

University of Rome “La Sapienza”

Master in Artificial Intelligence and Robotics

Machine Learning

A.Y. 2018/2019

Prof. Luca Iocchi

Sapienza University of Rome, Italy
Master in Artificial Intelligence and Robotics
Machine Learning (2018/19)

4. Probability and Bayes Networks

Luca Iocchi

Outline

- Uncertainty
- Probability
- Syntax and Semantics
- Inference
- Independence and Bayes' Rule

Uncertainty

Consider action $A_t = \text{leave for airport } t \text{ minutes before flight.}$

Will A_t get me there on time?

Problems:

- partial observability (road state, other drivers' plans, etc.)
- noisy sensors (KCBS traffic reports)
- uncertainty in action outcomes (flat tire, etc.)
- immense complexity of modelling and predicting traffic

Uncertainty

Hence a purely logical approach either

- risks falsehood: “ A_{25} will get me there on time” or
- leads to conclusions that are too weak for decision making: “ A_{25} will get me there on time if there’s no accident on the bridge and it doesn’t rain and my tires remain intact etc etc.”
- leads to non-optimal decisions (A_{1440} might reasonably be said to get me there on time, but I’d have to stay overnight in the airport ...)

Probability

Representation of uncertainty with probabilities.

Given the available evidence, A_{25} will get me at the airport on time with probability 0.04

Given the available evidence, A_{60} will get me at the airport on time with probability 0.85

Given the available evidence, A_{1440} will get me at the airport on time with probability 0.999

Probability

Sample space

- Ω *sample space* (set of possibilities)
- $\omega \in \Omega$ is a *sample point/possible world/atomic event/outcome of a random process/...*

Probability space (or probability model)

- Function $P : \Omega \mapsto \mathbb{R}$, such that
 - $0 \leq P(\omega) \leq 1$
 - $\sum_{\omega \in \Omega} P(\omega) = 1$

Example: rolling a dice

$$\Omega = \{1, 2, 3, 4, 5, 6\}$$

$$P(\omega) = \{1/6, 1/6, 1/6, 1/6, 1/6, 1/6\}$$

Event

An *event* A is any subset of Ω

Probability of an event A is a function assigning to A a value in $[0, 1]$

$$P(A) = \sum_{\{\omega \in A\}} P(\omega)$$

Example 1: $A_1 = \text{"dice roll"} < 4$, $A_1 = \{1, 2, 3\} \subset \Omega$

$$P(A_1) = P(1) + P(2) + P(3) = 1/6 + 1/6 + 1/6 = 1/2$$

Example 2: $A_2 = \text{"dice roll"} = 4$, $A_2 = \{4\}$, $P(A_2) = 1/6$

Example 3: $A_3 = \text{"dice roll"} > 6$, $A_3 = \emptyset$, $P(A_3) = 0$

Example 4: $A_4 = \text{"dice roll"} \leq 6$, $A_4 = \Omega$, $P(A_4) = 1$

Random variables

A *random variable* (outcome of a random phenomenon) is a function from the sample space Ω to some range (e.g., the reals or Booleans) $X : \Omega \mapsto B$.

Example: $Odd : \Omega \mapsto Boolean$.

X is a variable and a function !

$X = x_i$: the random variable X has the value $x_i \in B$

$X = x_i$ is equivalent to $\{\omega \in \Omega | X(\omega) = x_i\}$

Example: $Odd = true \equiv \{1, 3, 5\}$

Random variables

P induces a *probability distribution* for a random variable X :

$$P(X = x_i) = \sum_{\{\omega \in \Omega | X(\omega) = x_i\}} P(\omega)$$

Example

$$P(Odd = true) = P(1) + P(3) + P(5) = 1/6 + 1/6 + 1/6 = 1/2$$

Propositions

A proposition is the event (subset of Ω) where the proposition is true.

Notation for Boolean random variables: $a \equiv A = \text{true}$, $\neg a \equiv A = \text{false}$.

Given Boolean random variables A and B :

- event $a \equiv A = \text{true} \equiv \{\omega \in \Omega | A(\omega) = \text{true}\}$
- event $\neg a \equiv A = \text{false} \equiv \{\omega \in \Omega | A(\omega) = \text{false}\}$
- event $a \wedge b$ = points ω where $A(\omega) = \text{true}$ and $B(\omega) = \text{true}$
- event $\neg a \vee b$ = points ω where $A(\omega) = \text{false}$ or $B(\omega) = \text{true}$

$$P(\neg a \vee b) = \sum_{\{\omega \in \Omega | A(\omega) = \text{false} \vee B(\omega) = \text{true}\}} P(\omega)$$

Syntax for propositions

- *Propositional* or *Boolean* random variables
e.g., *Cavity* (do I have a cavity?).
Cavity = true is a proposition, also written *cavity*
- *Discrete* random variables (*finite* or *infinite*)
e.g., *Weather* is one of $\langle \text{sunny}, \text{rain}, \text{cloudy}, \text{snow} \rangle$.
Weather = rain is a proposition
Values must be exhaustive and mutually exclusive
- *Continuous* random variables (*bounded* or *unbounded*)
e.g., *Temp* = 21.6, *Temp* < 22.0.
- Arbitrary Boolean combinations of basic propositions
e.g., *cavity* \wedge *Weather* = *rain* \wedge *Temp* < 22.0.

Prior Probability

Prior or unconditional probabilities of propositions correspond to belief prior to arrival of any (new) evidence.

Examples:

$$P(\text{Odd} = \text{true}) = 0.5$$

$$P(\text{Cavity} = \text{true}) = 0.1$$

$$P(\text{Weather} = \text{sunny}) = 0.72$$

Probability distribution

A *probability distribution* is the set of probability values for all possible assignments of a random variable. Note: sum of all values must be 1.

Examples:

$$P(\text{Odd}) = \langle 0.5, 0.5 \rangle$$

$$P(\text{Cavity}) = \langle 0.1, 0.9 \rangle$$

$$P(\text{Weather}) = \langle 0.72, 0.1, 0.08, 0.1 \rangle$$

Note: for real valued random variable X , $P(X)$ is a continuous function.

Joint probability distribution

Joint probability distribution for a set of random variables gives the probability of every atomic joint event on those random variables (i.e., every sample point in the joint sample space).

Joint probability distribution of the random variables *Weather* and *Cavity*: $P(\text{Weather}, \text{Cavity})$ = a 4×2 matrix of values:

<i>Weather</i> =	<i>sunny</i>	<i>rain</i>	<i>cloudy</i>	<i>snow</i>
<i>Cavity</i> = <i>true</i>	0.144	0.02	0.016	0.02
<i>Cavity</i> = <i>false</i>	0.576	0.08	0.064	0.08

Every question about a domain can be answered by the joint distribution because every event is a sum of sample points

Conditional/Posterior Probability

Belief after the arrival of some evidence.

I know the outcome of a random variable, how this affects probability of other random variables?

Example:

I know that today *Weather* = *sunny*, how this information affects the random variable *Cavity*?

Notation:

$P(\text{Cavity} = \text{true} | \text{Weather} = \text{sunny})$: conditional/posterior probability

Conditional/Posterior Probability

In general, conditional/posterior probabilities are different from joint probabilities and from prior probabilities.

$$\begin{aligned} P(Cavity = true | Weather = sunny) &\neq \\ P(Cavity = true, Weather = sunny) &\neq \\ P(Cavity = true) \end{aligned}$$

Conditional/Posterior Probability

Consider another Boolean random variable Toothache.
Given that I have a toothache, how this affects the event of having a cavity?

Example:

$P(cavity) = 0.2$: prior

$P(cavity | toothache) = 0.6$: posterior

Conditional Probability Distributions

Conditional probability distributions: representation of all the values of a conditional probabilities of random variables.

Example:

$P(\text{Cavity}|\text{Toothache}) = 2\text{-element vector of } 2\text{-element vectors}$

Conditional probability

Definition of conditional probability:

$$P(a|b) \equiv \frac{P(a \wedge b)}{P(b)} \text{ if } P(b) \neq 0$$

Product rule

$$P(a \wedge b) = P(a|b)P(b) = P(b|a)P(a)$$

A general version holds for whole distributions, e.g.,
 $P(\text{Weather}, \text{Cavity}) = P(\text{Weather}|\text{Cavity})P(\text{Cavity})$

Total probabilities

For a Boolean random variable B

$$P(a) = P(a|b)P(b) + P(a|\neg b)P(\neg b)$$

In general, for a random variable Y accepting mutually exclusive values y_i

$$P(X) = \sum_{y_i \in \mathcal{D}(Y)} P(X|Y = y_i)P(Y = y_i)$$

$\mathcal{D}(Y)$: set of values for variable Y

Chain Rule

- *Chain rule* is derived by successive application of product rule:

$$P(X_1, X_2) = P(X_1)P(X_2|X_1)$$

$$\begin{aligned} P(X_1, \dots, X_n) &= P(X_1, \dots, X_{n-1})P(X_n|X_1, \dots, X_{n-1}) \\ &= P(X_1, \dots, X_{n-2})P(X_{n-1}|X_1, \dots, X_{n-2})P(X_n|X_1, \dots, X_{n-1}) \\ &= \dots \\ &= \prod_{i=1}^n P(X_i|X_1, \dots, X_{i-1}) \end{aligned}$$

Inference by Enumeration

Start with the joint distribution:

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

For any proposition ϕ , sum the atomic events where it is true:

$$P(\phi) = \sum_{\omega: \omega \models \phi} P(\omega)$$

Inference by Enumeration

Start with the joint distribution:

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

For any proposition ϕ , sum the atomic events where it is true:

$$P(\phi) = \sum_{\omega: \omega \models \phi} P(\omega)$$

$$P(\text{toothache}) = 0.108 + 0.012 + 0.016 + 0.064 = 0.2$$

Inference by Enumeration

Start with the joint distribution:

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

For any proposition ϕ , sum the atomic events where it is true:

$$P(\phi) = \sum_{\omega: \omega \models \phi} P(\omega)$$

$$P(\text{cavity} \vee \text{toothache}) = 0.108 + 0.012 + 0.072 + 0.008 + 0.016 + 0.064 = 0.28$$

Inference by Enumeration

Start with the joint distribution:

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

Can also compute conditional probabilities:

$$\begin{aligned}
 P(\neg \text{cavity} | \text{toothache}) &= \frac{P(\neg \text{cavity} \wedge \text{toothache})}{P(\text{toothache})} \\
 &= \frac{0.016 + 0.064}{0.108 + 0.012 + 0.016 + 0.064} = 0.4
 \end{aligned}$$

Normalization

Start with the joint distribution:

	toothache		\neg toothache	
	catch	\neg catch	catch	\neg catch
cavity	.108	.012	.072	.008
\neg cavity	.016	.064	.144	.576

Denominator can be viewed as a *normalization constant* α

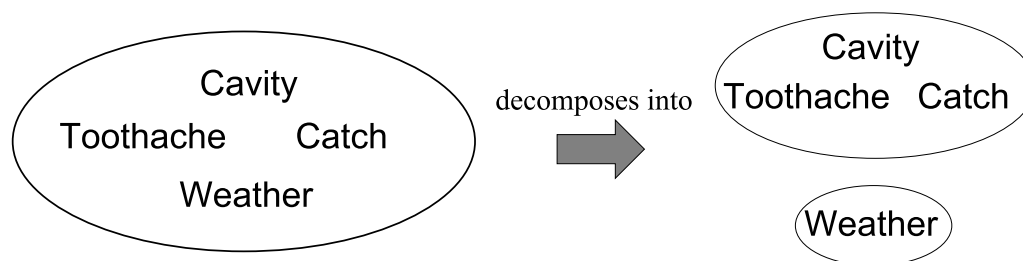
$$\begin{aligned}
 P(\text{Cavity} | \text{toothache}) &= \alpha P(\text{Cavity}, \text{toothache}) \\
 &= \alpha [P(\text{Cavity}, \text{toothache}, \text{catch}) + P(\text{Cavity}, \text{toothache}, \neg \text{catch})] \\
 &= \alpha [\langle 0.108, 0.016 \rangle + \langle 0.012, 0.064 \rangle] \\
 &= \alpha \langle 0.12, 0.08 \rangle = \langle 0.6, 0.4 \rangle
 \end{aligned}$$

General idea: compute distribution on query variable
by fixing *evidence variables* and summing over *hidden variables*

Independence

A and B are *independent* iff

$$P(A|B) = P(A) \quad \text{or} \quad P(B|A) = P(B) \quad \text{or} \quad P(A, B) = P(A)P(B)$$



$$\begin{aligned}
 P(\text{Toothache}, \text{Catch}, \text{Cavity}, \text{Weather}) &= \\
 &P(\text{Toothache}, \text{Catch}, \text{Cavity})P(\text{Weather})
 \end{aligned}$$

Independence

$P(\text{Toothache}, \text{Catch}, \text{Cavity}, \text{Weather})$ has 32 entries

$P(\text{Toothache}, \text{Catch}, \text{Cavity})$ and $P(\text{Weather})$ have $8 + 4 = 12$ entries

Example: n independent biased coins, reduced size from 2^n to n

Absolute independence powerful, but rare.

Complex systems have hundreds of variables, none of which are independent.

Conditional independence

- $P(\text{Toothache}, \text{Cavity}, \text{Catch})$ has $2^3 - 1 = 7$ independent entries
- If I have a cavity, the probability that the probe catches in it does not depend on whether I have a toothache:
(1) $P(\text{catch}|\text{toothache}, \text{cavity}) = P(\text{catch}|\text{cavity})$
- The same independence holds if I haven't got a cavity:
(2) $P(\text{catch}|\text{toothache}, \neg\text{cavity}) = P(\text{catch}|\neg\text{cavity})$
- *Catch* is *conditionally independent* of *Toothache* given *Cavity*:
 $P(\text{Catch}|\text{Toothache}, \text{Cavity}) = P(\text{Catch}|\text{Cavity})$
- Equivalent statements:
 $P(\text{Toothache}|\text{Catch}, \text{Cavity}) = P(\text{Toothache}|\text{Cavity})$
 $P(\text{Toothache}, \text{Catch}|\text{Cavity}) = P(\text{Toothache}|\text{Cavity})P(\text{Catch}|\text{Cavity})$

Conditional independence

General formulation:

X conditionally independent from Y given Z iff $P(X|Y, Z) = P(X|Z)$

$$P(X, Y|Z) = P(X|Y, Z)P(Y|Z) = P(X|Z)P(Y|Z)$$

$$P(Y_1, \dots, Y_n|Z) = P(Y_1|Y_2, \dots, Y_n, Z)P(Y_2|Y_3 \dots Y_n, Z) \cdots P(Y_n|Z)$$

Y_i conditionally independent from Y_j given Z

$$P(Y_1, \dots, Y_n|Z) = P(Y_1|Z)P(Y_2|Z) \cdots P(Y_n|Z)$$

Conditional independence

Chain rule + Conditional independence

$$\begin{aligned} P(X, Y, Z) &= P(X|Y, Z)P(Y, Z) = P(X|Y, Z)P(Y|Z)P(Z) \\ &= P(X|Z)P(Y|Z)P(Z) \end{aligned}$$

$$\begin{aligned} &P(\textit{Toothache}, \textit{Catch}, \textit{Cavity}) \\ &= P(\textit{Toothache}|\textit{Catch}, \textit{Cavity})P(\textit{Catch}, \textit{Cavity}) \\ &= P(\textit{Toothache}|\textit{Catch}, \textit{Cavity})P(\textit{Catch}|\textit{Cavity})P(\textit{Cavity}) \\ &= P(\textit{Toothache}|\textit{Cavity})P(\textit{Catch}|\textit{Cavity})P(\textit{Cavity}) \\ &2 + 2 + 1 = 5 \text{ independent numbers (instead of } 2^3 - 1) \end{aligned}$$

In most cases, the use of conditional independence reduces the size of the representation of the joint distribution from exponential in n to linear in n .

Bayes' Rule

- Product rule $P(a \wedge b) = P(a|b)P(b) = P(b|a)P(a)$

$$\Rightarrow \text{Bayes' rule } P(a|b) = \frac{P(b|a)P(a)}{P(b)}$$

or in distribution form

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)} = \alpha P(X|Y)P(Y)$$

Useful for assessing *diagnostic* probability from *causal* probability:

$$P(\text{Cause}|\text{Effect}) = \frac{P(\text{Effect}|\text{Cause})P(\text{Cause})}{P(\text{Effect})}$$

Bayes' Rule and conditional independence

Bayes rule

$$P(Z|Y_1, \dots, Y_n) = \alpha P(Y_1, \dots, Y_n|Z) P(Z)$$

Y_1, \dots, Y_n conditionally independent each other given Z

$$P(Z|Y_1, \dots, Y_n) = \alpha P(Y_1|Z) \cdots P(Y_n|Z) P(Z)$$

Effects conditionally independent each other given a cause.

$$P(\text{Cause}|\text{Effect}_1, \dots, \text{Effect}_n) = \alpha P(\text{Cause}) \prod_i P(\text{Effect}_i|\text{Cause})$$

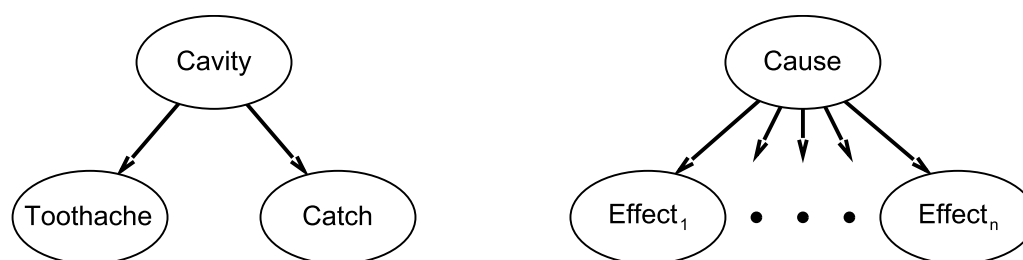
Total number of parameters is *linear* in n

Bayesian networks



A simple, graphical notation for conditional independence assertions and hence for compact specification of full joint distributions.

Bayesian networks



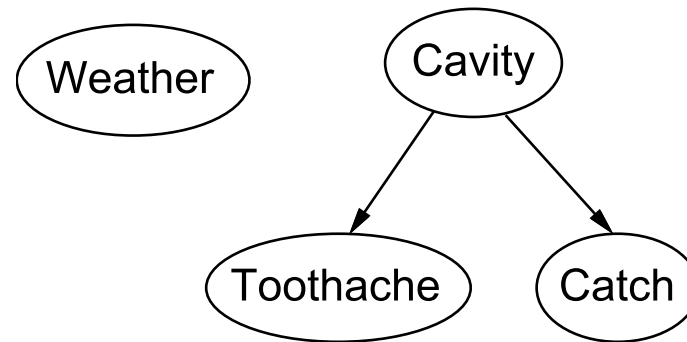
Syntax:

- a set of nodes, one per variable
- a directed, acyclic graph (link \approx “directly influences”)
- a conditional distribution for each node given its parents:
 $P(X_i | \text{Parents}(X_i))$

In the simplest case, conditional distribution represented as a *conditional probability table* (CPT) giving the distribution over X_i for each combination of parent values.

Dentist BN Example

Topology of network encodes conditional independence assertions:

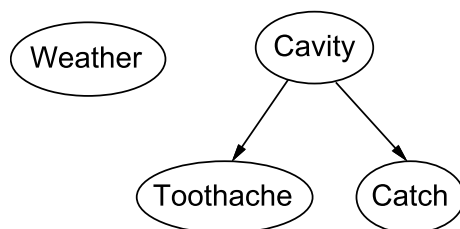


Weather is independent of the other variables

Toothache and *Catch* are conditionally independent given *Cavity*

Dentist BN Example

BN model given by the set of CPT $P(X_i | \text{Parents}(X_i))$ for each variable X_i



$P(\text{Weather})$
 $P(\text{Cavity})$
 $P(\text{Toothache} | \text{Cavity})$
 $P(\text{Catch} | \text{Cavity})$

All the joint probabilities can be computed from this model.
How many independent values?

Burglar BN Example

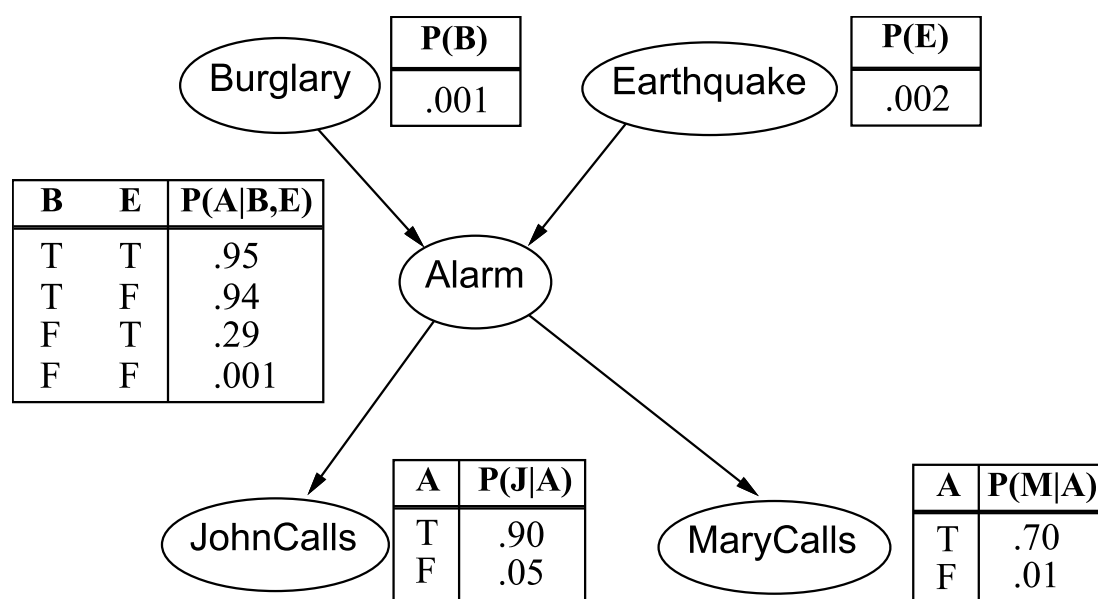
I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometimes the alarm is set off by minor earthquakes. Is there a burglar?

Variables: *Burglar*, *Earthquake*, *Alarm*, *JohnCalls*, *MaryCalls*

Network topology reflects "causal" knowledge:

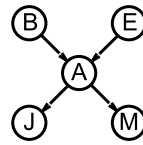
- A burglar can set the alarm
- An earthquake can set the alarm
- The alarm can cause Mary to call
- The alarm can cause John to call

Burglar BN Example



Compactness

A CPT for Boolean variable X_i with k Boolean parents has 2^k rows for the combinations of parent values



Each row requires one number p for $X_i = \text{true}$ (the number for $X_i = \text{false}$ is just $1 - p$)

If each variable has no more than k parents, the complete network requires $O(n \cdot 2^k)$ numbers

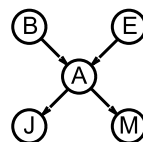
I.e., grows linearly with n , vs. $O(2^n)$ for the full joint distribution

For burglary net, $1 + 1 + 4 + 2 + 2 = 10$ numbers (vs. $2^5 - 1 = 31$)

Computing joint probabilities

All joint probabilities computed with the chain rule:

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{Parents}(X_i))$$



e.g., $P(j \wedge m \wedge a \wedge \neg b \wedge \neg e)$

$$\begin{aligned} &= P(j|a)P(m|a)P(a|\neg b, \neg e)P(\neg b)P(\neg e) \\ &= 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998 \\ &\approx 0.00063 \end{aligned}$$

Classification as Probabilistic estimation

Given target function $f : X \rightarrow V$, dataset D and a new instance x' , best prediction $\hat{f}(x') = v^*$

$$v^* = \operatorname{argmax}_{v \in V} P(v|x', D)$$

More general formulation: compute the probability distribution over V

$$P(V|x', D)$$

Learning as Probabilistic estimation

Given dataset D and hypothesis space H , compute a probability distribution over H given D .

$$P(H|D)$$

Bayes rule

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

- $P(h)$ = prior probability of hypothesis h
- $P(D)$ = prior probability of training data D
- $P(h|D)$ = probability of h given D
- $P(D|h)$ = probability of D given h