

# Adapted from AIMA slides

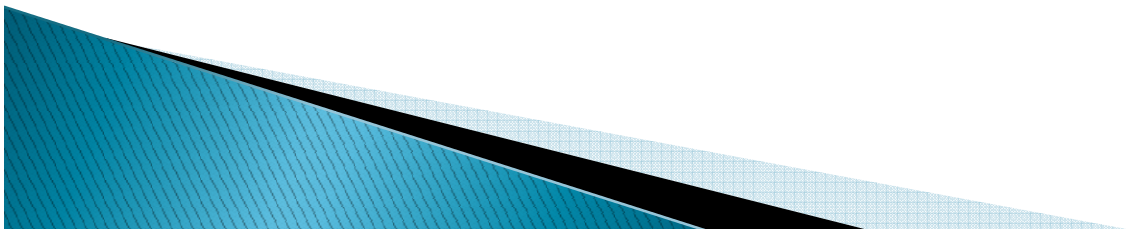
## Simple probabilistic models

Peter Antal

[antal@mit.bme.hu](mailto:antal@mit.bme.hu)

# Outline

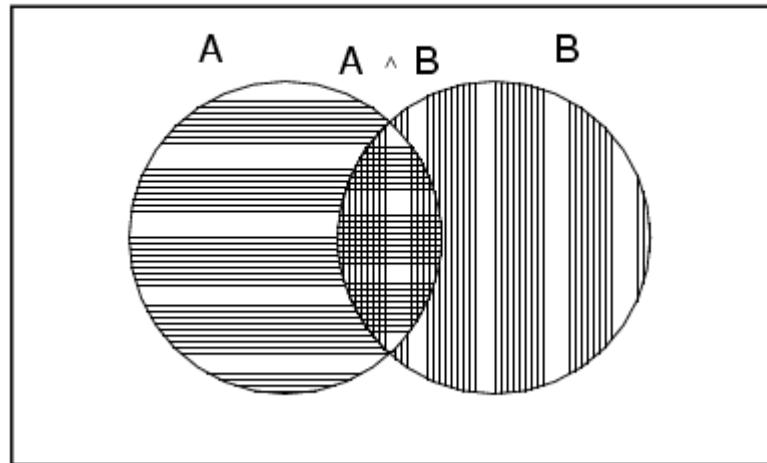
- ▶ Basics of probability theory
- ▶ Relation of two-valued vs probabilistic logic
  - Truth vs belief
  - Proofs vs inference
- ▶ Naïve Bayesian networks
- ▶ SPAM filter
- ▶ Special local models
  - Noisy-OR
  - Decision tree CPDs
  - Decision graph CPDs



# Axioms of probability

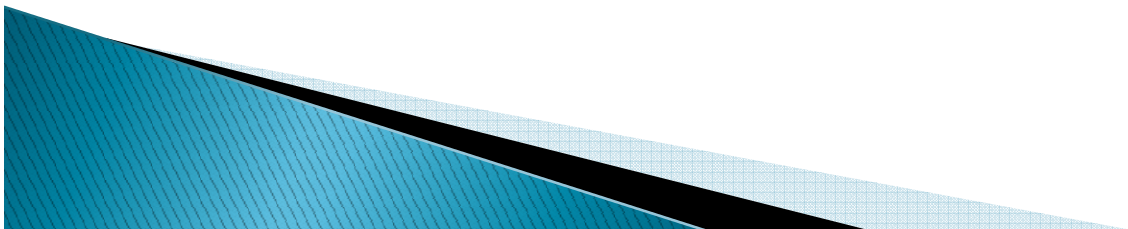
- ▶ For any propositions  $A, B$
- ▶
  - $0 \leq P(A) \leq 1$
  - $P(\text{true}) = 1$  and  $P(\text{false}) = 0$
  - $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$
  -

True



# Probability theory: concepts for the course

- ▶ Joint distribution (“omic-ness”)
  - (“omic-ness”: “comprehensiveness” + “query-free”)
- ▶ Conditional probability (“simple inference”)
- ▶ Chain rule (“factorization”)
- ▶ Bayes’ rule (“inversion”)
- ▶ Marginalization/expansion (“complex inference”)
- ▶ [Conditional] independence (“simplification”)



# Bayes' Rule

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

▶

⇒ **Bayes' rule:**  $P(a | b) = P(b | a) P(a) / P(b)$

▶ or in distribution form

▶

$$P(Y|X) = P(X|Y) P(Y) / P(X) = \alpha P(X|Y) P(Y)$$

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

▶

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

◦

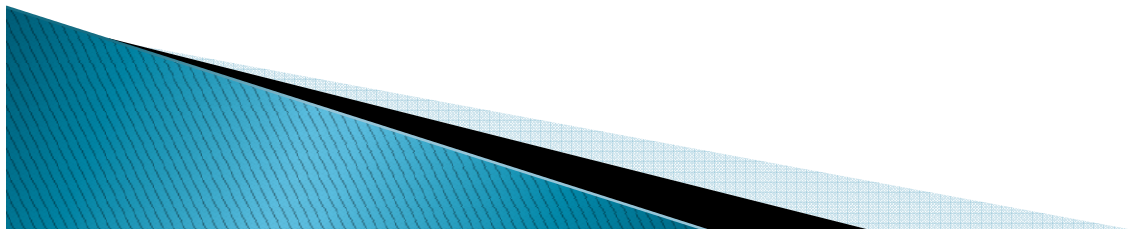
◦ E.g., let  $M$  be meningitis,  $S$  be stiff neck:

◦

$$P(m|s) = P(s|m) P(m) / P(s) = 0.8 \times 0.0001 / 0.1 = 0.0008$$

◦ Note: posterior probability of meningitis still very small!

◦



# Bayes rule

An algebraic triviality

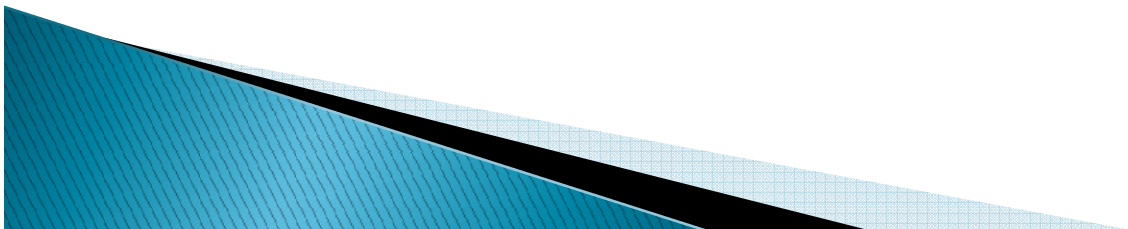
$$p(X | Y) = \frac{p(Y | X)p(X)}{p(Y)} = \frac{p(Y | X)p(X)}{\sum_x p(Y | X)p(X)}$$

A scientific research paradigm

$$p(\textit{Model} | \textit{Data}) \propto p(\textit{Data} | \textit{Model})p(\textit{Model})$$

A practical method for inverting causal knowledge to diagnostic tool.

$$p(\textit{Cause} | \textit{Effect}) \propto p(\textit{Effect} | \textit{Cause}) \times p(\textit{Cause})$$



# Inference by enumeration

Every question about a domain can be answered by the joint distribution.

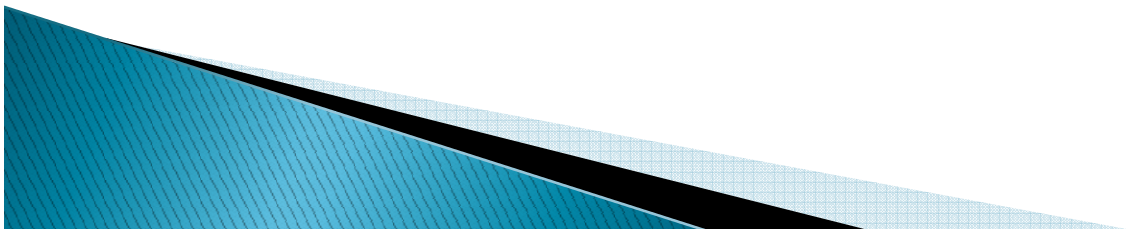
Typically, we are interested in the posterior joint distribution of the **query variables**  $Y$  given specific values  $e$  for the **evidence variables**  $E$

Let the **hidden variables** be  $H = X - Y - E$

Then the required summation of joint entries is done by summing out the hidden variables:

$$P(Y \mid E = e) = \alpha P(Y, E = e) = \alpha \sum_h P(Y, E = e, H = h)$$

- ▶ The terms in the summation are joint entries because  $Y$ ,  $E$  and  $H$  together exhaust the set of random variables
- ▶ Obvious problems:
  1. Worst-case time complexity  $O(d^n)$  where  $d$  is the largest arity
  2. Space complexity  $O(d^n)$  to store the joint distribution
  3. How to find the numbers for  $O(d^n)$  entries?



# Decision theory = probability theory + utility theory

## ▶ Decision situation:

- Actions
- Outcomes
- Probabilities of outcomes
- Utilities/losses of outcomes
  - QALY, micromort
- Maximum Expected Utility Principle (MEU)
  - Best action is the one with maximum expected utility

$$a_i$$

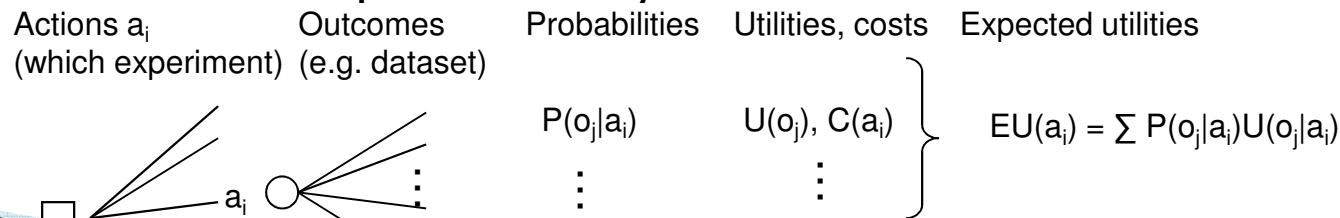
$$o_j$$

$$p(o_j | a_i)$$

$$U(o_j | a_i)$$

$$EU(a_i) = \sum_j U(o_j | a_i) p(o_j | a_i)$$

$$a^* = \arg \max_i EU(a_i)$$





# About the event space

- ▶ Atomic events are mutually exclusive and exhaustive.
- ▶ The single variable case.
  - *Weather* is one of  $\langle \text{sunny, rainy, cloudy, snow} \rangle$
  - $P((\text{Weather} = \text{sunny}) \vee (\text{Weather} = \text{rainy}))$
- ▶ Challenges in the multivariate case.
  - *Weather* is one of  $\langle \text{sunny, rainy, cloudy, snow} \rangle$
  - *TemperatureofRain* is one of  $\langle \text{icy, cold, warm} \rangle$ 
    - NONE?

# Classical vs probabilistic logic: truth and beliefs

$P_1$	...	$P_3$	KB	s	pKB	P(query evidence)
F	F	F	F	T	.01	.1
F	F	T	T	F	.12	.2
F	T	F	F	T	.35	.3
F	T	T	F	F	..	..
T	F	F	F	T	..	..
T	F	T	T	T	...	...
T	T	F	F	T	..	..
T	T	T	F	T	..	..

# Classical vs probabilistic logic: provability and inference

$P( KB \vdash_i \alpha = \text{sentence } \alpha \text{ can be derived from } KB \text{ by procedure } i / pKB )$

*“Belief propagation in networks”*

# Conditional independence



„Probability theory=measure theory+independence”

$I_p(X;Y|Z)$  or  $(X \perp\!\!\!\perp Y|Z)_p$  denotes that  $X$  is independent of  $Y$  given  $Z$ :

$$P(X;Y|z) = P(Y|z) P(X|z) \text{ for all } z \text{ with } P(z) > 0.$$

(Almost) alternatively,  $I_p(X;Y|Z)$  iff

$$P(X|Z,Y) = P(X|Z) \text{ for all } z,y \text{ with } P(z,y) > 0.$$

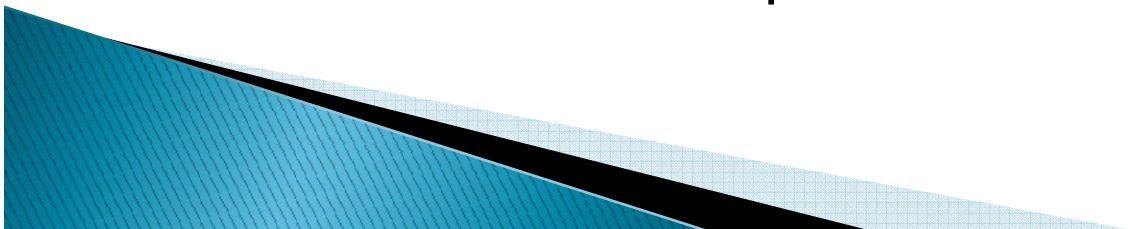
Other notations:  $D_p(X;Y|Z) = \text{def} = \neg I_p(X;Y|Z)$

Contextual independence: for not all  $z$ .

Homeworks:

Intransitivity: show that it is possible that  $D(X;Y)$ ,  $D(Y;Z)$ , but  $I(X;Z)$ .

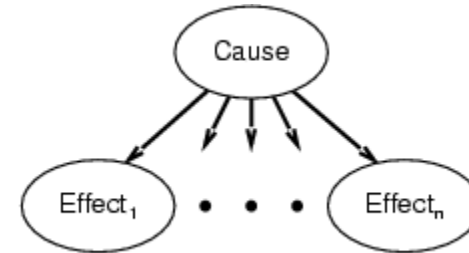
order : show that it is possible that  $I(X;Z)$ ,  $I(Y;Z)$ , but  $D(X,Y;Z)$ .



# Naive Bayesian network

Assumptions:

- 1, Two types of nodes: a cause and effects.
- 2, Effects are conditionally independent of each other given their cause.



## Variables (nodes)

Flu: present/absent

FeverAbove38C: present/absent

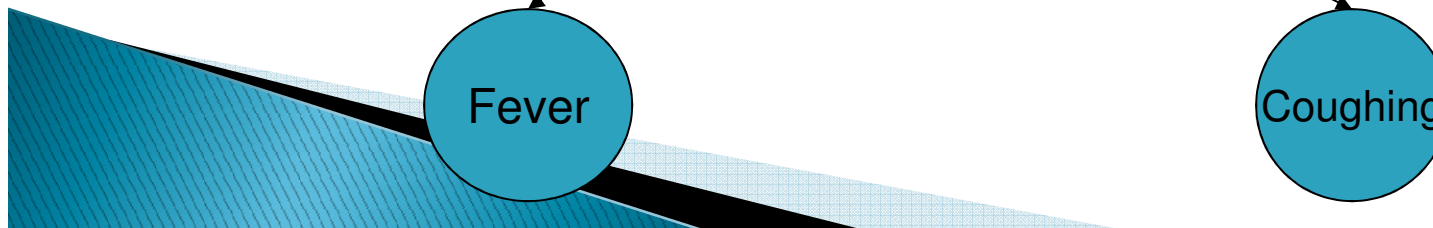
Coughing: present/absent

## Model

$P(\text{Fever}=\text{present}|\text{Flu}=\text{present})=0.6$   
 $P(\text{Fever}=\text{absent}|\text{Flu}=\text{present})=1-0.6$   
 $P(\text{Fever}=\text{present}|\text{Flu}=\text{absent})=0.01$   
 $P(\text{Fever}=\text{absent}|\text{Flu}=\text{absent})=1-0.01$

$P(\text{Flu}=\text{present})=0.001$   
 $P(\text{Flu}=\text{absent})=1-P(\text{Flu}=\text{present})$

$P(\text{Coughing}=\text{present}|\text{Flu}=\text{present})=0.3$   
 $P(\text{Coughing}=\text{absent}|\text{Flu}=\text{present})=1-0.3$   
 $P(\text{Coughing}=\text{present}|\text{Flu}=\text{absent})=0.02$   
 $P(\text{Coughing}=\text{absent}|\text{Flu}=\text{absent})=1-0.02$



# Naive Bayesian network (NBN)

Decomposition of the joint:

$$\begin{aligned} P(Y, X_1, \dots, X_n) &= P(Y) \prod_i P(X_i | Y, X_1, \dots, X_{i-1}) && // \text{by the chain rule} \\ &= P(Y) \prod_i P(X_i | Y) && // \text{by the N-BN assumption} \end{aligned}$$

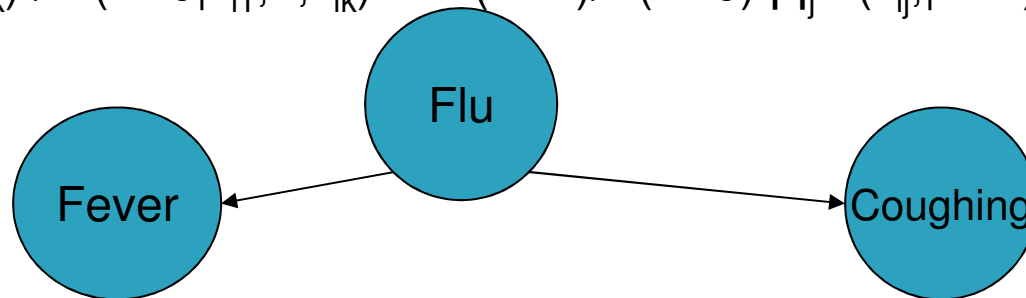
$2n+1$  parameteres!

Diagnostic inference:

$$P(Y | x_{i1}, \dots, x_{ik}) = P(Y) \prod_j P(x_{ij} | Y) / P(x_{i1}, \dots, x_{ik})$$

If  $Y$  is binary, then the odds

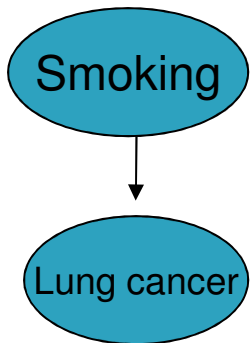
$$P(Y=1 | x_{i1}, \dots, x_{ik}) / P(Y=0 | x_{i1}, \dots, x_{ik}) = P(Y=1) / P(Y=0) \prod_j P(x_{ij} | Y=1) / P(x_{ij} | Y=0)$$



$$p(\text{Flu} = \text{present} \mid \text{Fever} = \text{absent}, \text{Coughing} = \text{present})$$

$$\propto p(\text{Flu} = \text{present}) p(\text{Fever} = \text{absent} \mid \text{Flu} = \text{present}) p(\text{Coughing} = \text{present} \mid \text{Flu} = \text{present})$$

# Conditional probabilities, odds, odds ratios



	$\neg S$	S	
$\neg LC$	$P(\neg S, \neg LC)$	$P(S, \neg LC)$	$P(\neg LC)$
LC	$P(\neg S, LC)$	$P(S, LC)$	$P(LC)$
	$P(\neg S)$	$P(S)$	

## Probability:

$P(LC)$

## Conditional probabilities (e.g., probability of LC given S):

$P(LC | \neg S) = ???$   $P(LC | S) = ???$   $P(LC)$

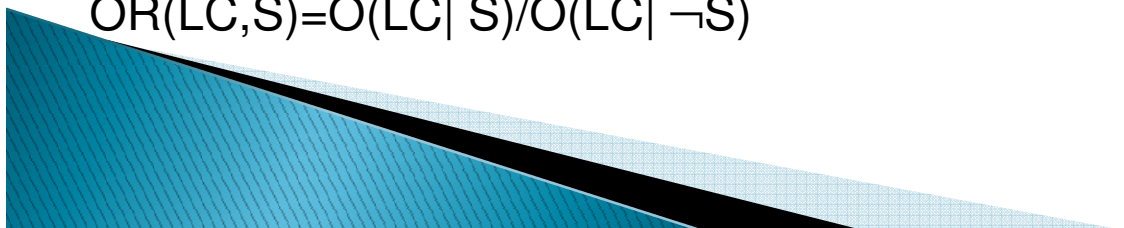
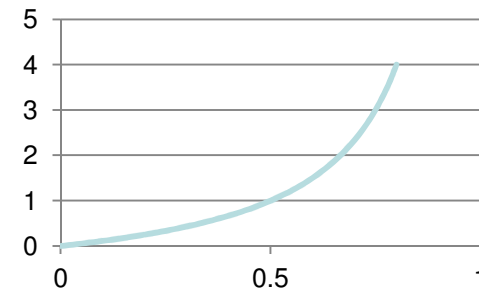
## Odds:

$[0, 1] \rightarrow [0, \infty]$ :  $Odds(p) = p / (1 - p)$

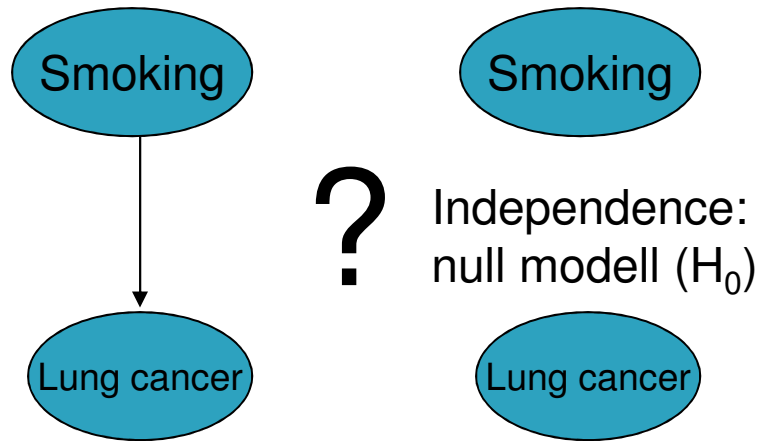
$O(LC | \neg S) = ???$   $O(LC | S)$

## Odds Ratio (OR) Independent of prevalence!

$OR(LC, S) = O(LC | S) / O(LC | \neg S)$



# Probabilities, odds, odds ratios



	$\neg S$	S	
$\neg LC$	8	7	15
LC	1	4	5
	9	11	20

Contingency table with marginals

	$\neg S$	S	
$\neg LC$	.4	.35	.75
LC	.05	.2	.25
	.45	.55	

## Conditional probabilities:

$P(LC | \neg S) = .11$  ???  $P(LC | S) = .36$  ???  $P(LC) = .25$

## Odds:

$[0, 1] \rightarrow [0, \infty]$ :  $Odds(p) = p / (1 - p)$

$O(LC | \neg S) = .12$  ???  $O(LC | S) = .56$

## Odds Ratio (OR):

$OR(LC, S) = O(LC | S) / O(LC | \neg S) = 4.6$

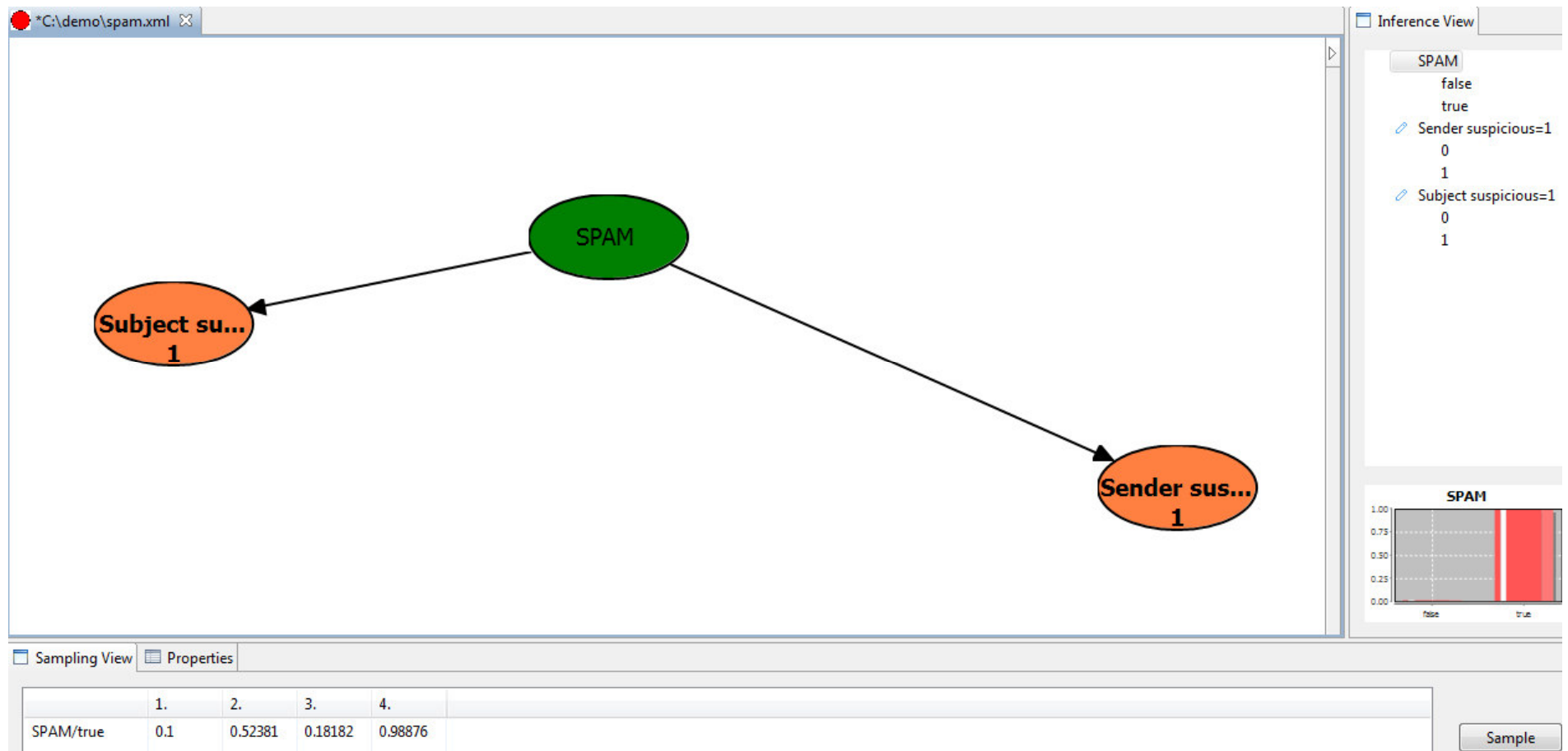


# BAYES CUBE (~BAYES EYE)



<http://redmine.genagrid.eu/projects/bayescubedownload/wiki/Wiki>

# Example: Construct a spam filter



# Summary

- ▶ Naïve Bayesian networks (N-BNs) demonstrate the use of independencies to cope with
  - model complexity ( $\sim$ space complexity, number of parameters)
  - inferential complexity ( $\sim$ time complexity).
- ▶ The assumption of conditional independence of the effects given their common cause allows
  - the efficient representation of the joint distribution
    - (in the discrete, multinomial case: linear number of parameters instead of exponential),
  - the efficient computation of the diagnostic posterior  $p(Y|X)$ 
    - (linear number of steps instead of exponential),
- ▶ Odds, log odds are popular transformations of probabilities.
- ▶ N-BNs are robust knowledge engineering and data analysis tools.
- ▶
- ▶ **Suggested reading:**
  - Druzdzell: Building Probabilistic Networks: Where Do the Numbers Come From?, IEEE Transactions on Knowledge and data engineering, 2000

