

---

# 13.

## Explanation and Diagnosis

### Abductive reasoning

---

So far: reasoning has been primarily *deductive*:

- given KB, is  $\alpha$  an implicit belief?
- given KB, for what  $x$  is  $\alpha[x]$  an implicit belief?

Even default / probabilistic reasoning has a similar form

Now consider a new type of question:

Given KB, and an  $\alpha$  that I do *not* believe,

what would be sufficient to make me believe that  $\alpha$  was true?

or what else would I have to believe for  $\alpha$  to become an implicit belief?

or what would *explain*  $\alpha$  being true?

Deduction: given  $(p \supset q)$ , from  $p$ , deduce  $q$

Abduction: given  $(p \supset q)$ , from  $q$ , abduce  $p$

$p$  is sufficient for  $q$  or one way for  $q$  to be true is for  $p$  to be true

Also induction: given  $p(t_1), q(t_1), \dots, p(t_n), q(t_n)$ , induce  $\forall x (p(x) \supset q(x))$

Can be used for causal reasoning: (*cause*  $\supset$  *effect*)

# Diagnosis

One simple version of diagnosis uses abductive reasoning

KB has facts about symptoms and diseases

including:  $(Disease \wedge Hedges \supset Symptoms)$

Goal: find disease(s) that best explain observed symptoms

Observe: we typically do not have knowledge of the form

$(Symptom \wedge \dots \supset Disease)$

so reasoning is not deductive

Example:

(tennis-elbow $\supset$ sore-elbow)
(tennis-elbow $\supset$ tennis-player)
(arthritis $\wedge$ untreated $\supset$ sore-joints)
(sore-joints $\supset$ sore-elbow $\wedge$ sore-hip)

Explain: sore-elbow

Want: tennis-elbow,  
(arthritis  $\wedge$  untreated),  
...

Non-uniqueness: multiple equally good explanations

+ logical equivalences:  $(untreated \wedge \neg \neg arthritis)$

## Adequacy criteria

Given KB, and  $\beta$  to be explained, we want an  $\alpha$  such that

1.  $\alpha$  is sufficient to account for  $\beta$

$KB \cup \{\alpha\} \models \beta$  or  $KB \models (\alpha \supset \beta)$

2.  $\alpha$  is not ruled out by KB

$KB \cup \{\alpha\}$  is consistent or  $KB \not\models \neg \alpha$

otherwise  $(p \wedge \neg p)$  would count as an explanation

3.  $\alpha$  is as simple as possible

parsimonious : as few *terms* as possible  
explanations should not unnecessarily  
strong or unnecessarily weak

e.g.  $KB = \{(p \supset q), \neg r\}$  and  $\beta = q$

$\alpha = (p \wedge s \wedge \neg t)$  is too strong

$\alpha = (p \vee r)$  is too weak

4.  $\alpha$  is in the appropriate vocabulary

atomic sentences of  $\alpha$  should be drawn  
from **H**, possible hypotheses in terms of  
which explanations are to be phrased

e.g. diseases, original causes

e.g. sore-elbow explains sore-elbow  
trivial explanation

sore-joints explains sore-elbow  
may or may not be suitable

Call such  $\alpha$  an explanation of  $\beta$  wrt KB

## Some simplifications

---

From criteria of previous slide, we can simplify explanations in the propositional case, as follows:

- To explain an arbitrary wff  $\beta$ , it is sufficient to choose a new letter  $p$ , add  $(p \equiv \beta)$  to KB, and then explain  $p$ .

$$\text{KB} \models (E \supset \beta) \quad \text{iff} \quad \text{KB} \cup \{(p \equiv \beta)\} \models (E \supset p)$$

- Any explanation will be (equivalent to) a conjunction of literals (that is, the negation of a clause)

Why? If  $\alpha$  is a purported explanation, and  $\text{DNF}[\alpha] = (d_1 \vee d_2 \vee \dots \vee d_n)$  then each  $d_i$  is also an explanation that is no less simple than  $\alpha$

A simplest explanation is then the negation of a clause with a *minimal* set of literals

So: to explain a literal  $\rho$ , it will be sufficient to find the minimal clauses  $C$  (in the desired vocabulary) such that

- |   |            |
|---|------------|
| 1. $\text{KB} \models (\neg C \supset \rho)$ or $\text{KB} \models (C \cup \{\rho\})$ | sufficient |
| 2. $\text{KB} \not\models C$  | consistent |

## Prime implicates

---

A clause  $C$  is a prime implicate of a KB iff

- |   |  |
|---|--|
| 1. $\text{KB} \models C$                            | Note: For any clause $C$ , if $\text{KB} \models C$ , then some subset of $C$ is a prime implicate |
| 2. For no $C^* \subset C$ , $\text{KB} \models C^*$ |  |

Example:  $\text{KB} = \{(p \wedge q \wedge r \supset g), (\neg p \wedge q \supset g), (\neg q \wedge r \supset g)\}$

Prime implicates:

$(p \vee \neg q \vee g),$   
 $(\neg r \vee g),$  and  
 $(p \vee \neg p), (g \vee \neg g), \dots$

Note: tautology  $(a \vee \neg a)$  is always a prime implicate unless  $\text{KB} \models a$  or  $\text{KB} \models \neg a$

For explanations:

- want minimal  $C$  such that  $\text{KB} \models (C \cup \{\rho\})$  and  $\text{KB} \not\models C$
- so: find prime implicates  $C$  such that  $\rho \in C$ ; then  $\neg(C - \rho)$  must be an explanation for  $\rho$

Example: explanations for  $g$  in example above

- 3 prime implicates contain  $g$ , so get 3 explanations:  $(\neg p \wedge q), r,$  and  $g$

# Computing explanations

Given KB, to compute explanations of literal  $\rho$  in vocabulary  $\mathbf{H}$ :

calculate the set  $\{\neg(C - \rho) \mid C \text{ is a prime implicate and } \rho \in C\}$   
prime implicates containing  $\rho$

But how to compute prime implicates?

Can prove: Resolution is complete for non-tautologous prime implicates

$KB \models C$  iff  $KB \rightarrow C$       completeness for  $\square$  is a special case!

So: assuming KB is in CNF, generate *all* resolvents in language  $\mathbf{H}$ , and retain those containing  $\rho$  that are minimal

Could pre-compute all prime implicates, but there may be *exponentially* many, even for a Horn KB

Example: atoms:  $p_i, q_i, E_i, O_i, 0 \leq i < n + E_n, O_n$

wffs:  $E_i \wedge p_i \supset O_{i+1}, E_i \wedge q_i \supset E_{i+1},$   
 $O_i \wedge p_i \supset E_{i+1}, O_i \wedge q_i \supset O_{i+1},$   
 $E_0, \neg O_0$

explain:  $E_n$

## Circuit example

### Components

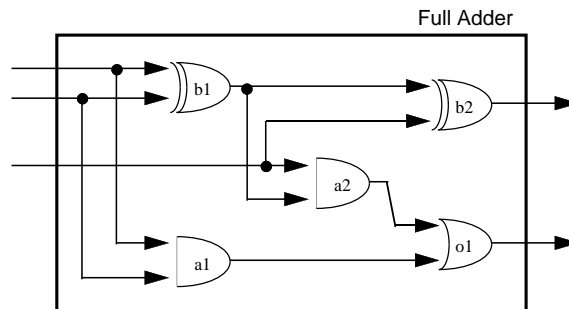
Gate( $x$ )  $\equiv$  Andgate( $x$ )  $\vee$  Orgate( $x$ )  $\vee$  Xorgate( $x$ )

Andgate(a1), Andgate(a2),

Orgate(o1),

Xorgate(b1), Xorgate(b2)

Fulladder(f) the whole circuit



### Connectivity

$in1(b1) = in1(f), in2(b1) = in2(f)$

$in1(b2) = out(b1), in2(b2) = in3(f)$

$in1(a1) = in1(f), in2(a1) = in2(f)$

$in1(a2) = in3(f), in2(a2) = out(b1)$

$in1(o1) = out(a2), in2(o1) = out(a1)$

$out1(f) = out(b2), out2(f) = out(o1)$

# Circuit behaviour

---

## Truth tables for logical gates

$\text{and}(0,0) = 0, \text{ and}(0,1) = 0, \dots \quad \text{or}(0,0) = 0, \text{ or}(0,1) = 1, \dots$   
 $\text{xor}(0,0) = 0, \text{ xor}(0,1) = 1, \dots$

## Normal behaviour

$\text{Andgate}(x) \wedge \neg \text{Ab}(x) \supset \text{out}(x) = \text{and}(\text{in1}(x), \text{in2}(x))$   
 $\text{Orgate}(x) \wedge \neg \text{Ab}(x) \supset \text{out}(x) = \text{or}(\text{in1}(x), \text{in2}(x))$   
 $\text{Xorgate}(x) \wedge \neg \text{Ab}(x) \supset \text{out}(x) = \text{xor}(\text{in1}(x), \text{in2}(x))$

## Abnormal behaviour: fault models

### Examples

$[\text{Orgate}(x) \vee \text{Xorgate}(x)] \wedge \text{Ab}(x) \supset \text{out}(x) = \text{in2}(x)$  (short circuit)

### Other possibilities ...

- some abnormal behaviours may be inexplicable
- some may be compatible with normal behaviour on certain inputs

# Abductive diagnosis

---

## Given KB as above + input settings

e.g.  $\text{KB} \cup \{\text{in1}(f) = 1, \text{in2}(f) = 0, \text{in3}(f) = 1\}$

## we want to explain observations at outputs

e.g.  $(\text{out1}(f) = 1 \wedge \text{out2}(f) = 0)$

## in the language of Ab

We want conjunction of Ab literals  $\alpha$  such that  
 $\text{KB} \cup \text{Settings} \cup \{\alpha\} \models \text{Observations}$

## Compute by “propositionalizing”:

For the above,  $x$  ranges over 5 components and  $u, v$  range over 0 and 1.

Easiest to do by preparing a table ranging over all Ab literals, and seeing which conjunctions entail the observations.

## Table for abductive diagnosis

	Ab(b1)	Ab(b2)	Ab(a1)	Ab(a2)	Ab(o1)	Entails observation?
1.	Y	Y	Y	Y	Y	N
2.	Y	Y	Y	Y	N	N
3.	Y	Y	Y	N	Y	N
4.	Y	Y	Y	N	N	N
5.	Y	Y	N	Y	Y	Y
6.	Y	Y	N	Y	N	N
7.	Y	Y	N	N	Y	Y
8.	Y	Y	N	N	N	Y
...						
25.	N	N	Y	Y	Y	N
26.	N	N	Y	Y	N	N
27.	N	N	Y	N	Y	N
28.	N	N	Y	N	N	N
29.	N	N	N	Y	Y	N
30.	N	N	N	Y	N	N
31.	N	N	N	N	Y	N
32.	N	N	N	N	N	N

## Example diagnosis

Using the table, we look for minimal sets of literals.

For example, from line (5), we have that

$$Ab(b1) \wedge Ab(b2) \wedge \neg Ab(a1) \wedge Ab(a2) \wedge Ab(o1)$$

entails the observations. However, lines (5), (7), (13) and (15) together lead us to a smaller set of literals (the first explanation below).

The explanations are

1.  $Ab(b1) \wedge \neg Ab(a1) \wedge Ab(o1)$
2.  $Ab(b1) \wedge \neg Ab(a1) \wedge \neg Ab(a2)$
3.  $Ab(b2) \wedge \neg Ab(a1) \wedge Ab(o1)$

Note: not all components are mentioned since for these settings, get the same observations whether or not they are working

but for this fault model only

Can narrow down diagnosis by looking at a number of different settings  
differential diagnosis

# Diagnosis revisited

---

## Abductive definition has limitations

- often only care about what is not working
- may not be able to characterize all possible failure modes
- want to prefer diagnoses that claim as few broken components as possible

## Consistency-based diagnosis:

Assume KB uses the predicate  $Ab$  as before, but perhaps only characterizes the normal behaviour

e.g.  $Andgate(x) \wedge \neg Ab(x) \supset out(x) = and(in1(x), in2(x))$

Want a minimal set of components  $D$ , such that

$\{Ab(c) \mid c \in D\} \cup \{\neg Ab(c) \mid c \notin D\}$

can use table as before  
with last column changed  
to "consistency"

is consistent with  $KB \cup Settings \cup Observations$

In previous example, get 3 diagnoses:  $\{b1\}$ ,  $\{b2, a2\}$  and  $\{b2, o1\}$

Note: more complex to handle non-minimal diagnoses

# Some complications

---

## 1. negative evidence

- allow for missing observations  
e.g. ensure that  $KB \cup \{\alpha\} \not\models fever$

## 2. variables and quantification

- same definition, modulo "simplicity", (but how to use Resolution?)
- useful to handle open wffs also  
 $KB \cup \{x = 3\} \models P(x)$  handles WH-questions

## 3. probabilities

- not all simplest explanations are equally likely
- also: replace  $(Disease \wedge \dots \supset Symptom)$  by a probabilistic version

## 4. defaults

- instead of requiring  $KB \cup \{\alpha\} \models \beta$ , would prefer that given  $\alpha$ , it is reasonable to believe  $\beta$   
e.g. being a bird explains being able to fly

# Other applications

---

## 1. object recognition

what scene would account for image elements observed?

what objects would account for collection of properties discovered?

## 2. plan recognition

what high-level goals of an agent would account for the actions observed?

## 3. hypothetical reasoning

instead of asking: what would I have to be told to believe  $\beta$ ?

ask instead: what would I learn if I was told that  $\alpha$ ?

Dual of explanation: want  $\beta$  such that

Solution: you learn  $\beta$  on being told  $\alpha$

iff

$\neg\beta$  is an explanation for  $\neg\alpha$

can use the abduction procedure

$KB \cup \{\alpha\} \models \beta$

$KB \not\models \beta$

simplicity, parsimony

using correct vocabulary