(\{BD\}, \{A, C\}, E)$ is False

A: $\text{IPr}(\{BD\}, \{A, C\}, E)$ is True

22.

Q6 - Consider the Bayesian network in the figure, and answer the questions below accordingly.

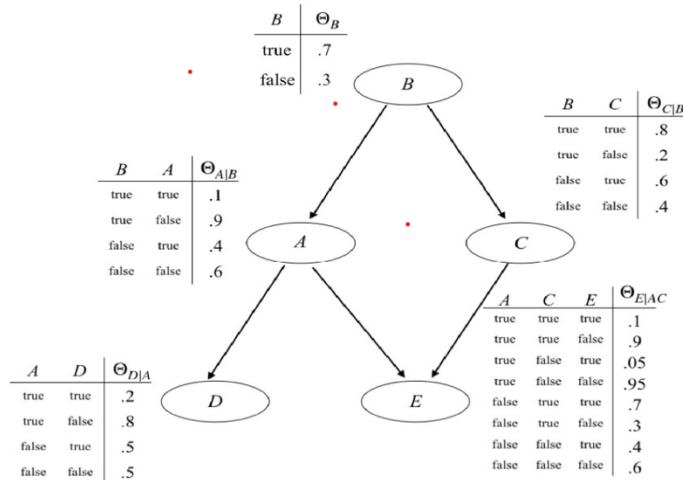


For this BN (as before), calculate the values for $\Pr(A=\text{true}, B=\text{false}, E=\text{true})$

A: 0.0096

23.

Q6 - Consider the Bayesian network in the figure, and answer the questions below accordingly.



Calculate the value for $\Pr(A=\text{false}, B=\text{true} | C=\text{true})$

A: 0.681

Probabilistic Graphical Models 4

Probabilistic Graphical Models 4

In this lecture, we will take a look at the problem of finding variable instantiations that have maximal probability under some given evidence.

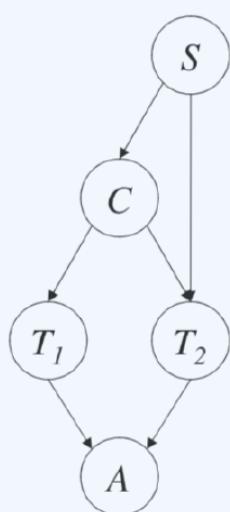
Consider the following exemplary scenario:

- A society (S) that is 55% male and 45% female.
- Members of this population can suffer from a medical condition (C) that is more likely in males.
- Two diagnostic tests are available for detecting this condition, T_1 and T_2 .
- T_2 is more effective on females while both tests are equally effective on males.

Example Scenario:

- A society (S) that is 55% male and 45% female.
- Members of this population can suffer from a medical condition (C) that is more likely in males.
- Two diagnostic tests are available for detecting this condition, T_1 and T_2 .
- T_2 is more effective on females while both tests are equally effective on males.

One possible modelling of this scenario:



A: T_1 and T_2 agrees.

S	C	T_1	$\theta_{c s}$	C	T_1	$\theta_{t_1 c}$
male	yes	+ve	.05	yes	+ve	.80
male	no	-ve	.95	yes	-ve	.20
female	yes	+ve	.01	no	+ve	.20
female	no	-ve	.99	no	-ve	.80

S	C	T_2	$\theta_{t_2 c,s}$	T_1	T_2	A	$\theta_{a t_1,t_2}$
male	yes	+ve	.80	+ve	+ve	yes	1
male	yes	-ve	.20	+ve	+ve	no	0
male	no	+ve	.20	+ve	-ve	yes	0
male	no	-ve	.80	+ve	-ve	no	1
female	yes	+ve	.95	-ve	+ve	yes	0
female	yes	-ve	.05	-ve	+ve	no	1
female	no	+ve	.05	-ve	-ve	yes	1
female	no	-ve	.95	-ve	-ve	no	0

One can partition this population into four different groups:

whether they are male or female,
and whether having the condition or not.

- Suppose that all we know is *a person took both tests and the test results agree* (i.e., $A = \text{true}$).

Then we may ask: *What is the most likely group that this person belong to?*

Network Pruning

A Tool of Simplification: Network Pruning

Bayesian Network Simplification: Network Pruning

Network structure has a major impact on the performance of variable elimination (and also on the performance of most algorithms).

Intuition: The complexity of inference can be very much affected by the number and location of query and evidence variables within the network structure.

Network Pruning is done in two parts: Node Pruning and Edge Pruning.

Node Pruning

Idea: Given a BN N and a query (\mathbf{Q}, \mathbf{e}) , one can remove any leaf node (with its CPT) from the network as long as it does not belong to variables $\mathbf{Q} \cup \mathbf{E}$.

This does not affect the ability of the network to answer the query correctly, thanks to the following result.

Theorem: Let N be a Bayesian network and let (\mathbf{Q}, \mathbf{e}) be a corresponding query. If $N' = \text{pruneNodes}(N, \mathbf{Q} \cup \mathbf{E})$, then $\Pr(\mathbf{Q}, \mathbf{e}) = \Pr'(\mathbf{Q}, \mathbf{e})$ where \Pr and \Pr' are the probability distributions induced by networks N and N' , respectively.

Remark that this can be done iteratively, possibly leading to further simplifications.

Edge Pruning

Idea: Given a Bayesian network N and a query (Q, e) , one can eliminate some of the network edges and reduce some of its CPTs without affecting its ability to compute the joint marginal $\Pr(Q, e)$ correctly.

Procedure: For each edge $U \rightarrow X$ that originates from a node U in E

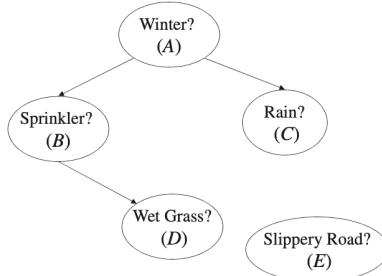
1. Remove the edge $U \rightarrow X$ from the network.
2. Replace the CPT $\Theta_{X|U}$ for node X by a smaller CPT, which is obtained from $\Theta_{X|U}$ by assuming the value u of parent U given in evidence e . This new CPT corresponds to $\sum_U \Theta_{X|U}^u$.

Theorem: Let N be a BN network and let e be an instantiation. If $N' = \text{pruneEdges}(N, e)$, then $\Pr(Q, e) = \Pr'(Q, e)$ where \Pr and \Pr' are the probability distributions induced by networks N and N' , respectively

Example: Pruning the network G w.r.t. evidence $C=\text{false}$

The two edges originating from node C were deleted and CPTs and were modified. $\Theta_{D|B,C}$ $\Theta_{E|C}$

Remark the pruned network is only good for answering queries of the form $\Pr(q, C = \text{false})$
Otherwise pruned network and the original network may disagree.

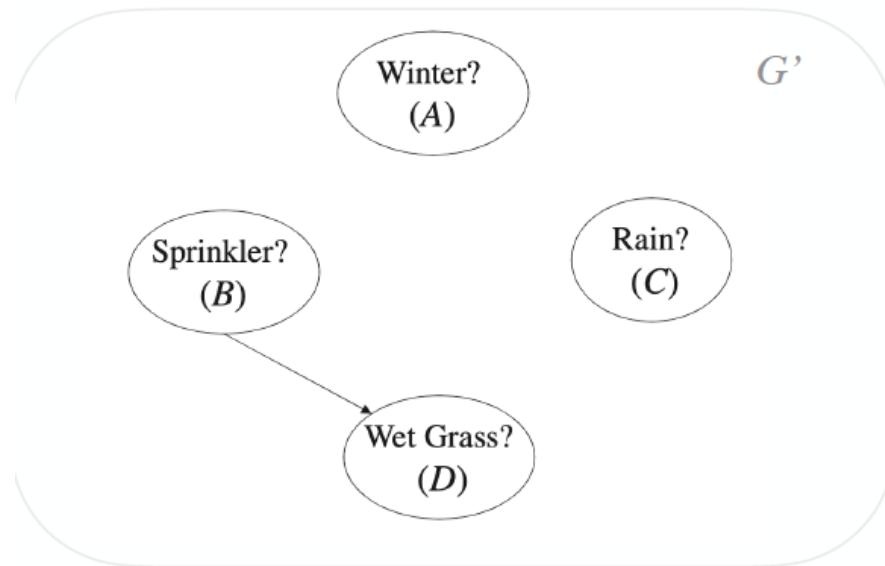


		A	B	$\Theta_{B A}$			A	C	$\Theta_{C A}$
		true	true	.2	true	true	.8		
		true	false	.8	true	false	.2		
		false	true	.75	false	true	.1		
		false	false	.25	false	false	.9		

		B	D	$\sum_C \Theta_{D BC}^{C=\text{false}}$			E	$\sum_C \Theta_{E C}^{C=\text{false}}$
		true	true	.9			true	0
		true	false	.1			false	1
		false	true	0				
		false	false	1				

Network Pruning: Example

Example: Pruning the network given the query $\mathbf{Q} = \{D\}$ and $\mathbf{e} : A=\text{true}, C = \text{false}$ Pruning leads to removing node E and the edges originating from nodes A and C, and modifying the CPTs, , and $\Theta_{B|A}$ $\Theta_{C|A}$ $\Theta_{D|BC}$



A	Θ_A	B	$\Theta'_B = \sum_A \Theta_{B A}^{A=\text{true}}$	C	$\Theta'_C = \sum_A \Theta_{C A}^{A=\text{true}}$
true	.6	true	.2	true	.8
false	.4	false	.8	false	.2
B	D		$\Theta'_{D B} = \sum_C \Theta_{D BC}^{C=\text{false}}$		
true	true		.9		
true	false		.1		
false	true		0		
false	false		1		

Most Likely Instantiations

MAP Queries

One can partition this society into four different groups:

Whether they are male or female,
and having the condition or not.

Suppose that all we know is:

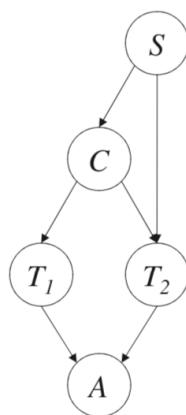
A person took both tests and the test results agree (i.e., $A = \text{yes}$).

Then we may ask things like:

What is the most likely group that this person belong to?

=

What are the most likely instantiation of variables S and C given evidence $A=\text{yes}$



S	θ_s			$\theta_{c s}$			$\theta_{t_1 c}$
		male	yes		yes	+ve	
male	.55	male	no	.95	yes	-ve	.20
		female	yes	.01	no	+ve	.20
female	.45	female	no	.99	no	-ve	.80

S	C	T_2	$\theta_{t_2 c,s}$	T_1	T_2	A	$\theta_{a t_1,t_2}$
male	yes	+ve	.80	+ve	+ve	yes	1
male	yes	-ve	.20	+ve	+ve	no	0
male	no	+ve	.20	+ve	-ve	yes	0
male	no	-ve	.80	+ve	-ve	no	1
female	yes	+ve	.95	-ve	+ve	yes	0
female	yes	-ve	.05	-ve	+ve	no	1
female	no	+ve	.05	-ve	-ve	yes	1
female	no	-ve	.95	-ve	-ve	no	0

What is the most likely group that this person belong to?

=

What are the most likely instantiation of variables S and C given evidence $A=\text{yes}$

$S=\text{male}$ and $C=\text{no}$ with a posterior probability of about 0.49.

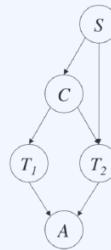
Such queries are known as **MAP** (*most a posteriori estimate*) **queries**,

And the answer to those queries are known as **MAP instantiations**.

MAP variables.

What are the most likely instantiation of variables S and C given evidence A=yes

$$\Pr(S=\text{male}, C=\text{no} \mid A = \text{yes}) \approx 0.49$$



S	C	$\theta_{c s}$	C	T1	$\theta_{t1 c}$
male	yes	.05	yes	+ve	.80
male	no	.95	yes	-ve	.20
female	yes	.01	no	+ve	.20
female	no	.99	no	-ve	.80

S	C	$\theta_{t1 c,s}$	T1	T2	A	$\theta_{a t1,t2}$
male	yes	.80	+ve	+ve	yes	1
male	yes	.20	+ve	+ve	no	0
male	no	.20	+ve	-ve	yes	0
male	-ve	.80	+ve	-ve	no	1
female	yes	.95	-ve	+ve	yes	0
female	-ve	.05	-ve	+ve	no	1
female	yes	.05	-ve	-ve	yes	1
female	-ve	.95	-ve	-ve	no	0

To calculate that, we can simply compute the posterior marginal (previous lecture) over MAP variables and then select the instantiation that has a maximum one.

A Special Case: MPE Queries

One can also consider the most likely instantiations of variables excluding the evidence (i.e., S, C, T₁, T₂):

$$S=\text{female}, \quad C=\text{no}, \quad T_1=-\text{ve}, \quad T_2=-\text{ve}, \quad \text{with the probability value of 0.47.}$$

Such queries are called **MPE queries** which can be considered as a special type of MAP queries.

Computing **MPE instantiations** is easier than computing MAP instantiations in general. (Hence, they have their own name.)

Remark:

Note that if we map the MPE to the previous variable, we likely to get a different answer

$$S=\text{female}, \quad C=\text{no}$$

which is different than the earlier answer:

$$S=\text{male}, \quad C=\text{no}$$

S	C	$\theta_{t1 c,s}$	T1	T2	A	$\theta_{a t1,t2}$
male	yes	.80	+ve	+ve	yes	1
male	-ve	.20	+ve	+ve	no	0
male	+ve	.20	+ve	-ve	yes	0
male	-ve	.80	+ve	-ve	no	1
female	yes	.95	-ve	+ve	yes	0
female	-ve	.05	-ve	+ve	no	1
female	+ve	.05	-ve	-ve	yes	1
female	-ve	.95	-ve	-ve	no	0

Remark:

Although likely to be different, sometimes this technique is used to approximate MAP instances.

MAP & MPE Queries

Formal Definition Given a Bayesian network with variables \mathbf{M}/\mathbf{Q} and evidence \mathbf{e} , MAP/MPE queries (values and instantiations) are defined as follows:

MAP	MPE
Value: $\text{MAP}_P(\mathbf{M}, \mathbf{e}) \stackrel{\text{def}}{=} \max_{\mathbf{m}} \Pr(\mathbf{m}, \mathbf{e}).$	Value: $\text{MPE}_P(\mathbf{e}) \stackrel{\text{def}}{=} \max_{\mathbf{q}} \Pr(\mathbf{q}, \mathbf{e}).$
Instantiation: $\text{MAP}(\mathbf{M}, \mathbf{e}) \stackrel{\text{def}}{=} \underset{\mathbf{m}}{\operatorname{argmax}} \Pr(\mathbf{m}, \mathbf{e}).$	Instantiation: $\text{MPE}(\mathbf{e}) \stackrel{\text{def}}{=} \underset{\mathbf{q}}{\operatorname{argmax}} \Pr(\mathbf{q}, \mathbf{e}).$

Remark:

Note that both MAP and MPE instantiations can also be characterised by $\Pr(\mathbf{q}|\mathbf{e})$ since, $\Pr(\mathbf{q}|\mathbf{e}) = \Pr(\mathbf{q}, \mathbf{e})/\Pr(\mathbf{e})$ and $\Pr(\mathbf{e})$ is independent of instantiation \mathbf{q} .

Computing MPE & MAP (by Variable Elimination)

Lets assume no evidence whatsoever, then

MPE instantiation:

	S	C	T_1	T_2	A	$\Pr(\cdot)$
31	female	no	-ve	-ve	yes	.338580

We can calculate this by variable elimination:

Idea: Instead of summing-out, we maximise-out.

Example: To maximise S out from the factor, $f(S, C, T_1, T_2, A)$ we construct a factor on the remaining variables : C, T_1, T_2, A

	S	C	T_1	T_2	A	$\Pr(\cdot)$
1	male	yes	+ve	+ve	yes	.017600
3	male	yes	+ve	-ve	no	.004400
5	male	yes	-ve	+ve	no	.004400
7	male	yes	-ve	-ve	yes	.001100
9	male	no	+ve	+ve	yes	.020900
11	male	no	+ve	-ve	no	.083600
13	male	no	-ve	+ve	no	.083600
15	male	no	-ve	-ve	yes	.334400
17	female	yes	+ve	+ve	yes	.003420
19	female	yes	+ve	-ve	no	.000180
21	female	yes	-ve	+ve	no	.000855
23	female	yes	-ve	-ve	yes	.000045
25	female	no	+ve	+ve	yes	.004455
27	female	no	+ve	-ve	no	.084645
29	female	no	-ve	+ve	no	.017820
31	female	no	-ve	-ve	yes	.338580

Even rows are omitted since they are all zeros e.g., $T_1 \neq T_2$ with $A = \text{yes}$.

	S	C	T_1	T_2	A	$\Pr(\cdot)$
1	male	yes	+ve	+ve	yes	.017600
17	female	yes	+ve	+ve	yes	.003420



C	T_1	T_2	A	$\Pr(\cdot)$
yes	+ve	+ve	yes	.017600 = max(.017600, .003420)

When we eliminate S , we call the new factor $\max_S f$.

Note that $\max_S f$ and f agree on MPE, hence it is as good as f in calculating this probability.

We continue until we are left with the trivial factor, and this gives us the MPE.

	S	C	T_1	T_2	A	$\Pr(\cdot)$
1	male	yes	+ve	+ve	yes	.017600
3	male	yes	+ve	-ve	no	.004400
5	male	yes	-ve	+ve	no	.004400
7	male	yes	-ve	-ve	yes	.001100
9	male	no	+ve	+ve	yes	.020900
11	male	no	+ve	-ve	no	.083600
13	male	no	-ve	+ve	no	.083600
15	male	no	-ve	-ve	yes	.334400
17	female	yes	+ve	+ve	yes	.003420
19	female	yes	+ve	-ve	no	.000180
21	female	yes	-ve	+ve	no	.000855
23	female	yes	-ve	-ve	yes	.000045
25	female	no	+ve	+ve	yes	.004455
27	female	no	+ve	-ve	no	.084645
29	female	no	-ve	+ve	no	.017820
31	female	no	-ve	-ve	yes	.338580

Even rows are omitted since they are all zeros e.g., $T_1 \neq T_2$ with $A = \text{yes}$.

Just like summing-out:

Maximising-Out

Given a factor f over variables \mathbf{X} , **maximising-out** variable $X \in \mathbf{X}$, from factor f results in another factor over variables $\mathbf{X} \setminus \{X\}$ and is defined as $\left(\max_X f \right)(\mathbf{y}) \stackrel{\text{def}}{=} \max_x f(x, \mathbf{y})$.

Maximising-out

$$\left(\max_X f \right)(\mathbf{y}) \stackrel{\text{def}}{=} \max_x f(x, \mathbf{y}).$$

Recall: Elimination continues until the trivial factor.

In case of MAP: We first sum-out all the non-MAP variables, and then maximise-out the MAP variables.

Remarks:

Similar to summation, maximisation is commutative,

Hence, no need to specify the order we maximise.

Moreover, it interacts with multiplication in a similar fashion to summation.

Theorem 4.1

If f_1 and f_2 are factors and if variable X appears only in f_2 , then

$$\max_X f_1 f_2 = f_1 \max_X f_2.$$

Extended Factors

We can compute MPE instantiations in addition to its probability:

Idea: The extended factors which assign to each instantiation both a number and an instantiation

Extended Factors

Let $f(\mathbf{X})$ be an extended factor, X be a variable in \mathbf{X} , and let $\mathbf{Y} = \mathbf{X} \setminus \{X\}$. We then define

$$\left(\max_X f \right) [\mathbf{y}] \stackrel{\text{def}}{=} x^* f[x^*, \mathbf{y}],$$

where $x^* = \operatorname{argmax}_{\mathbf{X}} f(\mathbf{x}, \mathbf{y})$. Moreover, let

$$\mathbf{Z} = \mathbf{X} \cup \mathbf{Y}$$

$$(f_1 f_2)[\mathbf{z}] \stackrel{\text{def}}{=} f_1[\mathbf{x}] f_2[\mathbf{y}],$$

where \mathbf{x} and \mathbf{y} are instantiation compatible with \mathbf{z} ($\mathbf{x} \sim \mathbf{z}$ and $\mathbf{y} \sim \mathbf{z}$).

MPE Algorithm

Algorithm 27 VE_MPE(\mathcal{N}, \mathbf{e})

input:

\mathcal{N} : a Bayesian network

\mathbf{e} : evidence

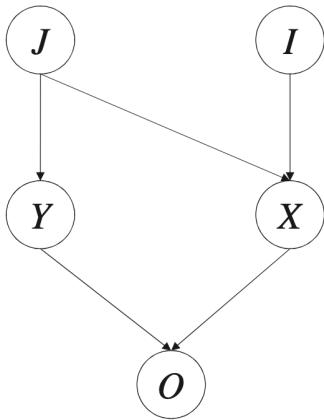
output: trivial factor f , where $f(\top)$ is the MPE probability of evidence \mathbf{e}

main:

- 1: $\mathcal{N}' \leftarrow \text{pruneEdges}(\mathcal{N}, \mathbf{e})$ {see Section 6.9.2}
 - 2: $\mathbf{Q} \leftarrow$ variables in network \mathcal{N}'
 - 3: $\pi \leftarrow$ elimination order of variables \mathbf{Q}
 - 4: $\mathcal{S} \leftarrow \{f^{\mathbf{e}} : f \text{ is a CPT of network } \mathcal{N}'\}$
 - 5: **for** $i = 1$ to $|\mathbf{Q}|$ **do**
 - 6: $f \leftarrow \prod_k f_k$, where f_k belongs to \mathcal{S} and mentions variable $\pi(i)$
 - 7: $f_i \leftarrow \max_{\pi(i)} f$
 - 8: replace all factors f_k in \mathcal{S} by factor f_i
 - 9: **end for**
 - 10: **return** trivial factor $\prod_{f \in \mathcal{S}} f$
-

Computing MPE

Example: Assume, we want to compute MPE instantiation given evidence $J=\text{true}$ and $O = \text{false}$.



I	Θ_I	J	Θ_J
true	.5	true	.5
false	.5	false	.5

J	Y	$\Theta_{Y J}$
true	true	.01
true	false	.99
false	true	.99
false	false	.01

I	J	X	$\Theta_{X IJ}$
true	true	true	.95
true	true	false	.05
true	false	true	.05
true	false	false	.95
false	true	true	.05
false	true	false	.95
false	false	true	.05
false	false	false	.95

X	Y	O	$\Theta_{O XY}$
true	true	true	.98
true	true	false	.02
true	false	true	.98
true	false	false	.02
false	true	true	.98
false	true	false	.02
false	false	true	.02
false	false	false	.98

Pruning edges and simplifying (omitting rows that are assigned with 0) gives the following.

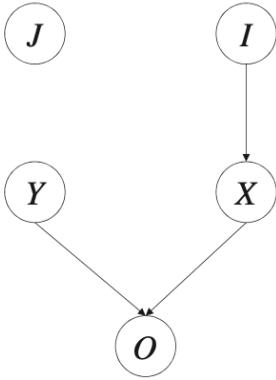
I	Θ_I^e
true	.5
false	.5

J	Θ_J^e
true	.5

Y	Θ_Y^e
true	.01
false	.99

I	X	$\Theta_{X I}^e$
true	true	.95
true	false	.05
false	true	.05
false	false	.95

X	Y	O	$\Theta_{O XY}^e$
true	true	false	.02
true	false	false	.02
false	true	false	.02
false	false	false	.98



Assuming the elimination order J, I, X, Y, O , Algorithm VE_MPE will then evaluate the following expression:

$$\begin{aligned} & \max_{J,I,X,Y,O} \Theta_I^e \Theta_J^e \Theta_Y^e \Theta_{X|I}^e \Theta_{O|XY}^e \\ &= \max_O \left(\max_Y \left(\max_X \left(\max_I \Theta_I^e \Theta_{X|I}^e \right) \Theta_{O|XY}^e \right) \Theta_Y^e \right) \left(\max_J \Theta_J^e \right). \end{aligned}$$

Continuing the example:

$$\begin{aligned} & \max_{J,I,X,Y,O} \Theta_I^e \Theta_J^e \Theta_Y^e \Theta_{X|I}^e \Theta_{O|XY}^e \\ &= \max_O \left(\max_Y \left(\max_X \left(\max_I \Theta_I^e \Theta_{X|I}^e \right) \Theta_{O|XY}^e \right) \Theta_Y^e \right) \left(\max_J \Theta_J^e \right). \end{aligned}$$

Intermediate steps:

				1/2
				$\max_j \Theta_j^e$
				$\top .5 \mid J = \text{true}$
I	X	$\Theta_I^e \Theta_{X I}^e$		
true	true	.475	\top	
true	false	.025	\top	
false	true	.025	\top	
false	false	.475	\top	
				$\max_i \Theta_i^e \Theta_{X i}^e$
			X	
			true	$.475 \mid I = \text{true}$
			false	$.475 \mid I = \text{false}$
X	Y	O		$\left(\max_I \Theta_I^e \Theta_{X i}^e \right) \Theta_{O XY}^e$
true	true	false	.0095	$I = \text{true}$
true	false	false	.0095	$I = \text{true}$
false	true	false	.0095	$I = \text{false}$
false	false	false	.4655	$I = \text{false}$
Y	O		$\max_x \left(\max_i \left(\max_I \Theta_I^e \Theta_{X i}^e \right) \Theta_{O XY}^e \right) \Theta_Y^e$	
true	false	.0095	$I = \text{true}, X = \text{true}$	
false	false	.4655	$I = \text{false}, X = \text{false}$	

				2/2
				$\left(\max_X \left(\max_I \Theta_I^e \Theta_{X i}^e \right) \Theta_{O XY}^e \right) \Theta_Y^e$
				$\top .000095 \mid I = \text{true}, X = \text{true}$
Y	O			
true	false	.000095	$I = \text{true}, X = \text{true}$	
false	false	.460845	$I = \text{false}, X = \text{false}$	
O		$\max_Y \left(\max_X \left(\max_I \Theta_I^e \Theta_{X i}^e \right) \Theta_{O XY}^e \right) \Theta_Y^e$		
false	.460845	$I = \text{false}, X = \text{false}, Y = \text{false}$		
		$\max_O \left(\max_Y \left(\max_X \left(\max_I \Theta_I^e \Theta_{X i}^e \right) \Theta_{O XY}^e \right) \Theta_Y^e \right)$		
		$.460845 \mid I = \text{false}, X = \text{false}, Y = \text{false}, O = \text{false}$		
		$\max_O \left(\max_Y \left(\max_X \left(\max_I \Theta_I^e \Theta_{X i}^e \right) \Theta_{O XY}^e \right) \Theta_Y^e \right) \left(\max_J \Theta_J^e \right)$		
		$.2304225 \mid J = \text{true}, I = \text{false}, X = \text{false}, Y = \text{false}, O = \text{false}$		

MPE Instance

Recall that in case of MAP: We first sum-out all the non-MAP variables, and then maximise-out the MAP variables |

Algorithm 30 VE MAP(\mathcal{N} , \mathbf{M} , \mathbf{e})

input:

- \mathcal{N} : Bayesian network
- \mathbf{M} : some variables in the network
- \mathbf{e} : evidence ($\mathbf{E} \cap \mathbf{M} = \emptyset$)

output: trivial factor containing the MAP probability $\text{MAP}_P(\mathbf{M}, \mathbf{e})$
main:

- 1: $\mathcal{N}' \leftarrow \text{pruneNetwork}(\mathcal{N}, \mathbf{M}, \mathbf{e})$
 - 2: $\pi \leftarrow$ a variable elimination order for \mathcal{N}' in which variables \mathbf{M} appear last
 - 3: $\mathcal{S} \leftarrow \{f^{\mathbf{e}} : f \text{ is a CPT of network } \mathcal{N}'\}$
 - 4: **for** $i = 1$ to length of order π **do**
 - 5: $f \leftarrow \prod_k f_k$, where factor f_k belongs to \mathcal{S} and mentions variable $\pi(i)$
 - 6: **if** $\pi(i) \in \mathbf{M}$ **then**
 - 7: $f_i \leftarrow \max_{\pi(i)} f$
 - 8: **else**
 - 9: $f_i \leftarrow \sum_{\pi(i)} f$
 - 10: **end if**
 - 11: replace all factors f_k in \mathcal{S} by factor f_i
 - 12: **end for**
 - 13: **return** trivial factor $\prod_{f \in \mathcal{S}} f$
-

Example: Assume the same BN. MAP variables $\mathbf{M} = \{I, J\}$

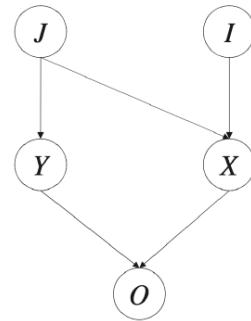
We use the order $\pi = O, Y, X, I, J$

Summing out variables O, Y , and X , we obtain the following set of factors, which represent a factored representation of the joint marginal $\Pr(I, J, O = \text{true})$ where “ $O = \text{true}$ ” is the evidence

I	J	f_1
true	true	.93248
true	false	.97088
false	true	.07712
false	false	.97088

I	f_2
true	.5
false	.5

J	f_3
true	.5
false	.5



Next, to eliminate I and J , we multiply their factors and then maximizing out.

First I ,	I	J	$f_1 f_2$			J	$\max_I f_1 f_2$
	true	true	.466240	T		true	.466240 $I = \text{true}$
	true	false	.485440	T		false	.485440 $I = \text{true}$
	false	true	.038560	T			
	false	false	.485440	T			

Next J ,	J	$(\max_I f_1 f_2) f_3$		$\max (\max_I f_1 f_2) f_3$	
	true	.233120	$I = \text{true}$		
	false	.242720	$I = \text{true}$.242720	$I = \text{true}, J = \text{false}$

MAP probability, $\Pr(I = \text{true}, J = \text{false}, O = \text{true})$



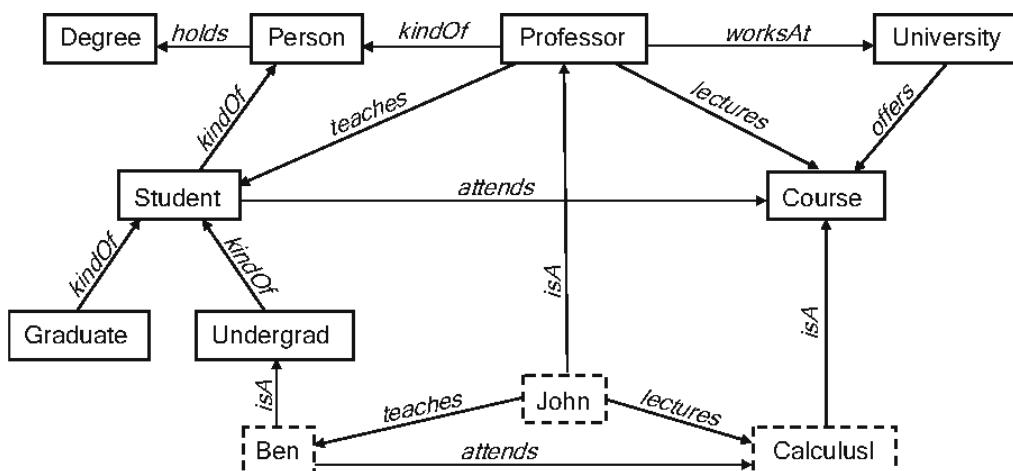
Ontologies

11. Ontologies: a way of organizing knowledge that is useful for *people*

What Is an ontology?

- “An ontology is a formal, explicit specification of a shared conceptualization.

Ontology/ Knowledge Graph



Definition(s) of an ontology

"An ontology defines the terms used to describe and represent an area of knowledge" -Jeff Heflin

"An ontology

- defines the basic terms and relations
- that form the vocabulary of a topic area
- as well as the rules for combining terms and relations
- to infer extensions to the vocabulary" -Neches ('91)

The goal of an Ontology = Knowledge Sharing and Reuse

- Communication between people
- Interoperability between software agent

Ontology Elements

Class (concept,type)

-a name and a set of properties that describe a certain set of individuals in the domain

Example: person + name, birthdate, ID-nr, ...

Instances

-the members of the sets defined by classes

Example: person45, name=Frank, birthdate=01-01-1960, ID-nr=12345

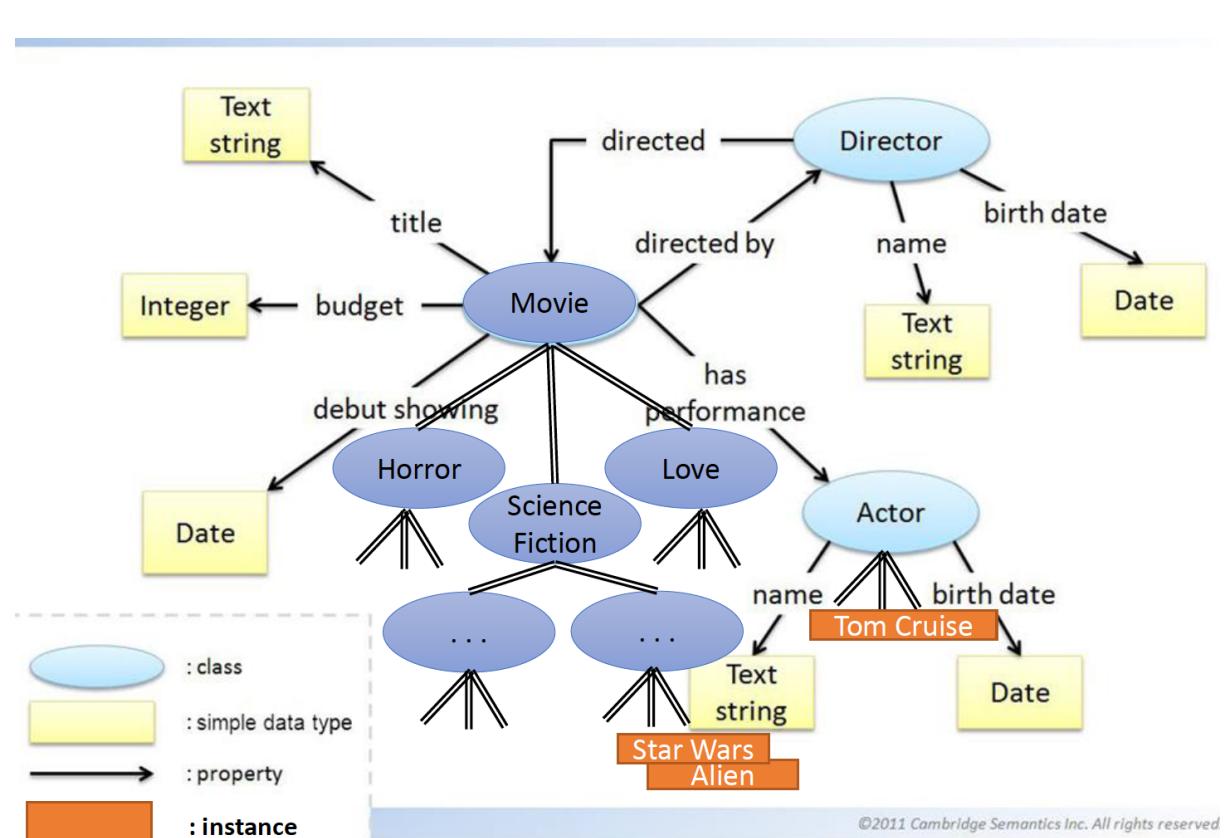
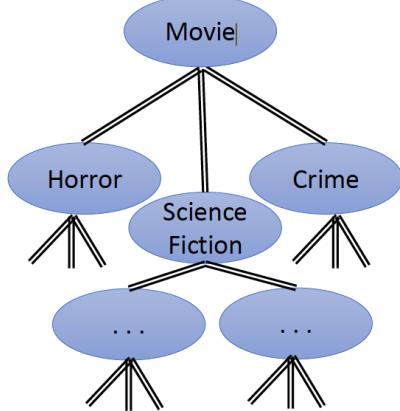
Property (relation)

-assert facts about the instances

Example: person45 is-father-of person256

An example of ontology

A Film Ontology



Real life examples

handcrafted

- music: Music ontology , MusicMoz(1073/7)
- biomedical: SNOMED(300k), Gene Ontology GO(43k terms)Emtree(45k+190kSystems biology
- ranging from lightweight
- UNSPC, Open directory (400k) to heavyweight (Cyc (300k))
- ranging from small (METAR) to large (UNSPC)

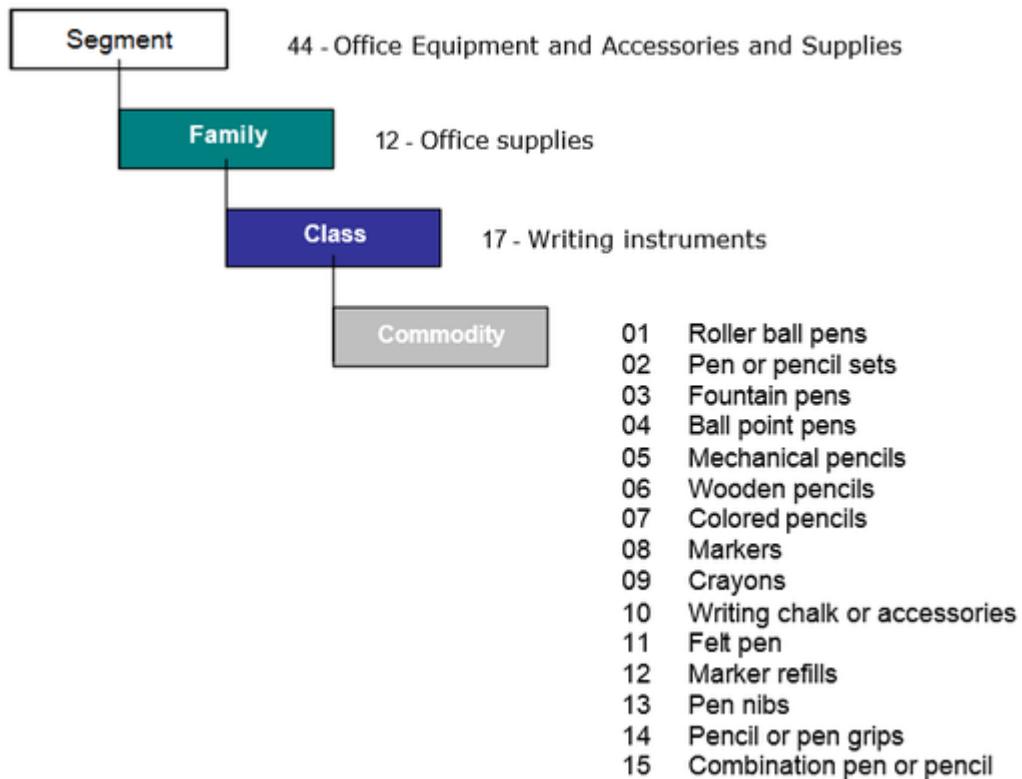
METAR is the most common format in the world for the transmission of observational weather data

METAR Class Hierarchy

- [UnitAbbreviation](#) ()
- [WeatherEvent](#) ([hasWeatherQualifier](#))
 - [CumulativeEvent](#) ()
 - [MaximumTemperature](#) ()
 - [SixHourMaximumTemp](#) ()
 - [TwentyFourHourMaximumTemp](#) ()
 - [MinimumTemperature](#) ()
 - [SixHourMinimumTemp](#) ()
 - [TwentyFourHourMinimumTemp](#) ()
 - [SixHourCumulativeEvent](#) ()
 - [SixHourMaximumTemp](#) ()
 - [SixHourMinimumTemp](#) ()
 - [ThreeHourCumulativeEvent](#) ()
 - [ThreeHourPressureTendency](#) ()
 - [TwentyFourHourCumulativeEvent](#) ()
 - [TwentyFourHourMaximumTemp](#) ()
 - [TwentyFourHourMinimumTemp](#) ()
 - [CurrentWeatherEvent](#) ()
 - [ObscurationEvent](#) ()
 - [Fog](#) ()
 - [Haze](#) ()
 - [Mist](#) ()
 - [Sand](#) ()
 - [Smoke](#) ()
 - [Spray](#) ()
 - [VolcanicAsh](#) ()
 - [WidespreadDust](#) ()

A METAR weather report is predominantly used by aircraft pilots, and by meteorologists, who use aggregated METAR information to assist in weather forecasting.

UNSPC Product and services



Level	Code	Description
Segment	44000000	Office Equipment, Accessories and Supplies
Family	44120000	Office supplies
Class	44121900	Ink and lead refills
Commodity	44121903	Pen refills

The **United Nations Standard Products and Services Code(UNSPSC)** is a taxonomy of products and services for use in eCommerce. It is a four-level hierarchy coded as an eight-digit number, with an optional fifth level adding two more digits.

By classifying their products & services, businesses can assist their customers with:

Finding and Purchasing

Product discovery

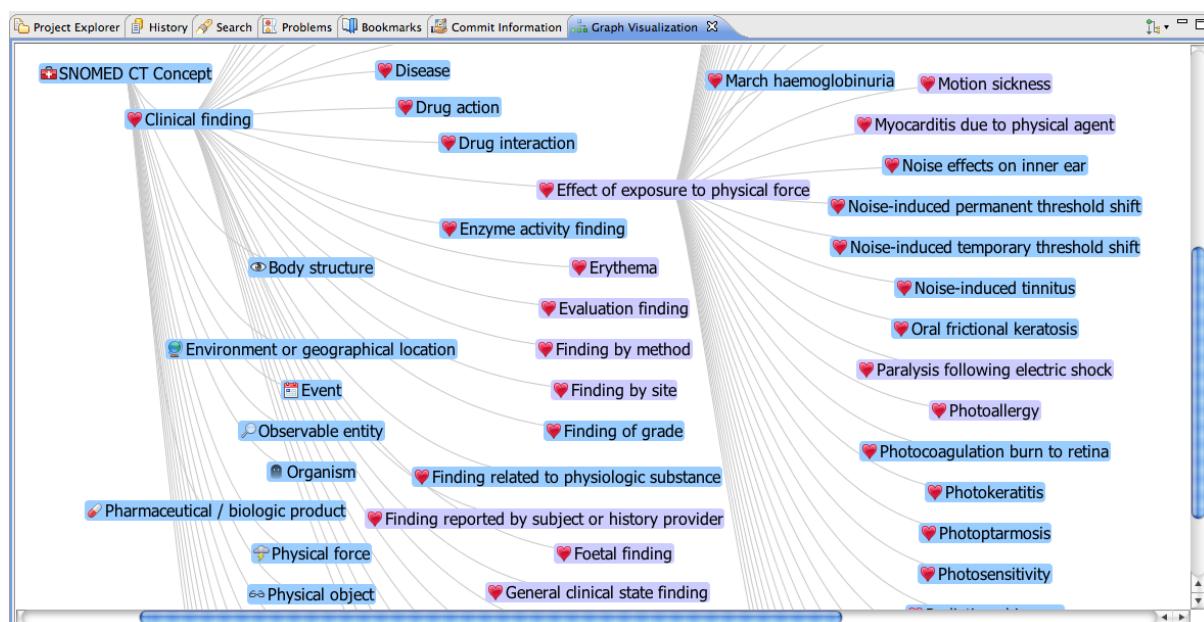
Facilitates expenditure analysis.

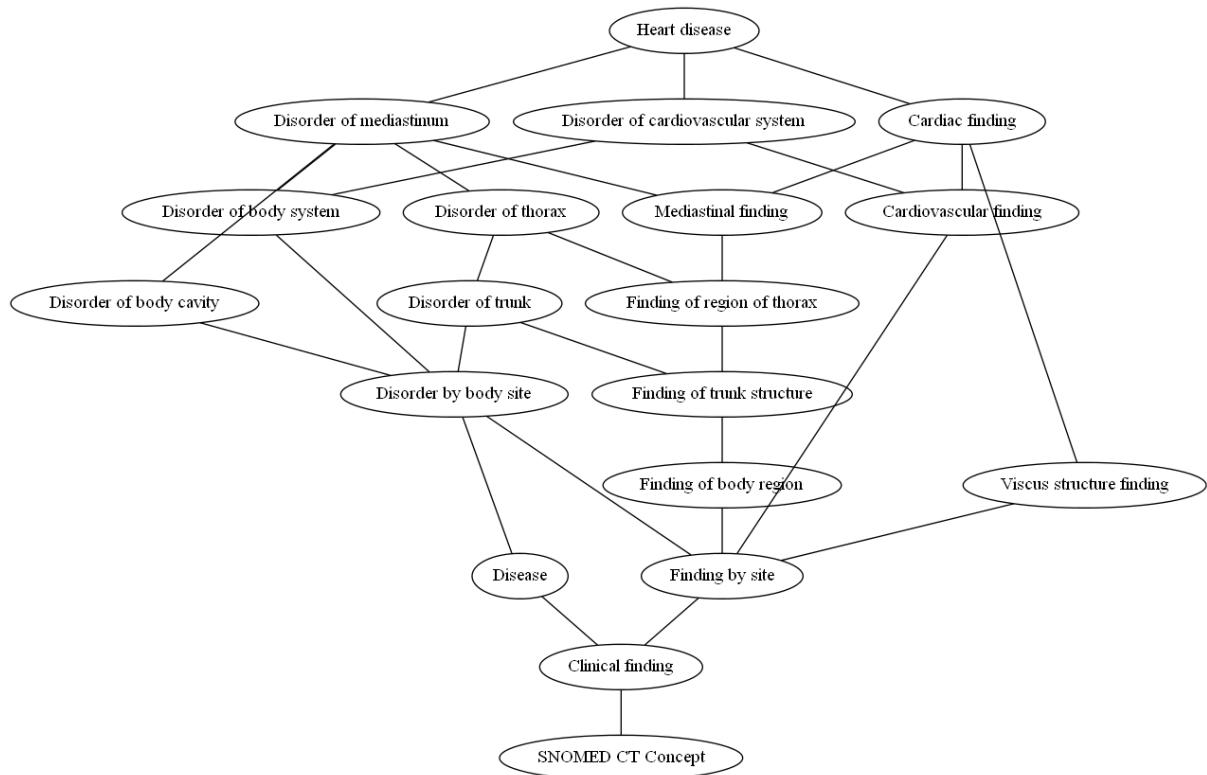
Control and uniformity across the company

UNSPC Product and services



SNOMED: used in healthcare





Biomedical ontologies (a few..)

- Mesh
- Medical Subject Headings, National Library of Medicine
- 22.000 descriptions
- EMTREE
- Commercial Elsevier, Drugs and diseases
- 45.000 terms, 190.000 synonyms
- UMLS
- Integrates 100 different vocabularies
- SNOMED
- 300.000 concepts, College of American Pathologists
- Gene Ontology
- 15.000 terms in molecular biology
- NCBI Cancer Ontology:
- 17,000 classes (about 1M definitions),

Gene Ontology

The Gene Ontology (GO) is a major bioinformatics initiative to unify the representation of gene and gene product attributes across all species

Three domains:

cellular component
molecular function
biological process

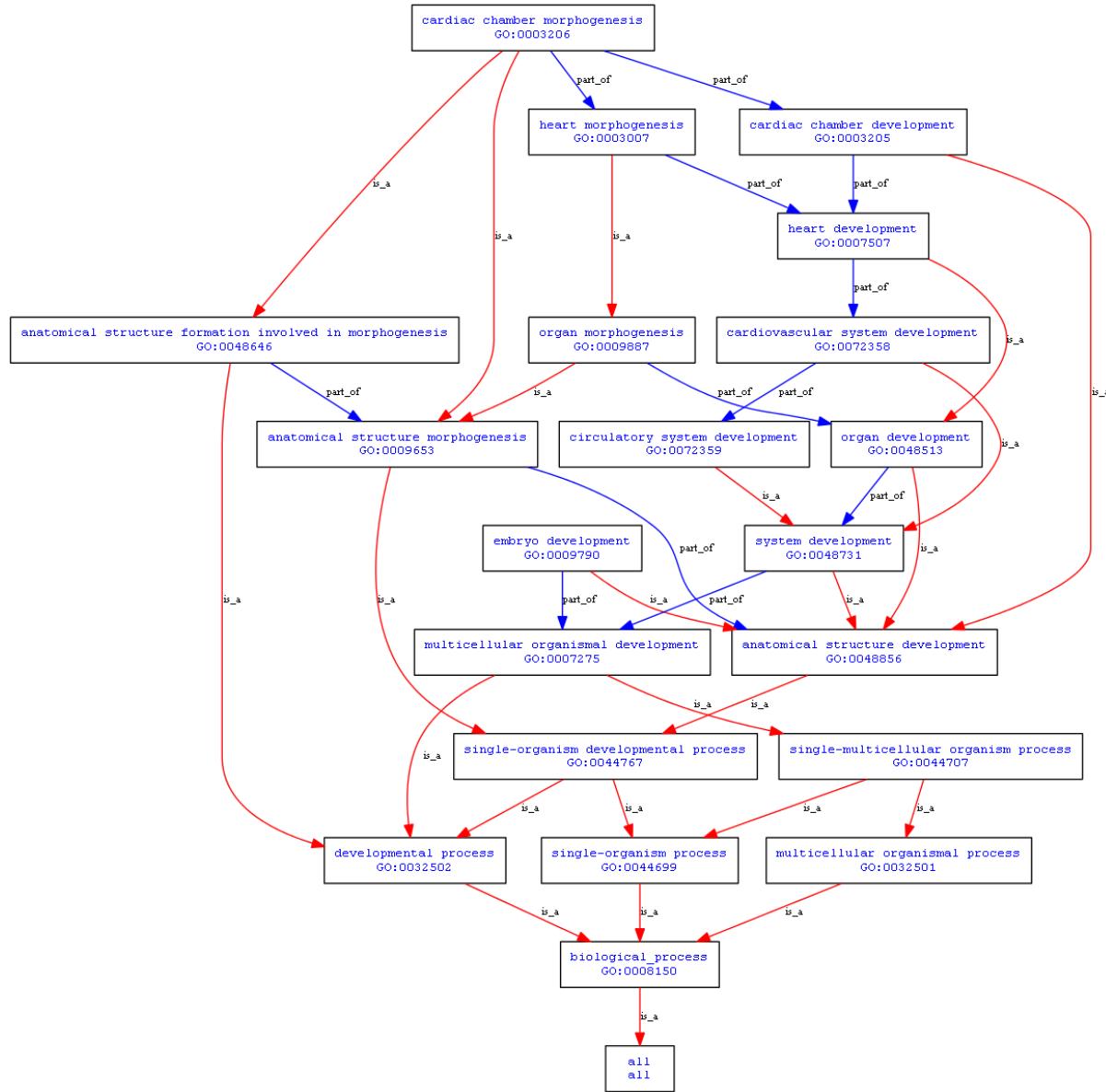
Gene product: Actin... UniProtKB:P68032

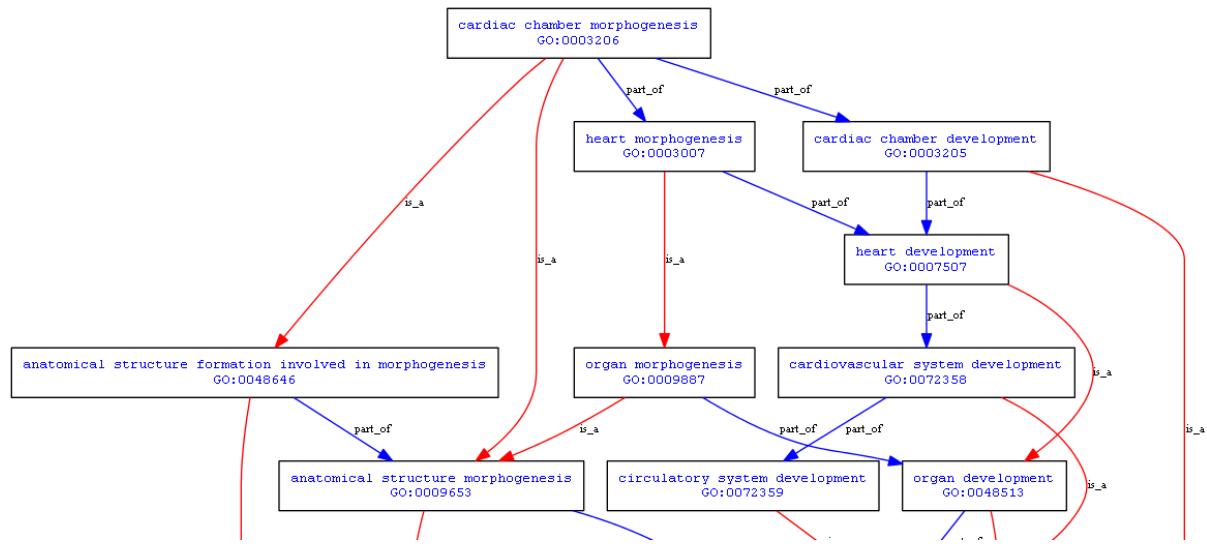
GO term: heartcontraction; GO:0060047 (biologicalprocess)

Evidencecode: InferredfromMutant Phenotype(IMP)

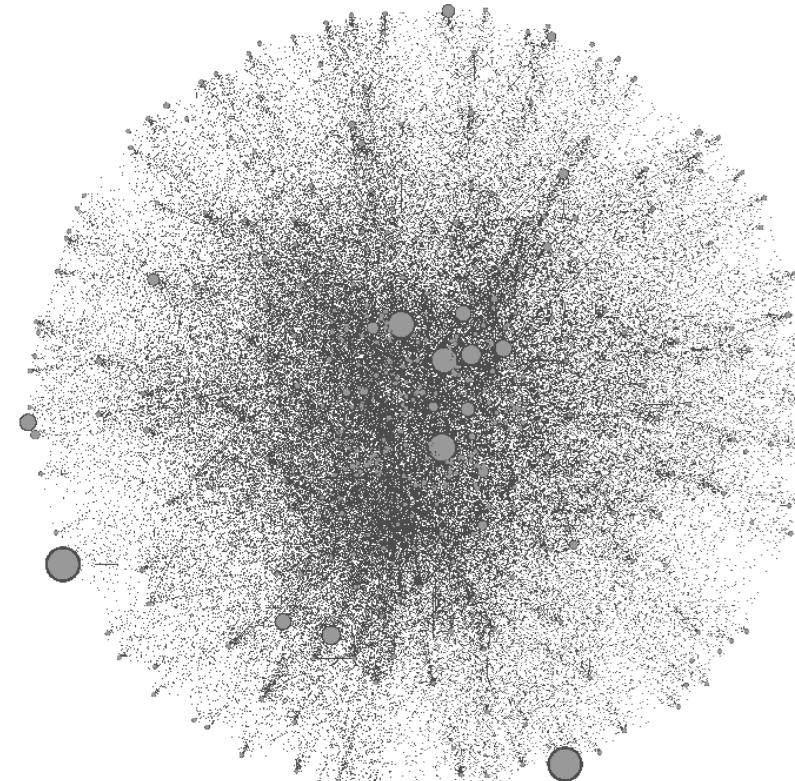
Reference: PMID 17611253

Assignedby: UniProtKB, June6, 2008





Bio-medical ontologies in Bio-portal at Stanford



A logic for Ontologies

A task for Logic Engineering (KR)

- What logic do we need for UMS PC?

What logic do we need for Mesh?

- What about GO?

F–Psychiatry and Psychology F01–behavior and behavior mechanisms

F02–psychological phenomena and processes

F03–mental disorders

What's inside an ontology?

Classes + class-hierarchy

instances

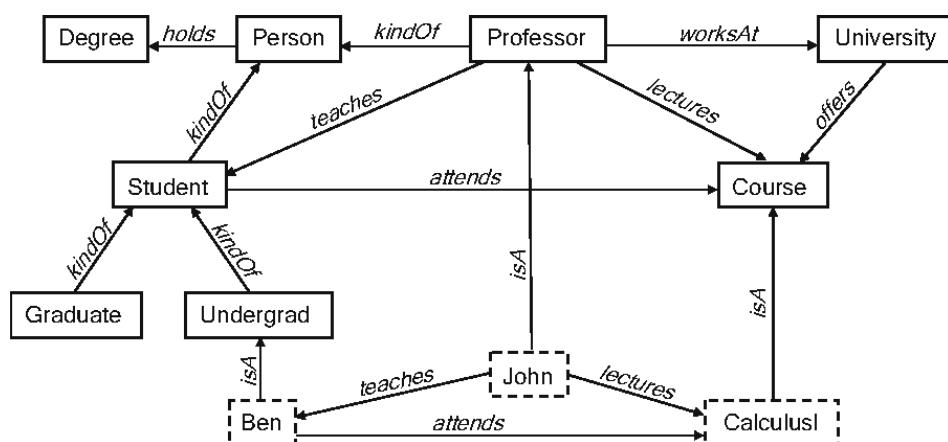
slots/values

inheritance

- restrictions on slots (type, cardinality)
- properties of slots (symm., trans., ...)
- relations between classes (disjoint, covers)
- reasoning tasks: classification, subsumption

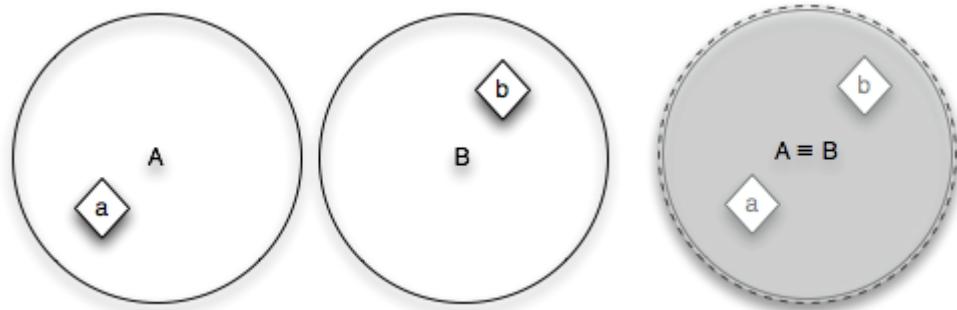
Operators for ontologies?

Ontology/ Knowledge Graph

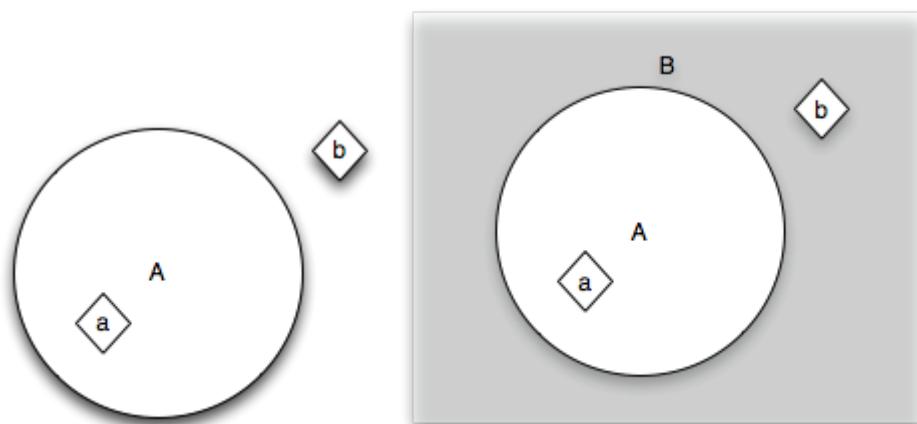


Class Axioms -Equivalence and Complement

Equivalent Classes contain the same individuals, and have the same definition

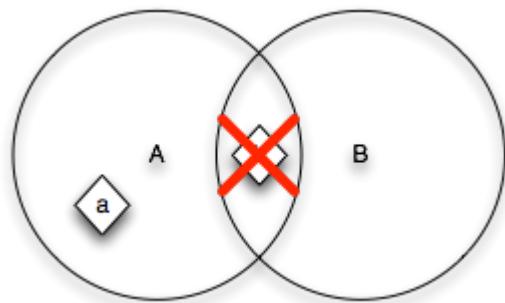


The **complement** of a class contains all individuals that are not in the class

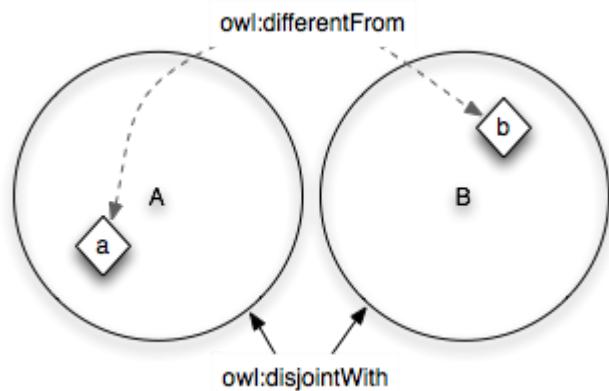


Class Axioms -Disjointness

Disjoint Classes do not contain the same individuals

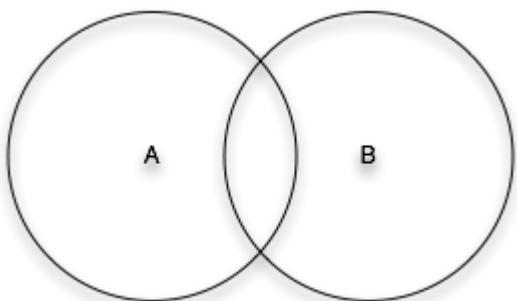
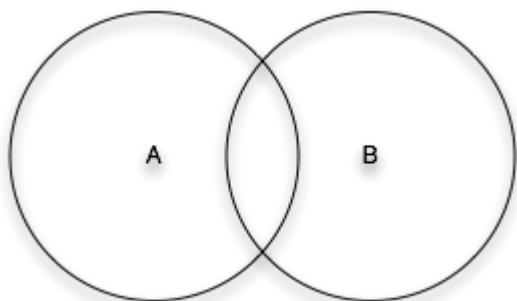


The **complement and disjointness** allow us to infer that individuals are different.

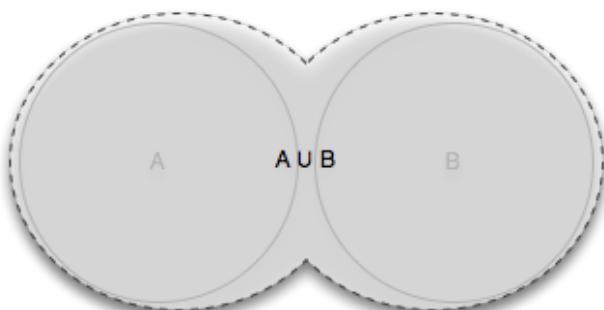


Class Axioms -Union and Disjoint Union

The Union Contains all individuals that belong to the classes of the union

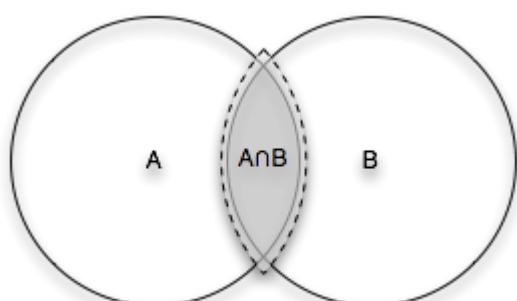


A **disjoint union** is a union of mutually **disjoint classes**

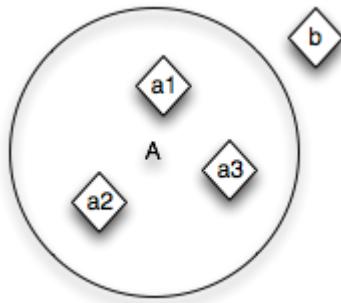


Class Axioms -Intersection and Enumeration

The Intersection contains individuals that each belong to **both classes**



You can **enumerate all** members of a class



Class Restrictions-Some/all(existential, universal)

All members of a class have at least **some** value from the specified class
"**NobelPrizeWinners** is the class of **Persons** who have **won a NobelPrize**"

All members of a class have only values from one specified class

A **NobelPrize** is the class of **Prizes** that only have been won by **NobelPrizeWinners**

That is not the same as: Every **NobelPrizeWinner** only **won NobelPrizes**

Necessary and Sufficient Conditions

All members of a class must also be members of another class (it is **necessary for a NobelPrizeWinner** to have been nominated for Nobel Price).

"**NobelPrizeWinners** is exactly the class of all **Persons** who have **won a NobelPrize**"

The class is precisely specified by some other class (it is a **necessary and sufficient condition**).

"**NobelPrizeWinners** is the subset of all **Persons** who were **nominated for a NobelPrize**"

Property Types -Symmetric Property

Used to specify that a property always holds in both directions

```
if      type(p,SymmetricProperty) and p(x,y)
then    p(y,x)
```

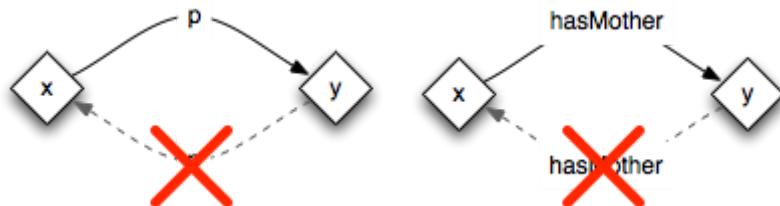


Property Types -Asymmetric Property

Used to specify that a property **never** holds in both directions

if type(p,AsymmetricProperty) **and** p(x,y)

then p(y,x) **is inconsistent**

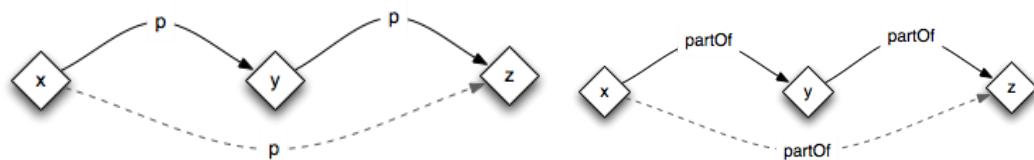


Property Types -Transitive Property

Used to specify that a property **propagates** over itself

if type(p,TransitiveProperty) **and** p(x,y) **and** p(y,z)

then p(x,z)

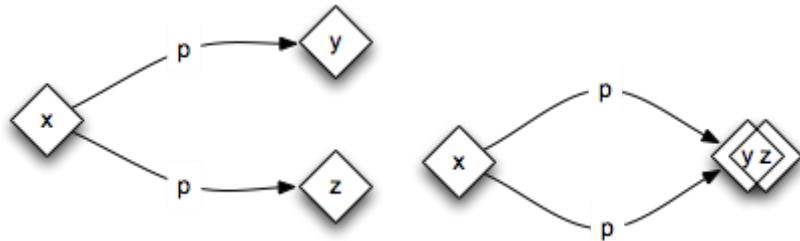


Property Types -Functional Property

Used to specify that a property has **only one** value for any particular instance

if type(p,FunctionalProperty) **andp(x,y)** **andp(x,z)**

then $y = z$

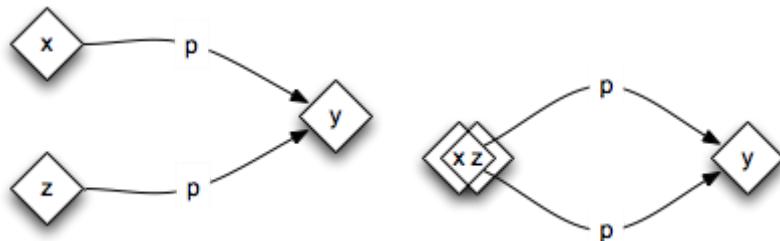


Property Types -Inverse Functional Property

Used to specify that a value for the property **uniquely identifies** an instance

if type(p,InverseFunctionalProperty) **andp(x,y)** **andp(z,y)**

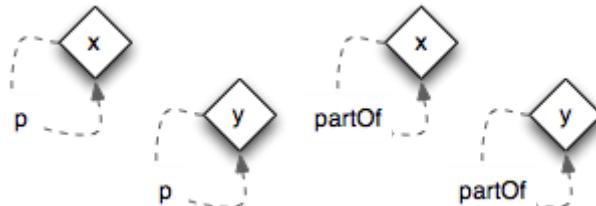
then $x = z$



Property Types -Reflexive Property

Used to specify that **every individual is always related to itself** by that property

if type(p,ReflexiveProperty)
then $\forall x p(x,x)$

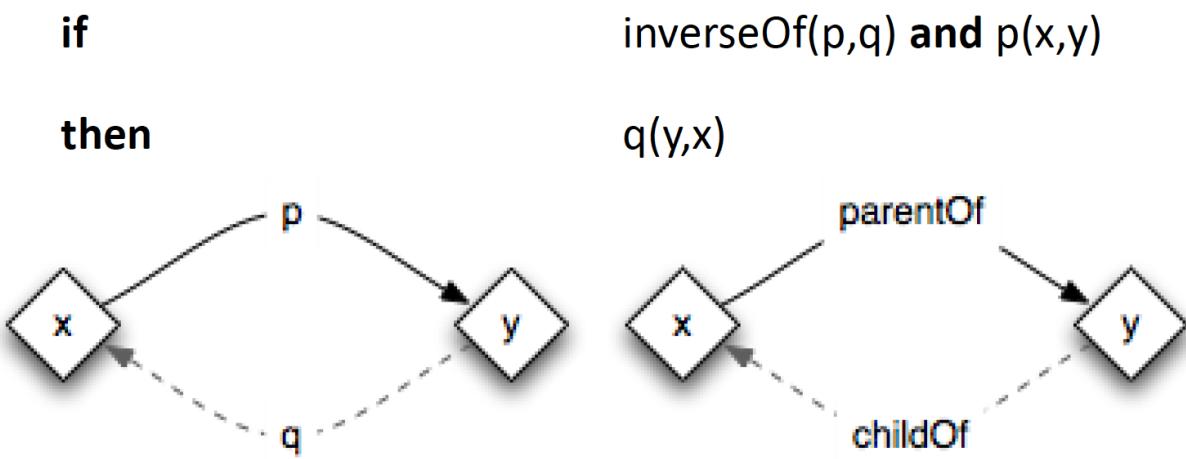


if type(p,IrreflexiveProperty) **and** $p(x,x)$
then $p(x,x)$ **is inconsistent**



Property Axioms -Inverse Property

Used to specify that one property is **always** the inverse of another property



Summary & what comes now?

- We have seen
- Ontologies: a **hierarchical representation** of types (“**classesinstances** (“individuals”) that belong to classes, and **properties between** the individuals.
- A number of realistic examples of ontologies in use
- These ontologies are easy to read by people
- **Now:** But can we also make a **logics for** ontologies, so that computers can reason automatically?

