

Week 21: Bayesian Networks

Martha Lewis
(based on slides from Raul Santos Rodriguez)

... Russell and Norvig (Ch. 14)

... David Barber's Bayesian reasoning and Machine Learning:

<http://www.cs.ucl.ac.uk/staff/d.barber/brml/>

This week we look at probabilistic graphical models. These are very powerful representations that enable us to scale probabilistic models to large real problems. The objective is to learn the following topics:

- Bayes's theorem
- Definition of a DAG (or Bayesian network)
- Conditional independence in DAGs
- Dynamic Bayesian networks
- Hidden Markov Models

Bayes's theorem: The Horrible Disease?

You are about to be tested for a rare disease. How worried should you be if the test result is positive? The doctor has bad news and good news...

Bad You tested positive for a serious disease, and that the test is 99% accurate (*i.e., the probability of testing positive given that you have the disease is 0.99, as is the probability of testing negative given that you don't have the disease*).

Good This is a rare disease, striking only 1 in 10,000 people.

What are the chances that you actually have the disease?

Bayes's theorem



$$\underline{P(A | B)} = \frac{\underline{P(B | A)}\underline{P(A)}}{\underline{P(B)}} = \frac{P(A \cap B)}{P(B)}$$

Bayes's theorem

Likelihood

Probability of collecting this data when our hypothesis is true

Bill Howe, UW

Prior

The probability of the hypothesis being true before collecting data

$$\underbrace{P(H|D)} = \frac{\underbrace{P(D|H)} \underbrace{P(H)}}{\underbrace{P(D)}}$$

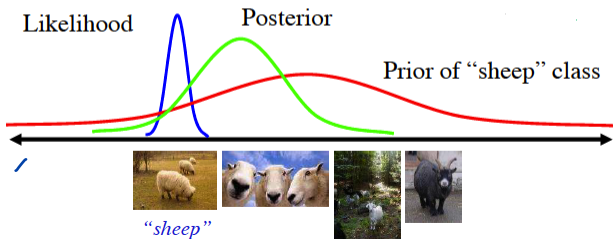
Posterior

The probability of our hypothesis being true given the data collected

Marginal

What is the probability of collecting this data under all possible hypotheses?

$$p(h|d) = \frac{p(d|h)p(h)}{\sum_{h' \in H} p(d|h')p(h')}$$



Bayes's theorem: example

The test is 99% accurate:

$$P(\underline{T} = 1 | \underline{D} = 1) = 0.99$$

$$P(T = 0 | D = 0) = 0.99$$

where T denotes test and D denotes disease.

The disease affects 1 in 10000:

$$P(D = 1) = 0.0001$$

$$P(\underline{D} = 1 | \underline{T} = 1) = \frac{\overset{= 0.99}{P(\underline{T} = 1 | \underline{D} = 1)} \overset{= 0.0001}{P(D = 1)}}{\underset{= 0.01}{P(\underline{T} = 1 | \underline{D} = 0)} \underset{= 0.9999}{P(D = 0)} + \underset{= 0.99}{P(T = 1 | D = 1)} \underset{= 0.0001}{P(D = 1)}}$$

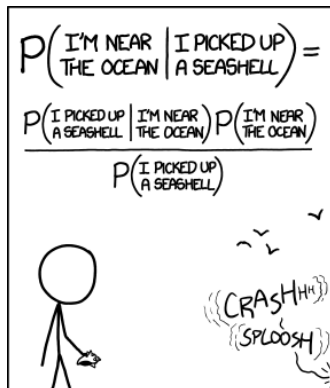
$P(D = 1 | T = 1) = 0.0098$

[https:

//www.intelligentinvestor.com.au/the-theory-that-cracked-the-enigma-code-1811436]



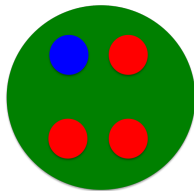
Bayes's theorem



STATISTICALLY SPEAKING, IF YOU PICK UP A SEASHELL AND DON'T HOLD IT TO YOUR EAR, YOU CAN PROBABLY HEAR THE OCEAN.

A simple graphical model: Example

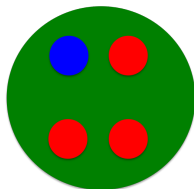
Problem: drawing balls from the set $\{r, r, r, b\}$



What are $\underline{P(B_1)}$, $\underline{P(B_2|B_1)}$ and $\underline{P(B_1, B_2)}$?

A simple graphical model: Example

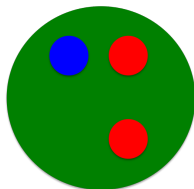
Problem: drawing balls from the set $\{r, r, r, b\}$



$$P(B_1) = \begin{array}{|c|c|} \hline \bullet B_1 & ? \\ \hline \bullet B_1 & ? \\ \hline \end{array}$$

A simple graphical model: Example

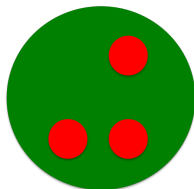
Problem: drawing balls from the set $\{r, r, r, b\}$



$$P(B_1) = \begin{array}{|c|c|} \hline r & 3/4 \\ \hline b & ? \\ \hline \end{array}$$

A simple graphical model: Example

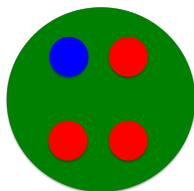
Problem: drawing balls from the set $\{r, r, r, b\}$



$$P(B_1) = \begin{array}{|c|c|} \hline B_1 & 3/4 \\ \hline B_1 & 1/4 \\ \hline \end{array}$$

A simple graphical model: Example

Problem: drawing balls from the set $\{r, r, r, b\}$

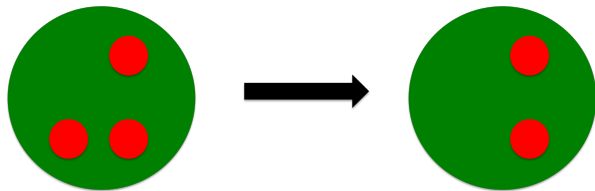


$$P(\underline{B_2} | \underline{B_1}) =$$

	B_2	B_2
B_1	?	?
B_1	?	?

A simple graphical model: Example

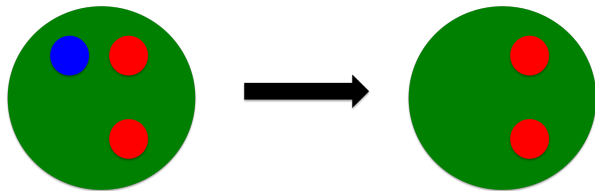
Problem: drawing balls from the set $\{r, r, r, b\}$



$$P(B_2|B_1) = \begin{array}{c|cc} & B_2 & B_2 \\ \hline B_1 & ? & ? \\ \hline B_1 & 1 & 0 \end{array}$$

A simple graphical model: Example

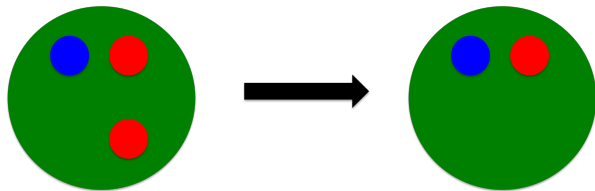
Problem: drawing balls from the set $\{r, r, r, b\}$



$$P(B_2|B_1) = \begin{array}{c|cc} & B_2 & B_2 \\ \hline B_1 & ? & 1/3 \\ \hline B_1 & 1 & 0 \end{array}$$

A simple graphical model: Example

Problem: drawing balls from the set $\{r, r, r, b\}$

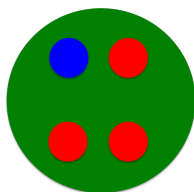


$$P(B_2|B_1) = \begin{array}{c|cc} & B_2 & B_2 \\ \hline B_1 & 2/3 & 1/3 \\ \hline B_1 & 1 & 0 \end{array}$$

A simple graphical model: Example

Problem: drawing balls from the set $\{r, r, r, b\}$

$$P(B_1, B_2) := P(B_1 \wedge B_2)$$


$$P(B_1, B_2) =$$

	B_2	B_2
B_1	$\boxed{1/2}$	$\boxed{1/4}$
B_1	$1/4$	0

A simple graphical model: Example

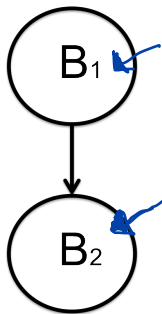
Problem: drawing balls from the set $\{r, r, r, b\}$

$$P(B_1) = \begin{array}{|c|c|} \hline B_1 & 3/4 \\ \hline B_1 & 1/4 \\ \hline \end{array}$$

$$P(B_2|B_1) = \begin{array}{|c|c|c|} \hline & B_2 & B_2 \\ \hline B_1 & 2/3 & 1/3 \\ \hline B_1 & 1 & 0 \\ \hline \end{array}$$

$$P(B_1, B_2) = \begin{array}{|c|c|c|} \hline & B_2 & B_2 \\ \hline B_1 & 1/2 & 1/4 \\ \hline B_1 & 1/4 & 0 \\ \hline \end{array}$$

$$P(B_1 \wedge B_2) = P(B_2|B_1)P(B_1)$$



- Bayes's theorem allows us to calculate the probability of a hypothesis being true given the data
- We looked at some worked examples

Directed probabilistic graphical models

A **Bayesian Network** is a directed graph in which each node is annotated with quantitative probability information.

Recall that we have $P(x_1, \dots, x_n) = \prod_i P(x_i | x_{i+1}, \dots, x_n)$. Calculating this becomes intractable.

Bayesian Network

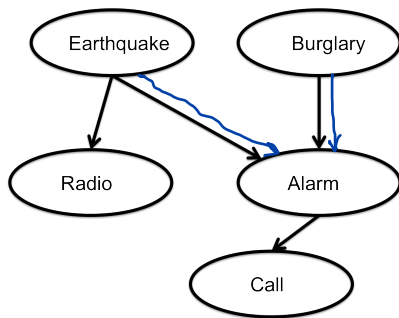
- Each **node** corresponds to a random variable (discrete or continuous).
- Each **edge** represents influence.
- The graph has no directed cycles: **directed acyclic graph** (DAG).
- Each node x_i has a conditional probability distribution $P(x_i | \text{Parents}(x_i))$ that quantifies the effect of the **parents** on the node.

The DAG tells us that if we have n variables x_i , the joint distribution of these variables factorises as follows:

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

Directed probabilistic graphical models: example

A **Bayesian Network** is a directed graph in which each node is annotated with quantitative probability information.



Russell+Norvig

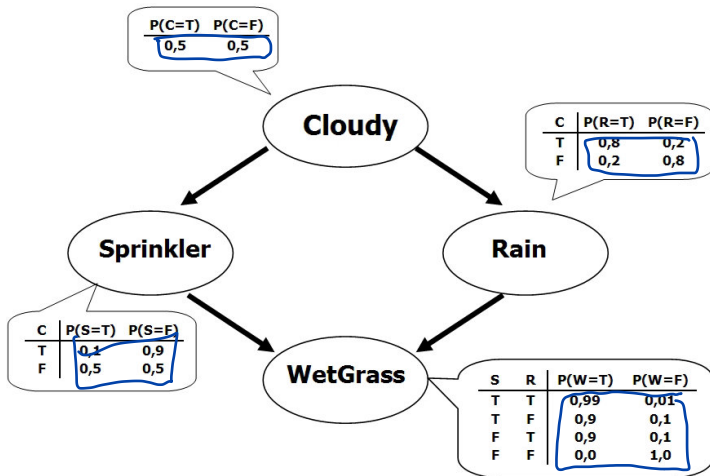
$$P(A, R, E, B) = P(A|R, E, B)P(R|E, B)P(E|B)P(B)$$

Can be reduced to:

$$P(A, R, E, B) = P(A|E, B)P(R|E)P(E)P(B)$$

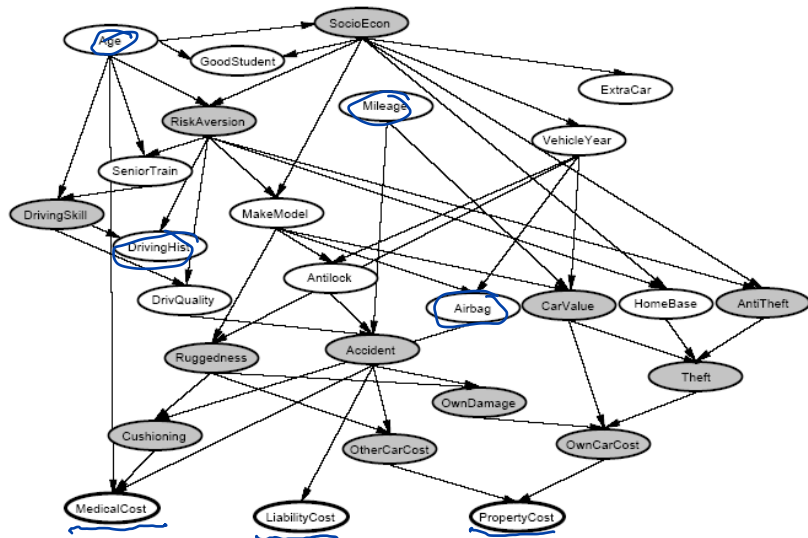
Joint vs Factorised joint distributions

$$p(C, S, R, W) = p(C)p(S|C)p(R|C)p(W|S, R)$$

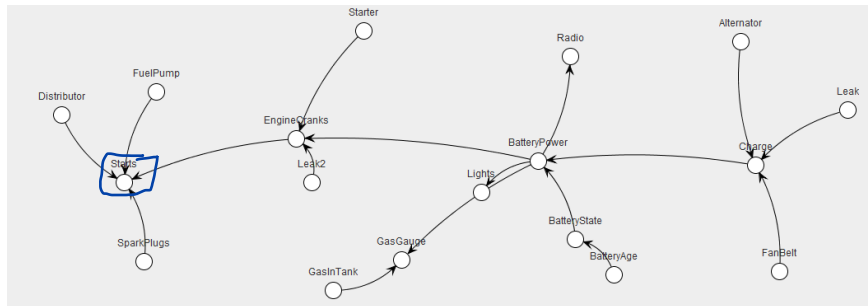


Joint distribution → 15 values vs Factorised joint distribution → 9 values

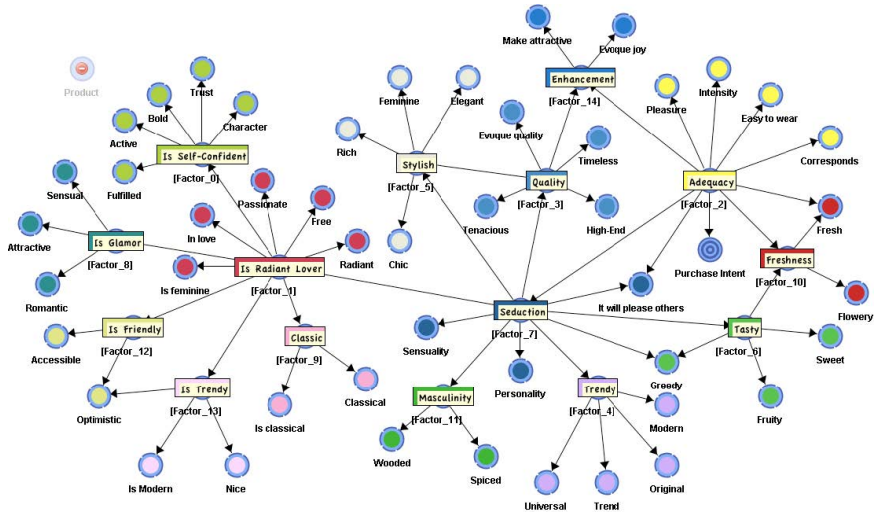
Example: insurance



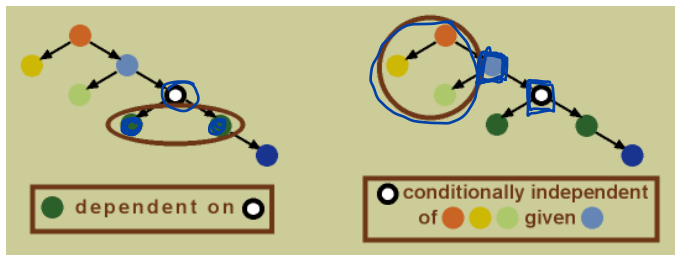
Example: car engine



Example: perfume market analysis

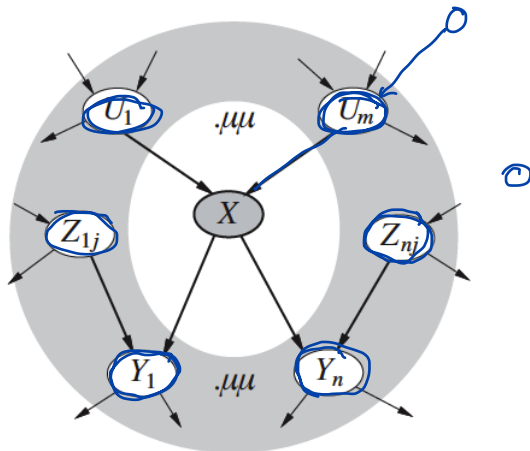


Conditional independence



Conditional independence: Markov Blanket

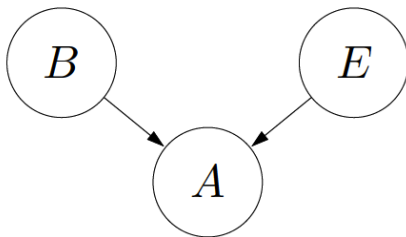
Given its **Markov Blanket**, X is independent of all other nodes



$$MB(X) = Parents(X) \cup Children(X) \cup Parents(Children(X))$$

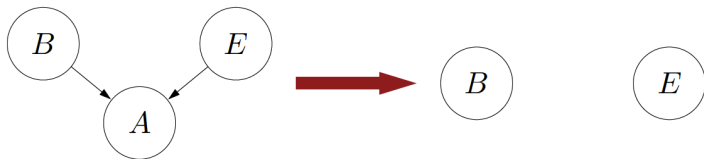
- A directed graphical model allows us to represent the dependencies in a Bayesian network.
- By factoring out the network we can calculate probabilities more efficiently.
- In applications this problem can become huge!

Given the following Bayesian Network,

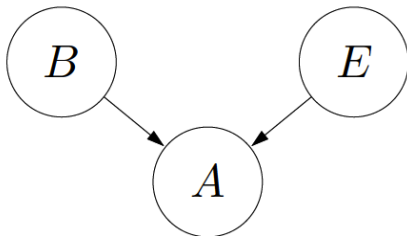


show that B and E are independent, given A

Marginalisation of a leaf node yields a Bayesian network without the node.

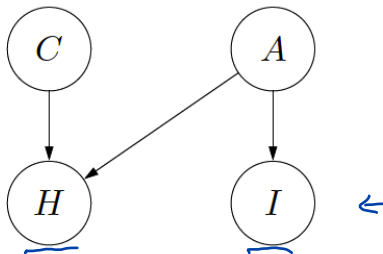


Suppose two causes positively influence an effect. Conditioned on the effect, conditioning on one cause reduces the probability of the other cause.



Example: cold or allergies

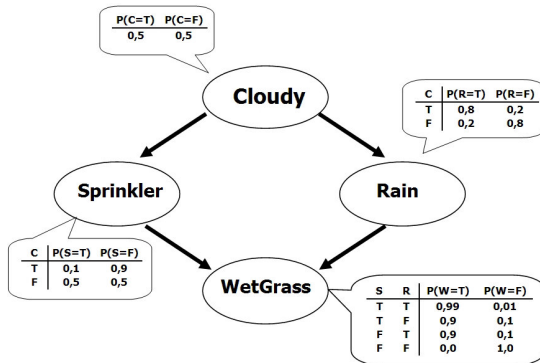
You are coughing and have itchy eyes. What do you have?



Variables: **C**old, **A**llergy, **Cough**, **I**chy eyes



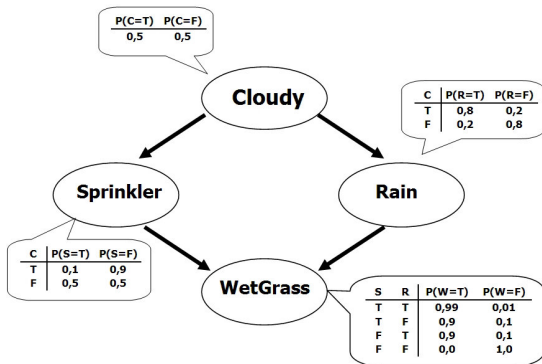
... or how can we use these models to efficiently answer probabilistic queries?



$$\underline{P(C, S, R, W)} = \underline{P(C)P(S|C)P(R|C)P(W|S, R)}$$

Inference: example

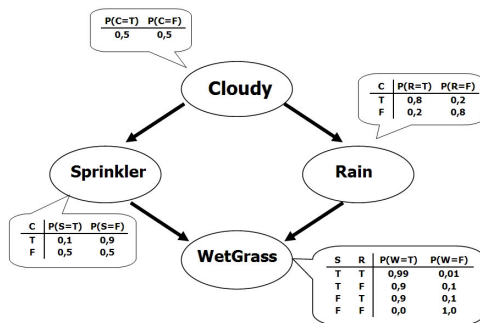
Let us use 0 to denote false and 1 to denote true.



What is the marginal probability, $P(S=1)$, that the sprinkler is on?

Inference: example

What is the marginal probability, $P(S=1)$, that the sprinkler is on?



$$\begin{aligned} P(S=1) &= \sum_C \sum_R \sum_W P(\underline{C}, \underline{R}, \underline{W}, \underline{S=1}) \\ &= \sum_C \sum_R \sum_W \underline{P(C)} \underline{P(S=1|C)} \underline{P(R|C)} \underline{P(W|S=1, R)} \end{aligned}$$

What is the marginal probability, $P(S=1)$, that the sprinkler is on?

$$\begin{aligned}
 P(S = 1) &= \sum_C \sum_R \sum_W P(W|S = 1, R)P(S = 1|C)P(R|C)P(C) \\
 &= P(\underline{W = 0}|S = 1, \underline{R = 0})P(\underline{S = 1}|\underline{C = 0})P(\underline{R = 0}|\underline{C = 0})P(\underline{C = 0}) \quad \begin{matrix} C & R & W \\ 0 & 0 & 0 \end{matrix} \\
 &+ P(W = 1|S = 1, R = 0)P(S = 1|C = 0)P(R = 0|C = 0)P(C = 0) \quad \begin{matrix} 0 & 0 & 1 \end{matrix} \\
 &+ P(W = 0|S = 1, R = 1)P(S = 1|C = 0)P(R = 1|C = 0)P(C = 0) \quad \begin{matrix} 0 & 1 & 0 \end{matrix} \\
 &+ \dots
 \end{aligned}$$

What is the marginal probability, $P(S=1)$, that the sprinkler is on?

Pseudocode

```
Prod = 0
For C=0:1:1
  For R=0:1:1
    For W=0:1:1
      Prod = Prod + P(C)P(R|C)P(S = 1|C)P(W|S = 1, R)
    end
  end
end
```

What is the marginal probability, $P(S=1)$, that the sprinkler is on?

$$\begin{aligned}
 P(S=1) &= \sum_C \sum_R \sum_W \underbrace{P(W|S=1, R)} P(S=1|C) P(R|C) P(C) \\
 &= \sum_C \sum_R P(S=1|C) P(R|C) P(C) \underbrace{\sum_W P(W|S=1, R)}_{1} \\
 &= \sum_C \sum_R P(S=1|C) P(R|C) P(C) \cdot 1 \\
 &= \sum_C P(S=1|C) P(C) \underbrace{\sum_R P(R|C)}_{1} \\
 &= \sum_C P(S=1|C) P(C) \cdot 1 \\
 &= \underbrace{P(S=1|C=0)P(C=0) + P(S=1|C=1)P(C=1)} = 0.3
 \end{aligned}$$

What is the marginal probability, $P(S=1)$, that the sprinkler is on?

Pseudocode

```
 $\Psi = 0$   
 $\Phi = 0$   
 $\Theta = 0$   
For  $W=0:1:1$   
     $\Phi = \Phi + P(W|S = 1, R)$   
end  
For  $R=0:1:1$   
     $\Psi = \Psi + P(R|C)\Phi$   
end  
For  $C=0:1:1$   
     $\Theta = \Theta + P(S = 1|C)\Psi$   
end
```


If we wish to compute several marginals at the same time, we can use **Dynamic Programming**: **avoid redundant computation**.

Exact algorithms work well for small graphs and for graphs that are trees or close to **trees** (have low tree-width). For large densely connected graphs we require the use of algorithms beyond the scope of this course:

- *Variational methods.*
- Sampling (Monte Carlo) methods: *Gibbs sampling*.

What is the posterior probability, $P(S = 1|W = 1)$, that the sprinkler is on given that the grass is wet?

What is the posterior probability, $P(S = 1|W = 1)$, that the sprinkler is on given that the grass is wet?

$$P(\underline{S = 1} | \underline{W = 1}) = \frac{\underline{P(S = 1, W = 1)}}{\underline{P(W = 1)}}$$

$$\underbrace{P(W = 1) = \sum_R \sum_C \sum_S P(W = 1, S, C, R)}$$

$$\underbrace{P(S = 1, W = 1) = \sum_R \sum_C P(W = 1, \underline{S = 1}, C, R)}$$

What is the posterior probability, $P(S = 1|W = 1, R = 1)$, that the sprinkler is on given that the grass is wet and it is raining?

What is the posterior probability, $P(S = 1|W = 1, R = 1)$, that the sprinkler is on given that the grass is wet and it is raining?

$$P(S = 1|W = 1, R = 1) = \frac{P(S = 1, W = 1, R = 1)}{P(W = 1, R = 1)}$$

$$P(W = 1, R = 1) = \sum_S \sum_C P(S, W = 1, R = 1, C)$$

$$P(S = 1, W = 1, R = 1) = \sum_C P(S = 1, W = 1, R = 1, C)$$



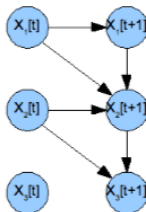
- Bayesian networks encode a joint distribution over the variables.
- We can calculate different probabilities by marginalising over the joint distribution.
- More sophisticated methods are needed for large graphs.

Temporal models

Dynamic Bayesian Networks (DBNs) are directed graphical models of stochastic processes

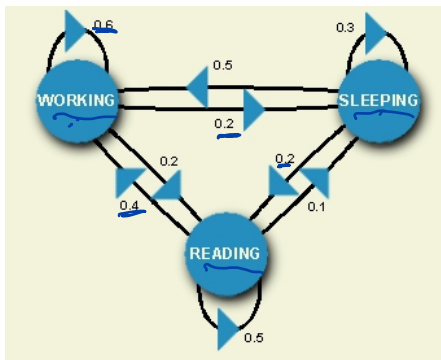
Dynamic Bayesian Networks...

- 1 Extend BNs to exploit temporal redundancy
 - Markov assumption: future independent of past given present
 - Stationarity: transition model is stable over time
- 2 State described via random variables
- 3 Each variable depends only on few other



Markov chain

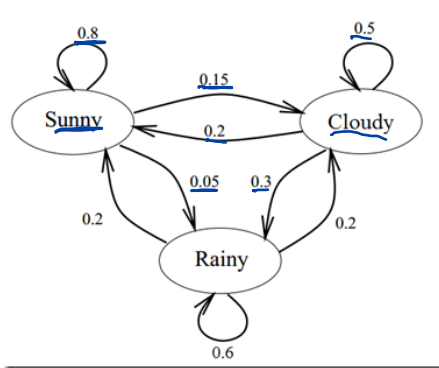
Markov chain: sequence of random variables each conditionally dependent only on the previous.



$$P(W, R, S) = P(W)P(R|W)P(S|R)$$

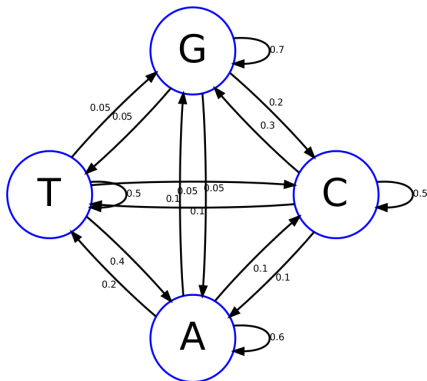
Markov chain

Markov chain: sequence of random variables each conditionally dependent only on the previous.



Markov chain

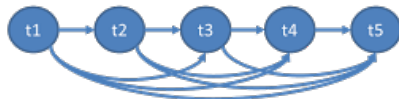
Markov chain: sequence of random variables each conditionally dependent only on the previous.



Markov chain

Fully joint distribution:

Probability distribution of a state depends on ALL preceding states



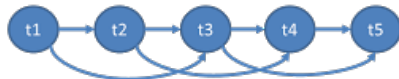
First order Markov Chain:

Probability distribution of a state depends ONLY on 1 immediate preceding state



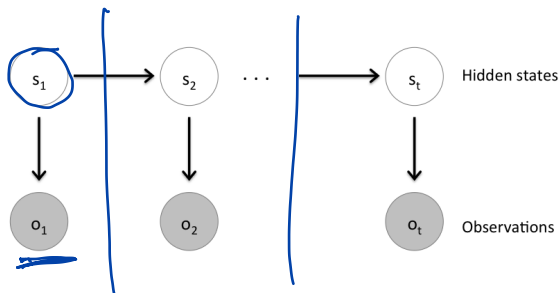
Second order Markov Chain:

Probability distribution of a state depends ONLY on 2 immediate preceding states



Hidden Markov Models (HMMs)

The simplest kind of DBN is a **Hidden Markov Model (HMM)**, which has one discrete hidden node and one discrete or continuous observed node per slice.

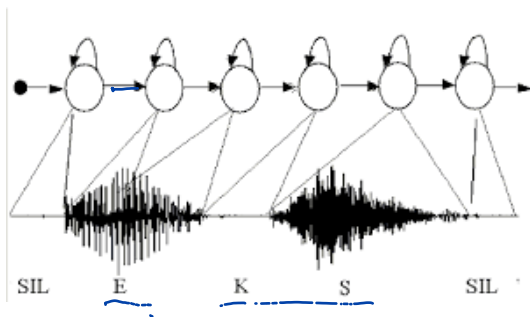


In general, we assume we have an **initial distribution** $P(s_1)$, a **transition model** $P(s_t | s_{t-1})$, and an **observation model** $P(o_t | s_t)$.

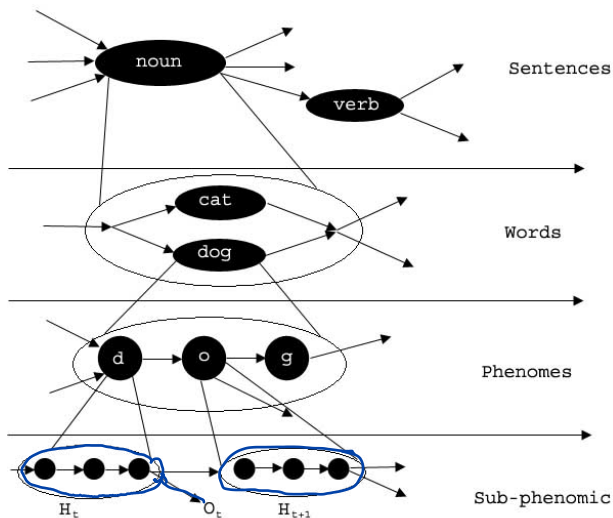
$$P(\text{words} \mid \text{sound}) \propto P(\text{sound} \mid \text{words}) P(\text{words})$$

Final beliefs Likelihood of data Prior language model
eg mixture of Gaussians eg unigrams

Hidden Markov Model (HMM)



HMM: language



HMM: tracking

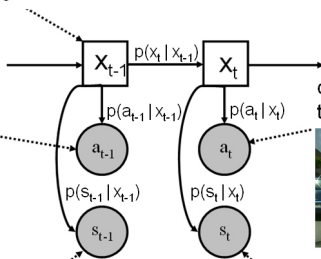
number of corresponding frame
in reference video to frame t-1 in
observed video

descriptor of
frame t-1 of
observed video



UTM GPS data at instant
t-1 of observed video

$\begin{bmatrix} x : 425804.05020800866 \\ y : 4594566.545587094 \end{bmatrix}$



descriptor of frame
t of obs. video



UTM GPS data at instant
t of observed video

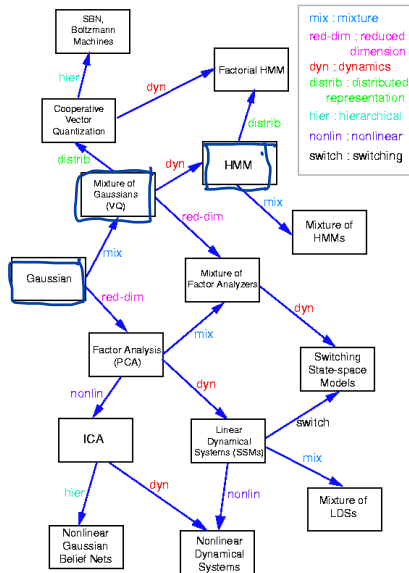
$\begin{bmatrix} x : 425804.05890800866 \\ y : 4594566.546087094 \end{bmatrix}$

Hidden
variables

Observed
variables

$y=(a,s)$

Relationships between graphical models



- Bayesian networks can be modified to include temporal information
- Markov chains include temporal information up to a given point
- Hidden Markov models allow us to infer the presence of hidden variables based on the output they produce