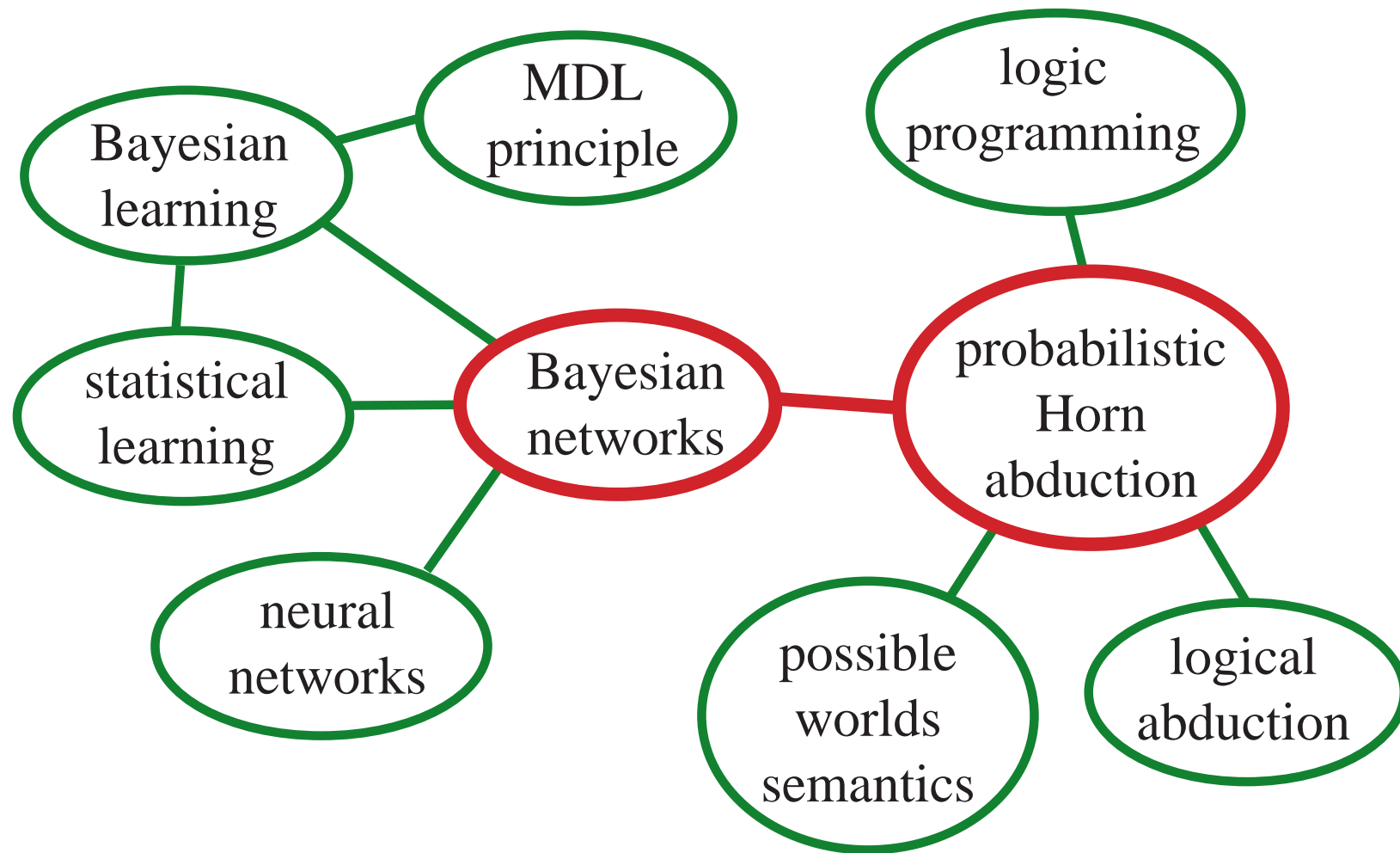


Learning, Bayesian Probability, Graphical Models, and Abduction

David Poole

University of British Columbia

Induction and abduction

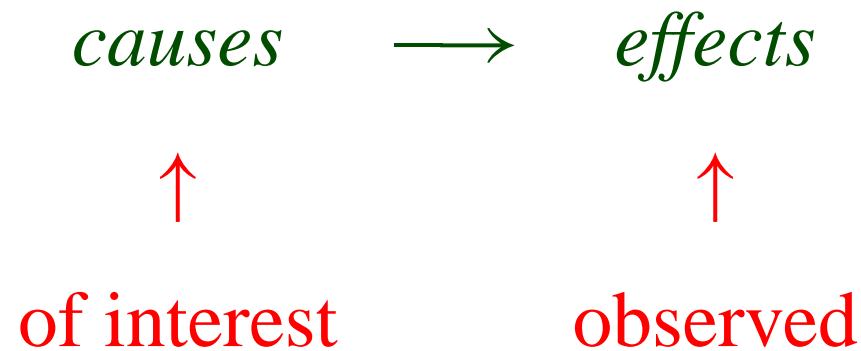


Overview

- Causal and evidential modelling & reasoning
- Bayesian networks, Bayesian conditioning & abduction
- Noise, overfitting, and Bayesian learning

Causal & Evidential Modelling

Causal modelling:



vision: *scene* \longrightarrow *image*

diagnosis: *disease* \longrightarrow *symptoms*

learning: *model* \longrightarrow *data*

Evidential modelling:

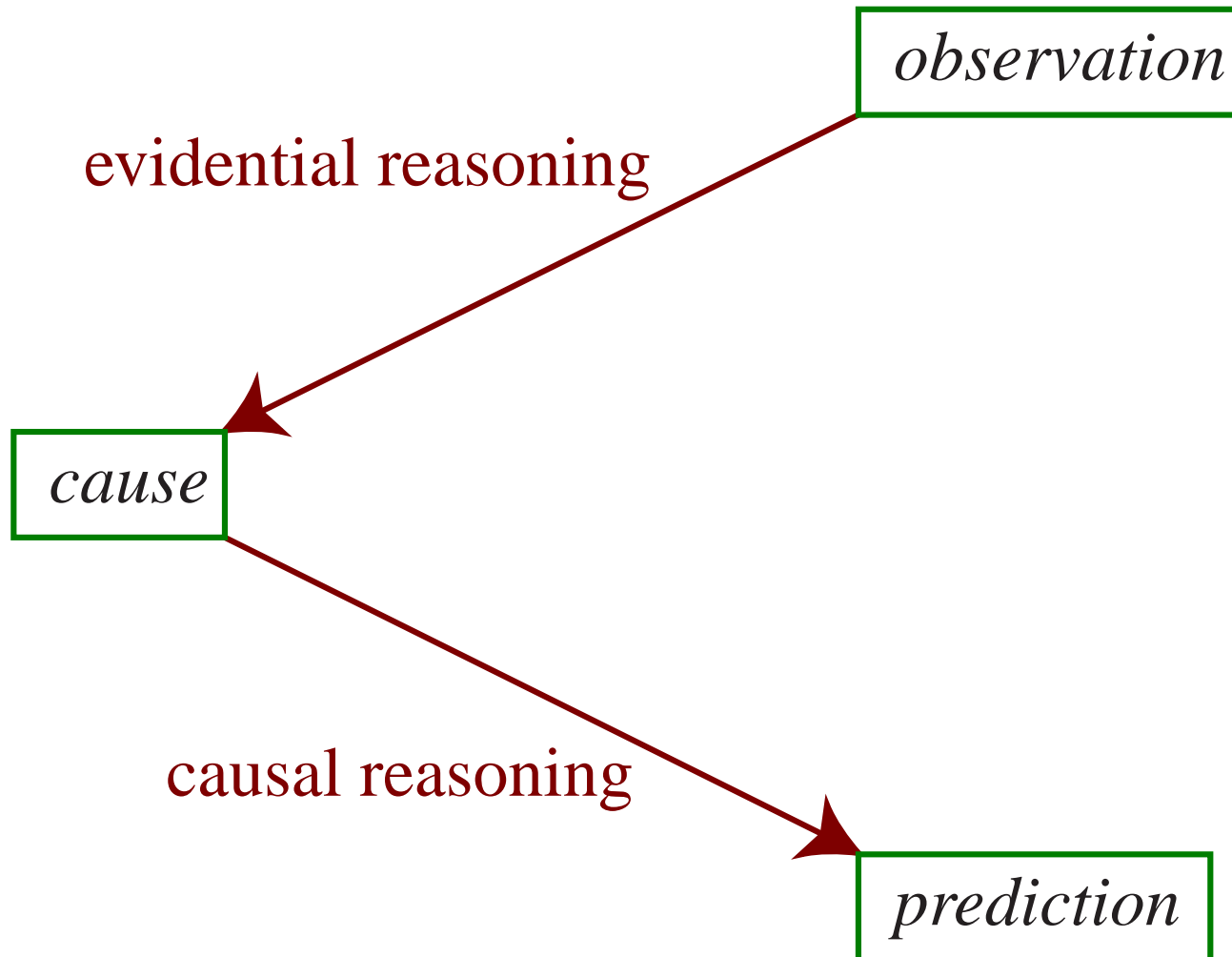
effects \longrightarrow *causes*

vision: *image* \longrightarrow *scene*

diagnosis: *symptoms* \longrightarrow *diseases*

learning: *data* \longrightarrow *model*

Causal & Evidential Reasoning



Reasoning & Modelling Strategies

How do we do causal and evidential reasoning, given modelling strategies?

- Evidential modelling & only evidential reasoning (Mycin, Neural Networks).
- Model evidentially + causally (problem: consistency, redundancy, knowledge acquisition)
- Model causally; use different reasoning strategies for causal & evidential reasoning. (deduction + abduction or Bayes' theorem)

Bayes' Rule

[de Moivre 1718, Bayes 1763, Laplace 1774]

$$P(h|e) = \frac{P(e|h)P(h)}{P(e)}$$

Proof:

$$\begin{aligned} P(h \wedge e) &= P(e|h)P(h) \\ &= P(h|e)P(e) \end{aligned}$$

Lesson #1

You should know the difference between

- evidential & causal modelling
- evidential & causal reasoning

There seems to be a relationship between Bayes' theorem and abduction — used for the same task.

Overview

- Causal and evidential modelling & reasoning
- Bayesian networks, Bayesian conditioning & abduction
- Noise, overfitting, and Bayesian learning

Bayesian Networks

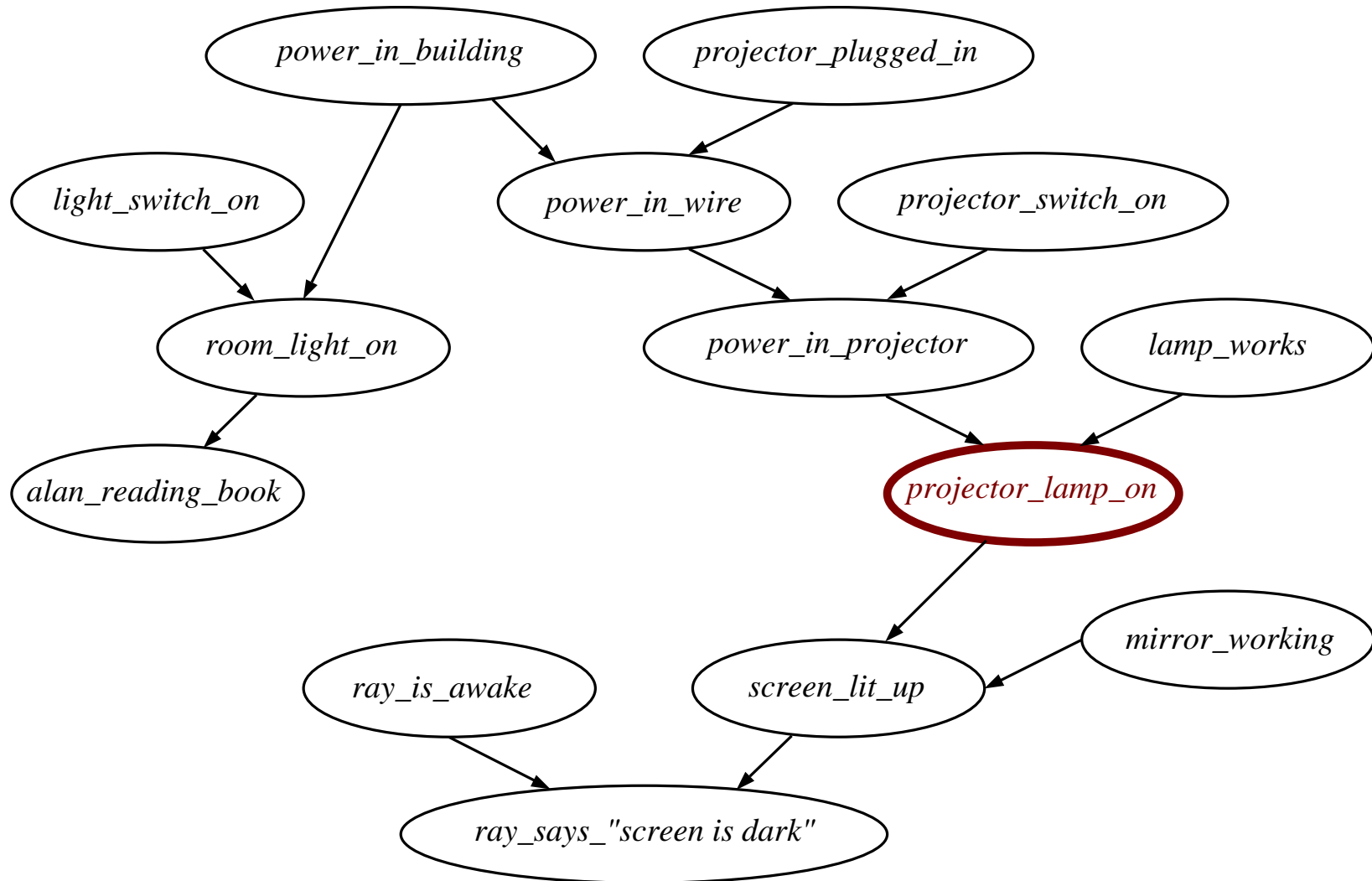
- Graphical representation of independence.
- DAGs with nodes representing random variables.
- Embed independence assumption:

If b_1, \dots, b_n are the parents of a then

$$P(a|b_1, \dots, b_n, V) = P(a|b_1, \dots, b_n)$$

if V is not a descendant of a .

Bayesian Network for Overhead Projector



Bayesian networks as logic programs

projector_lamp_on ←

power_in_projector ∧

lamp_works ∧

projector_working_ok. ← possible hypothesis
with associated probability

projector_lamp_on ←

power_in_projector ∧

\sim *lamp_works* ∧

working_with_faulty_lamp.

Probabilities of hypotheses

$P(\text{projector_working_ok})$

$= P(\text{projector_lamp_on} \mid$

$\text{power_in_projector} \wedge \text{lamp_works})$

— provided as part of Bayesian network

$P(\sim\text{projector_working_ok})$

$= 1 - P(\text{projector_working_ok})$

What do these logic programs mean?

- Possible world for each assignment of truth value to a possible hypothesis:

$\{projector_working_ok, working_with_faulty_lamp\}$

$\{projector_working_ok, \sim working_with_faulty_lamp\}$

$\{\sim projector_working_ok, working_with_faulty_lamp\}$

$\{\sim projector_working_ok, \sim working_with_faulty_lamp\}$

- Probability of a possible world is the product of the probabilities of the associated hypotheses.
- Logic program specifies what else is true in each possible world.

Probabilistic logic programs & abduction

Semantics is abductive in nature

— set of explanations of a proposition characterizes the possible worlds in which it is true.

(assume possible hypotheses and their negations).

$$P(g) = \sum_{E \text{ is an explanation of } g} P(E)$$

$$P(E) = \prod_{h \in E} P(h)$$

↑ given with logic program

Conditional Probabilities

$$P(g|e) = \frac{P(g \wedge e)}{P(e)} \quad \leftarrow \text{explain } g \wedge e$$
$$\quad \quad \quad \quad \quad \quad \quad \quad \leftarrow \text{explain } e$$

Given evidence e , explain e , then try to explain g from these explanations.

Lessons #2

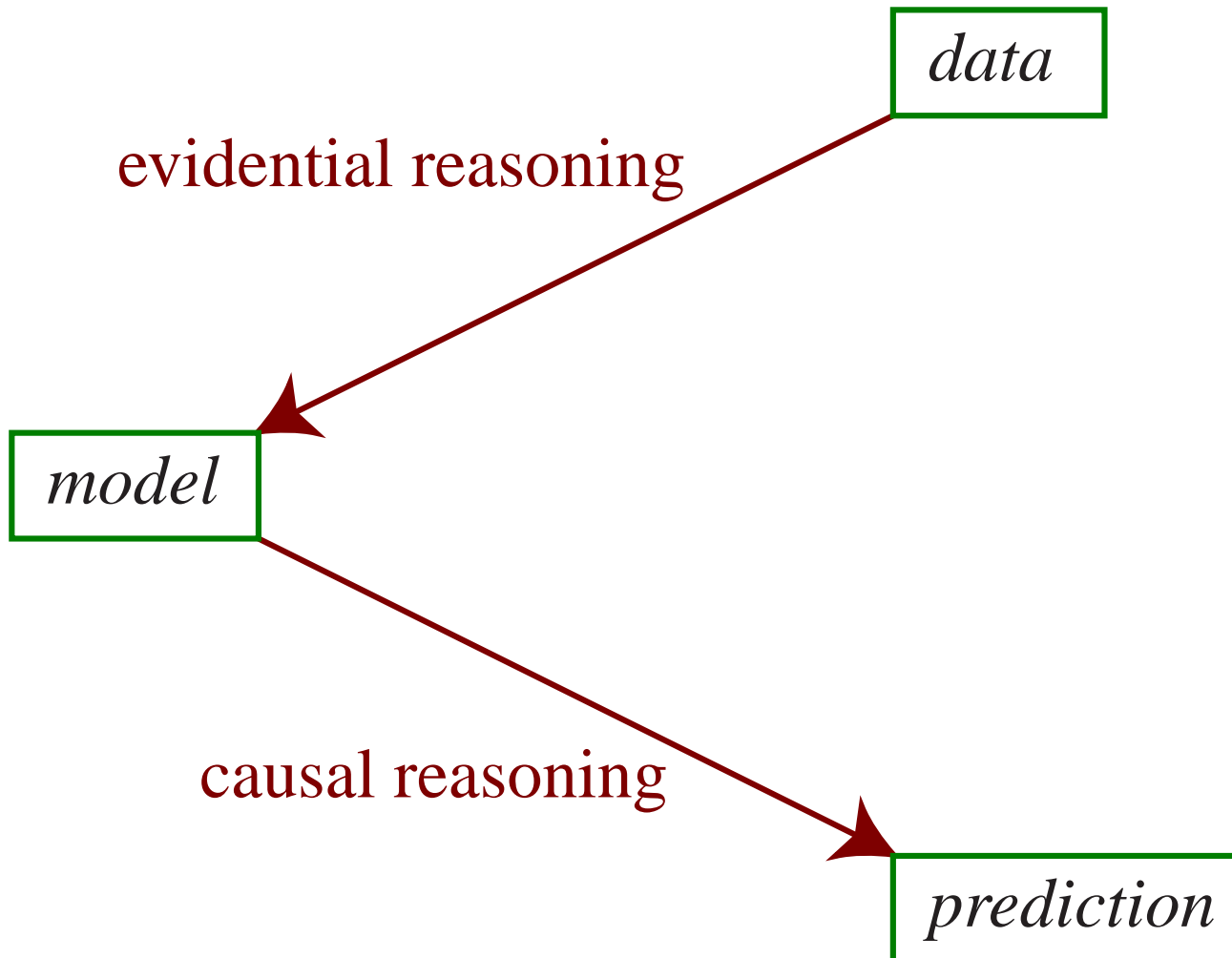
- Bayesian conditioning = abduction
- The evidence of Bayesian conditioning is what is to be explained.
 - Condition on all information obtained since the knowledge base was built.

$$P(h|e \wedge k) \longrightarrow P_k(h|e)$$

Overview

- Causal and evidential modelling & reasoning
- Bayesian networks, Bayesian conditioning & abduction
- Noise, overfitting, and Bayesian learning

Induction



Potential Confusion

Often the data is about some evidential reasoning task. e.g., classification, diagnosis, recognition ...

Example: We can do **Bayesian learning** where the hypotheses are decision trees, neural networks, parametrized distribution, or Bayesian networks.

Example: We can do **hill climbing learning with cross validation** where the hypotheses are decision trees, neural networks, parametrized distribution, or Bayesian networks.

Noise and Overfitting

Most data contains noise (errors, inadequate attributes, spurious correlations)

⇒ **overfitting** — the model learned fits random correlations in the data

Example: A more detailed decision tree *always* fits the data better, but usually smaller decision trees provides better predictive value.

Need tradeoff between

model simplicity + fit to data.

Overfitting and Bayes' theorem

fit to data

bias

$$P(h|e) = \frac{P(e|h)P(h)}{P(e)}$$

normalizing constant

Minimum description length principle

Choose best hypothesis given the evidence:

$$\arg \max_h P(h|e)$$

$$= \arg \max_h \frac{P(e|h)P(h)}{P(e)}$$

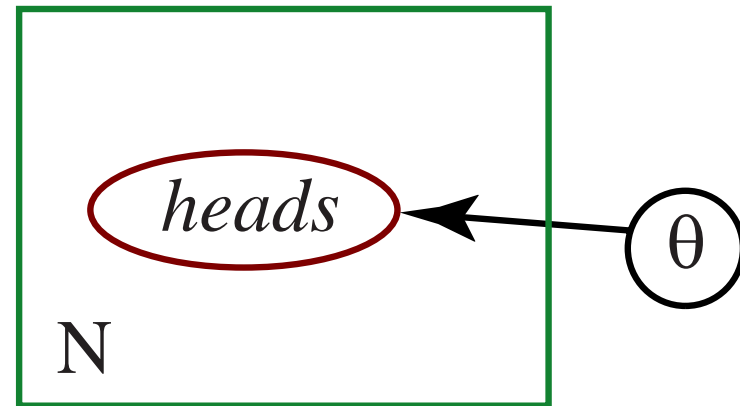
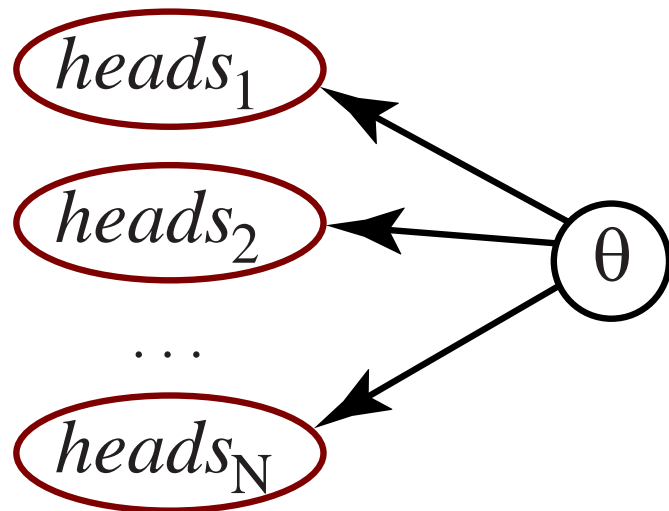
$$= \arg \max_h P(e|h)P(h)$$

$$= \arg \max_h \underbrace{-\log_2 P(e|h)}_{\substack{\text{the number of bits to} \\ \text{describe the data in} \\ \text{terms of the model}}} + \underbrace{-\log_2 P(h)}_{\substack{\text{the number of} \\ \text{bits to describe} \\ \text{the model}}}$$

Graphical Models for Learning

Idea: model \longrightarrow data

Example: parameter estimation for probability of heads (from [Buntine, JAIR, 94])



Abductive Version of Parameter Learning

$heads(C) \leftarrow$

$turns_heads(C, \Theta) \wedge prob_heads(\Theta).$

$tails(C) \leftarrow$

$turns_tails(C, \Theta) \wedge prob_heads(\Theta).$

$\forall C \forall \Theta \{turns_heads(C, \Theta), turns_tails(C, \Theta)\} \in \mathbf{C}$

$\{prob_heads(\theta) : \theta \in [0, 1]\} \in \mathbf{C}$

$Prob(turns_heads(C, \theta)) = \theta$

$Prob(turns_tails(C, \theta)) = 1 - \theta$

$Prob(prob_heads(\theta)) = 1 \quad \leftarrow \text{uniform on } [0, 1].$

Explaining Data

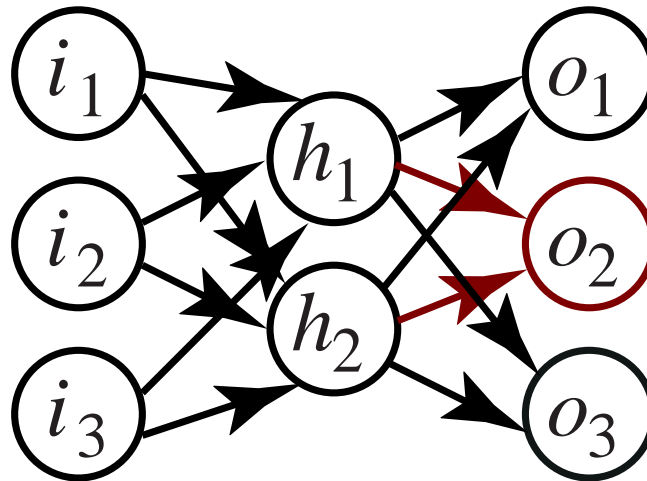
If you observe:

$heads(c_1), tails(c_2), tails(c_3), heads(c_4), heads(c_5), \dots$

For each $\theta \in [0, 1]$ there is an explanation:

$\{prob_heads(\theta), turns_heads(c_1, \theta), turns_tails(c_2, \theta),$
 $turns_tails(c_3, \theta), turns_heads(c_4, \theta), turns_heads(c_5, \theta),$
 $\dots\}$

Abductive neural-network learning



$prop(X, o_2, V) \leftarrow$

$prop(X, h_1, V_1) \wedge prop(X, h_2, V_2) \wedge$

$param(p_8, P_1) \wedge param(p_{10}, P_2) \wedge \leftarrow$ abducible

$V = \frac{1}{1 + e^{(V_1 P_1 + V_2 P_2)}} \cdot \leftarrow$ sigmoid additive

Lesson #3

Abductive view of Bayesian learning:

- rules imply data from parameters or possible representations
- parameter values or possible representations abducible
- rules contain logical variables for data items

Evidential versus causal modelling

	Neural Nets	Bayesian Nets
modelling	evidential sigmoid additive	causal linear gaussian
	— related by Bayes theorem	
evidential reasoning	direct	abduction
causal reasoning	none	direct
context changes	fragile	robust
learning	easier	more difficult

Conclusion

- Bayesian reasoning is abduction.
- Bayesian network is a causal modelling language — abduction + probability.
- What logic-based abduction can learn from Bayesians:
Handle noise, overfitting
Conditioning: explain everything
Algorithms: exploit sparse structure, exploit extreme distributions, or stochastic simulation.
- What the Bayesians can learn:
Richer representations.