



# Uncertainty

Professor Marie Roch

12.2-12.6,  
13-13.3.3.1  
(except 13.2.1),  
Rabiner's HMM Tutorial



# Uncertainty

---

*toothache*  $\Rightarrow$  *cavity*

What else can cause a toothache?

*toothache*  $\Rightarrow$  *cavity*  $\vee$  *cracked tooth*  $\vee$  *stuck popcorn*  $\vee$  ...

Logic can fail us:

- laziness – Too difficult to enumerate rules without exceptions
- theoretical ignorance – May not fully understand the system
- practical ignorance – System may not be fully observable



# An agent's view



Probabilities represent a level of belief in a world

# Decision-theoretic agents

**function** DT-AGENT(*percept*) **returns** an *action*

**persistent:** *belief\_state*, probabilistic beliefs about the current state of the world  
*action*, the agent's action

update *belief\_state* based on *action* and *percept*

calculate outcome probabilities for actions,

    given action descriptions and current *belief\_state*

select *action* with highest expected utility

    given probabilities of outcomes and utility information

**return** *action*



# Basic probability

- Random variables represent an outcome or world, e.g.  $X$  represents a die roll,  $W_{2,1}$  is the Wumpus in cave 2,1
- $P(X)$  is the probability of  $X$  happening
- It is common to use
  - CAPITALS to represent outcomes in general:  $P(X)$
  - lower case to denote specific outcomes:  $P(x=5)$
- Probability distributions characterize probability over all outcomes and require:

$$\forall x \in \text{domain}(X), 0 \leq P(X = x) \leq 1$$

$$\sum_{x \in \text{domain}(X)} P(X = x) = 1$$



# Posterior (conditional) probability

- Posterior probability is conditioned on another event

What is the probability that I have a cavity *given* that I have a toothache?

$$P(\text{cavity}|\text{toothache})$$

- In contrast, prior probabilities have no condition

$$P(\text{cavity})$$

- Definition:  $P(A|B) = \frac{P(A \wedge B)}{P(B)}$  or equivalently:  $P(A \wedge B) = P(A|B)P(B)$   
product rule

note:  $P(A \wedge B)$  is frequently written as  $P(A, B)$



# Propositions

- Let us consider random variable values as possible worlds (like our model checking in propositional logic)
- If we want to know P of proposition  $\phi$  holding:

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

- In addition

$$P(\neg\phi) = \sum_{\omega \in \neg\phi} P(\omega) = 1 - \sum_{\omega \in \phi} P(\omega)$$

and

$$P(\phi \vee \rho) = P(\phi) + P(\rho) - P(\phi \wedge \rho)$$

inclusion-exclusion principle



# Joint probabilities

- Probability of multiple things, e.g.  $P(A, B, C)$ .

- Can be decomposed with the product rule:

$$P(A, B, C) = P((A, B), C)$$

$$= P(C|A, B)P(A, B)$$

$$= P(C|A, B)P(B|A)P(A)$$

This is called the chain rule.

In general:  $P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i|x_{i-1}, x_{i-2}, \dots, x_1)$

- If A and B are independent, then  $P(B|A) = P(B)$
- In general, joint probabilities of independent variables can be multiplied:  $P(A, B, C) = P(A)P(B)P(C)$  (A/B/C independent)





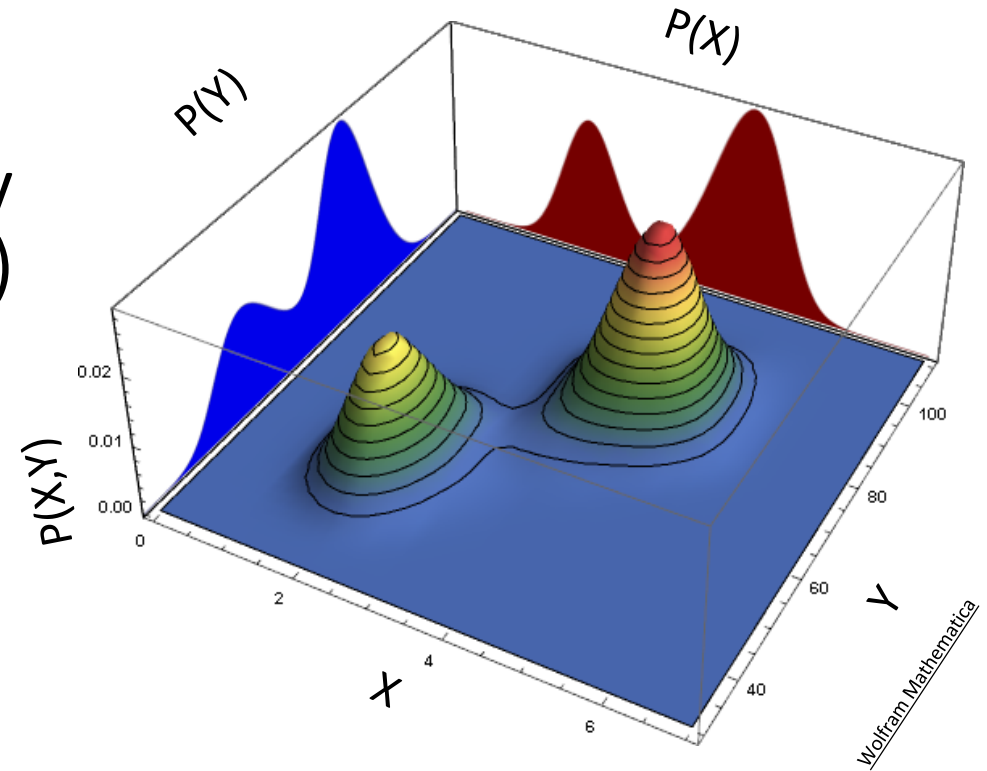
# Marginalization

- Suppose we know the joint probability between  $X$  and  $Y$ ,  $P(X,Y)$ , and want  $P(X)$

$$P(X) = \sum_y P(X, Y = y)$$

example:

$$P(\text{Eat}, \text{Rent}) = \sum_{r=0}^{\text{owed}} P(\text{Eat}, \text{Rent} = r)$$



# Bayes' rule

(of conditional probability)

Remember definition posterior probability  $P(A|B) = \frac{P(A,B)}{P(B)}$

$$P(B|A) = \frac{P(B,A)}{P(A)} \rightarrow P(B,A) = P(B|A)P(A)$$

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

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

also known as Bayes' Theorem/Bayes' Law)



# Bayes' rule: Why do we care?

Suppose we observe an effect.

- Knowing the cause can be difficult
- Simpler to estimate  $P(\text{effect} \mid \text{cause})$ ; Bayes' rule lets us turn this around:

$$P(Disease|Symptom) = \frac{P(Symptom|Disease)P(Disease)}{P(Symptom)}$$

If we are looking at multiple diseases, we do not need  $P(\text{Symptom})$  to make a choice between them.

We can treat  $1/P(\text{Symptom})$  as a constant  $\alpha$ :

$$P(Disease|Symptom) = \alpha P(Symptom|Disease)P(Disease)$$



# Bayes' rule example

- A symptom of meningitis is a stiff neck

$$P(s|m) = 0.7$$

but the case rate for stiff necks is low and meningitis very low

$$P(s) = 0.01, P(m) = 1/50000$$

$$P(m|s) = \frac{P(s|m)P(m)}{P(s)} = \frac{0.7 \times 1/50000}{0.01} = 0.0014$$

# Conditional independence

What do you see?

What does your neighbor see?

# Conditional independence

- Both Shyam and Monica observe the cloud
- If they haven't talked to each other about what they saw, the probability conditioned on a specific cloud is independent

$$P(S, M|c) = P(S|c)P(M|c)$$



# Naïve Bayes models

- Exploit conditional independence to make simple models
- If a cause has  $n$  effects that are conditionally independent:

$$P(\textit{cause}, e_1, e_2, \dots, e_n) = P(\textit{cause}) \prod_{i=1}^n P(e_i | \textit{cause})$$

- Naïve Bayes models imply that we don't really know if our effects are conditionally independent, but we assume so anyway.
- If we wanted to find  $P(\textit{cause} | e_1, e_2, \dots, e_n)$  we could use the product rule and conditional independence:

$$P(\textit{cause} | e_1, e_2, \dots, e_n) = \alpha P(\textit{cause}) \prod_{i=1}^n P(e_i | \textit{cause})$$

# Example: Sentence to category

*Disneyland raised its entrance price by thirty percent.*

We might ask the question: Is this about business or entertainment?

We could consider how often articles are about each of these categories (prior):

$$P(\text{business}) = .03$$

$$P(\text{entertainment}) = .04$$



# Example: Sentence to category

Bayes factors:

$$P(\text{Disneyland} | \text{business}) = .2$$

$$P(\text{Disneyland} | \text{entertainment}) = .8$$

$$P(\text{price} | \text{business}) = .9$$

$$P(\text{price} | \text{entertainment}) = .1$$

Prior probabilities:

$$P(\text{business}) = .03$$

$$P(\text{entertainment}) = .04$$

$$\begin{aligned} P(bz | Dland, \$) \\ &= P(Dland | bz)P(\$ | bz) \\ &= .03 \cdot .2 \cdot .9 = .0054 \end{aligned}$$

$$\begin{aligned} P(ent | Dland, \$) \\ &= P(ent)P(Dland | ent)P(\$ | ent) \\ &= .04 \cdot .8 \cdot .1 = .0032 \end{aligned}$$

We classify as the sentence as the category that maximizes P



# Probabilistic reasoning

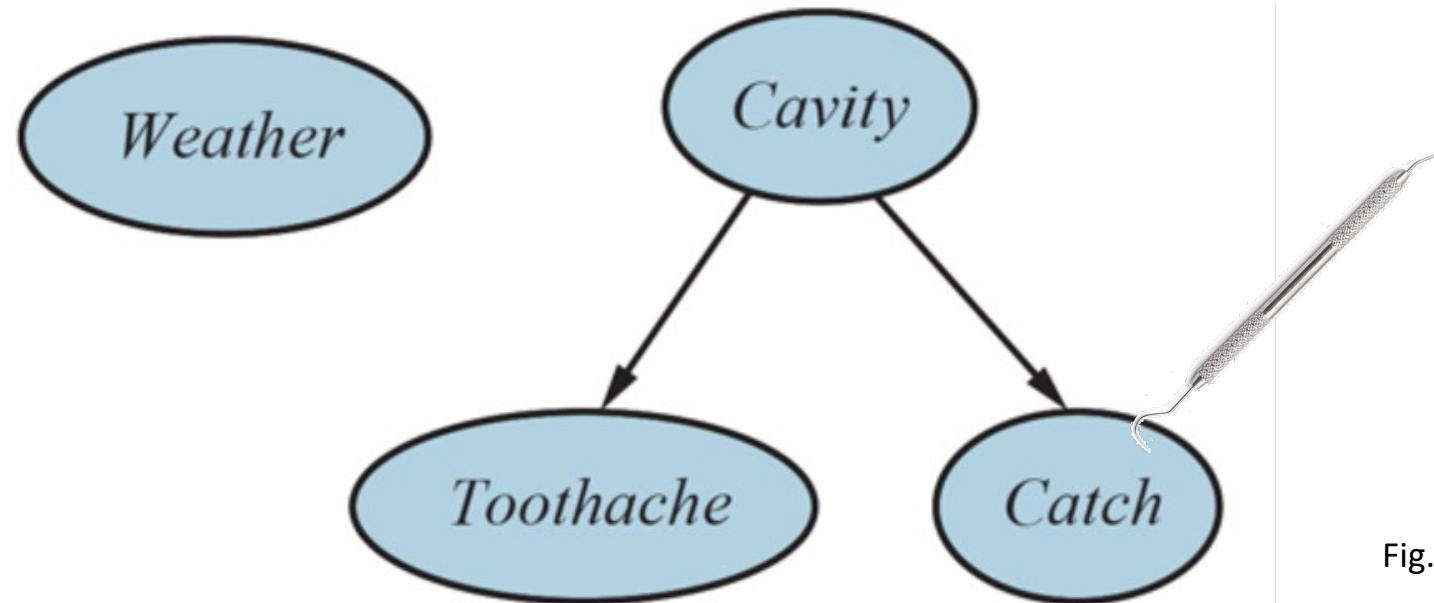


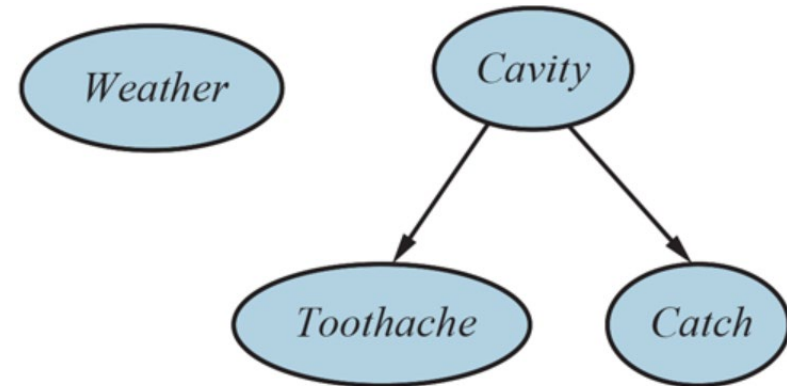
Fig. 13.1

## Bayesian network

- Having a cavity influences the likelihood of a toothache or a dentist's sickle probe to catch on your tooth
- Changes in weather do not cause toothaches or probe catches.

# Bayesian networks

- Nodes are random variables
- Variables can be connected by directed arcs that do not form cycles
- Each variable  $V$  has
  - prior probability (no parents):  $P(V)$
  - conditional probability  $P(V|\text{parents}(V))$
- Forms a directed acyclic graph



# Pearl's Bayesian network example

---

- Burglar alarm set off by
  - Burglar
  - Earthquake
- Neighbors Mary and John have agreed to let you know when they hear the alarm
  - Mary listens to headphones, and often misses the alarm
  - Your home telephone ringtone is similar to the alarm (silly you) and John sometimes calls you when your phone rings (yes, you still have a landline)



# Pearl's Bayesian network example

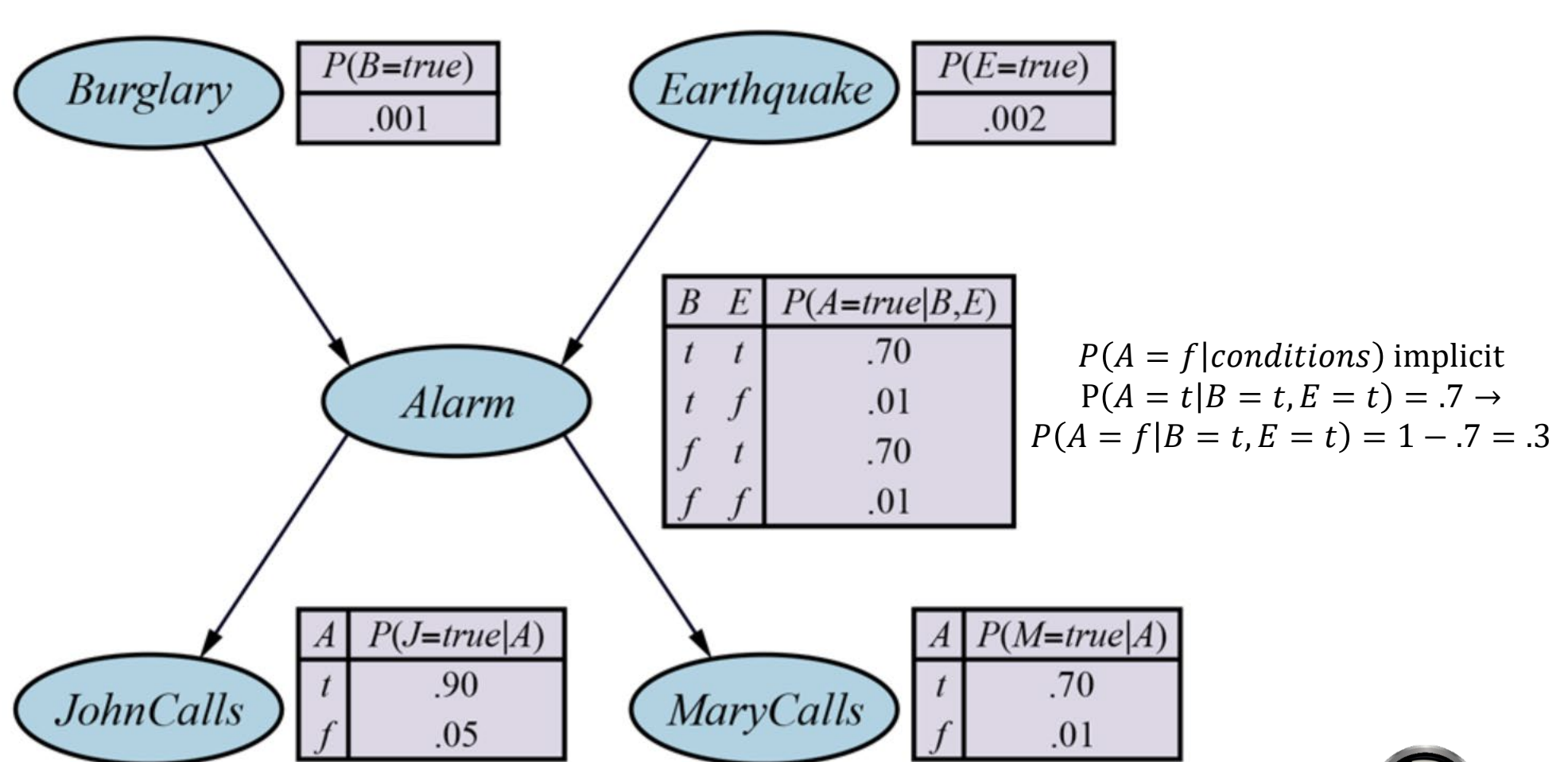


Fig. 13.2 R&N

# Bayes net semantics

- $P(x_1, x_2, \dots, x_n) = \sum_{i=1}^n P(x_i | \text{parents}(X_i))$

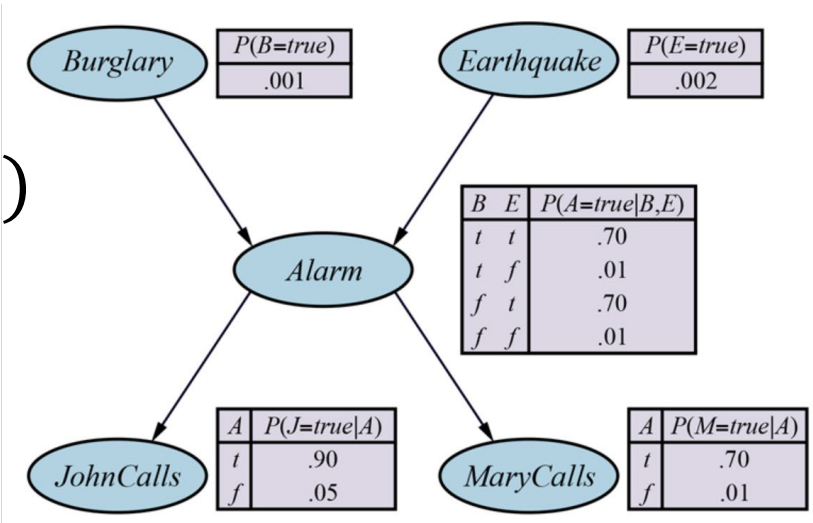


Fig. 13.2 R&N

- Consider: alarm sounds with neither burglary/earthquake and both neighbors call

$$\begin{aligned}
 P(j, m, a, \neg b, \neg e) &= P(j|a)P(m|a)P(a|\neg b \wedge \neg e)P(\neg b)P(\neg e) \\
 &= .9 \times .7 \times .01 \times .999 \times .998 = .00628
 \end{aligned}$$

# Bayes net semantics

- We can compute the marginal to answer just about any question related to this, e.g. John & Mary call when the alarm sounds and there is no burglary

$$P(j, m, a, \neg b) = \sum_{E \in e, \neg e} P(j|a)P(m|a)P(a|\neg b \wedge E)P(\neg b)P(E)$$

- Note that earthquake was not specified in the question; we computed the marginal probability to integrate/sum it out.

# Constructing a Bayes network

- Nodes

- Determine required random variables
- Number them  $X_1, X_2, \dots, X_n$  (better if causes precede effects)

- Network edges

for  $i = 1:n$

Find minimum  $parents(X_i)$  such that  $P(X_i|X_{i-1}, \dots, X_1) = P(X_i|parents(X_i))$

Add edges  $parents(X_i)$  to  $X_i$

Estimate conditional probability table  $P(X_i|parents(X_i))$

Note: We are only concerned about direct influence, so Alarm influences MaryCalls, but Burglary and Earthquake do not.



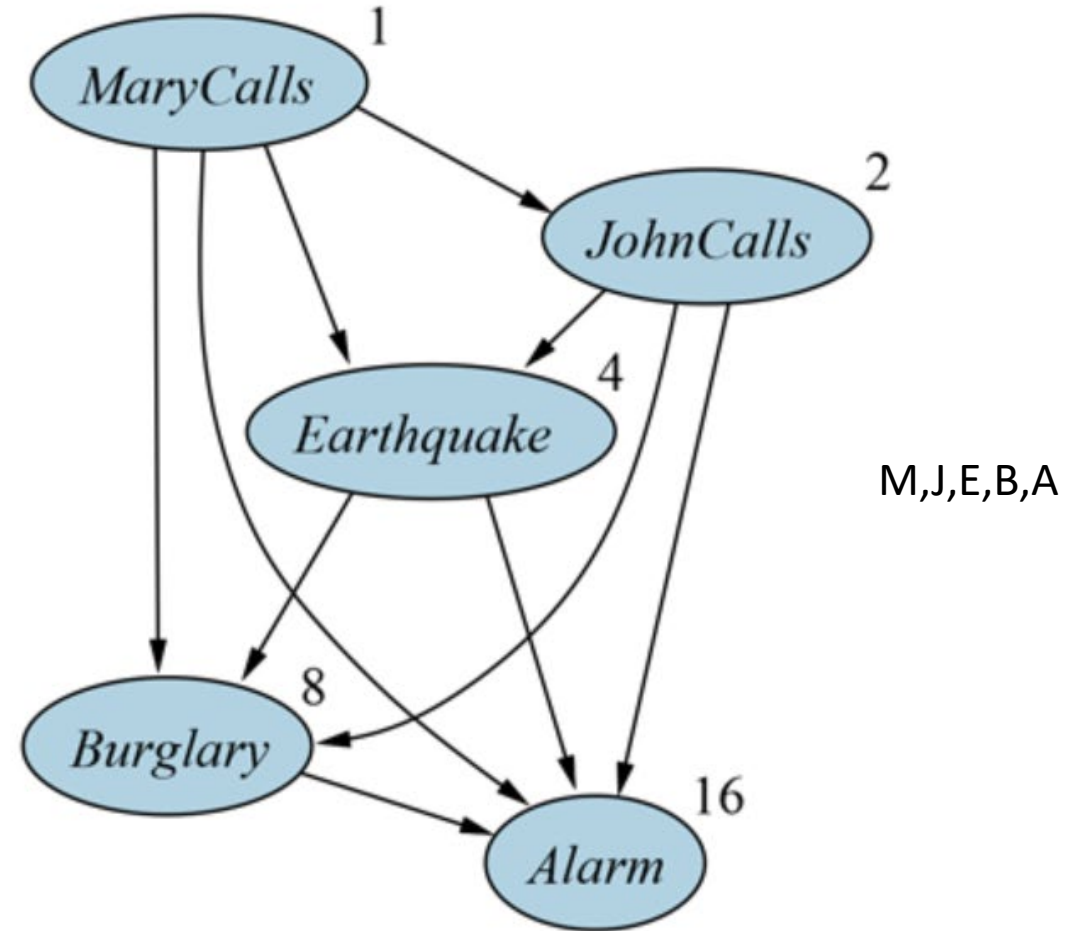
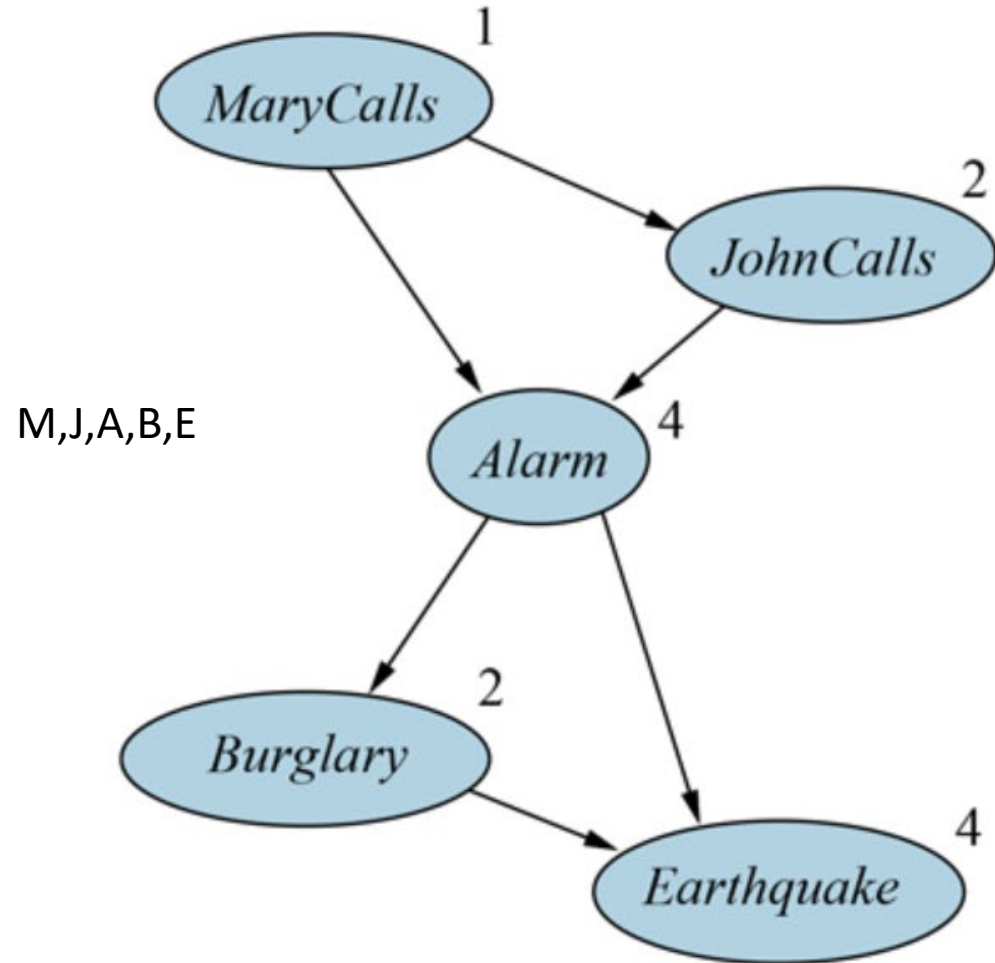
# Estimating the conditional probability table

- Estimated from training data
- Discrete – Use a frequentist model, e.g.,  $P(A = a|b) = \frac{\text{count}(a,b)}{\text{count}(b)}$
- Continuous - Fit distribution to data, e.g.  $P(A = a|b) \sim n(\mu, \sigma^2)$ , use the mean and variance of examples where  $B=b$
- Chapter 20 has more details on learning in probabilistic models. (See also any basic statistics book's chapter on maximum likelihood estimation)

# Variable order in Bayes network construction

- Construction depends on order
- Consider order: MaryCalls, JohnCalls, Alarm, Burglary, Earthquake
  - MC – no parents
  - JC – If MC, probably an earthquake, hence  $P(JC|MC)$
  - A – Alarm more likely if both MC and JC call, therefore they are parents of A.
  - B – If we know alarm state, MC & JC do not give us additional information about whether this was a burglary or earthquake, hence  $P(B|A)$   
[We assume that this is a minor earthquake, an example of laziness in uncertainty.]
  - E – If A, then it is more likely that there was an earthquake, but B would also cause an alarm and knowing this reduces the probability of E. Hence  $P(E|A,B)$

# Order of variables in construction matters



In general, better to order variables in what we think might be a causal manner. However, both networks will learn appropriate distributions.

# Efficient representations

- For a binary Bayes net with at most  $k$  parents, conditional probability tables (CPT) have  $O(2^k)$  entries.
- Many times, relationships fit into patterns that we call **canonical distributions**, and can specify the conditional probability tables with the canonical name and a few parameters.

# Canonical distribution examples

- deterministic nodes – are not probabilistic but can be represented by a function, e.g. *ReservoirLevelChange* might be the sum of inputs from rivers – evaporation
- context-specific independence – Parents might be independent when other parents have specific values

Example:  $P(\text{Damage}|\text{Ruggedness},\text{Accident}) = d_1$  if  $\text{Accident} == \text{false}$   
else  $d_2(\text{Ruggedness})$

$$P_{d_1}(\text{Damage}) = \begin{cases} .995 & \text{Damage} = \text{false} \\ .005 & \text{Damage} = \text{true} \end{cases} \quad \leftarrow \text{indep. of ruggedness}$$

$$P_{d_2}(\text{Damage}|\text{Accident}) = \begin{cases} f_{\text{Ruggedness}}(.20) & \text{Damage} = \text{false} \\ f_{\text{Ruggedness}}(.80) & \text{Damage} = \text{true} \end{cases}$$

# Canonical distribution examples

- noisy-or – Permits uncertainty in causation  
e.g., In propositional logic we might state:  $Fever \Leftrightarrow Cold \vee Flu \vee Malaria$   
If you have one of these, you have a fever.
- Suppose disease  $i$  occurs without fever with frequency  $q_i$ :

$$\begin{aligned}q_{cold} &= P(\neg fever | cold, \neg flu, \neg malaria) = 0.6 \\q_{flu} &= P(\neg fever | \neg cold, flu, \neg malaria) = 0.2 \\q_{malaria} &= P(\neg fever | \neg cold, \neg flu, malaria) = 0.1\end{aligned}$$

Noisy-or would make fever true as follows:

$$P(fever | parents(fever)) = 1 - \prod_{\substack{j: j=true \wedge \\ j \in parents(fever)}} q_j$$

We can think of this as 1 – the joint probability that everything you have that is making you sick did not cause a fever.

# Noisy-or example

$$q_{\text{cold}} = P(\neg \text{fever} | \text{cold}, \neg \text{flu}, \neg \text{malaria}) = 0.6,$$

$$q_{\text{flu}} = P(\neg \text{fever} | \neg \text{cold}, \text{flu}, \neg \text{malaria}) = 0.2,$$

$$q_{\text{malaria}} = P(\neg \text{fever} | \neg \text{cold}, \neg \text{flu}, \text{malaria}) = 0.1.$$

<i>Cold</i>	<i>Flu</i>	<i>Malaria</i>	$P(\text{fever}   \cdot)$	$P(\neg \text{fever}   \cdot)$	Fig. 13.5
<i>f</i>	<i>f</i>	<i>f</i>	0.0	1.0	
<i>f</i>	<i>f</i>	<i>t</i>	0.9	<b>0.1</b>	
<i>f</i>	<i>t</i>	<i>f</i>	0.8	<b>0.2</b>	
<i>f</i>	<i>t</i>	<i>t</i>	0.98	$0.02 = 0.2 \times 0.1$	
<i>t</i>	<i>f</i>	<i>f</i>	0.4	<b>0.6</b>	
<i>t</i>	<i>f</i>	<i>t</i>	0.94	$0.06 = 0.6 \times 0.1$	
<i>t</i>	<i>t</i>	<i>f</i>	0.88	$0.12 = 0.6 \times 0.2$	
<i>t</i>	<i>t</i>	<i>t</i>	0.988	$0.012 = 0.6 \times 0.2 \times 0.1$	

# Bayesian nets with continuous variables

- Several options
  - Discretize – split into discrete values based on range
  - Use a parametric distribution, e.g., normal distribution
  - Non-parametric options possible, but beyond our scope
- Linear-Gaussian conditional distribution
  - Most common parametric distribution
  - Variance fixed, mean dependent on a continuous parent



# Hybrid Bayesian nets

- Contain both discrete and continuous variables

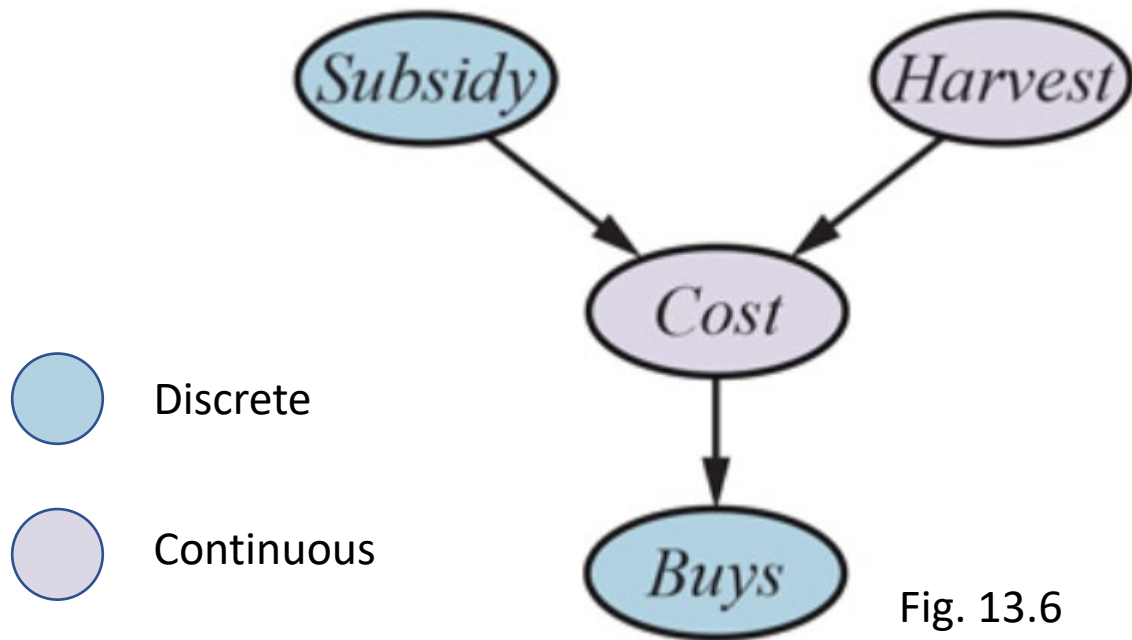


Fig. 13.6

Whether or not a consumer purchases fruit depends on its cost.

The cost depends on the harvest and whether or not a government subsidy was provided.

# Linear-Gaussian example

$$P(c|h, \text{subsidy}) = N(c; a_t h + b_t, \sigma_t^2) = \frac{1}{\sigma_t \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{c - (a_t h + b_t)}{\sigma_t} \right)^2}$$

$$P(c|h, \neg \text{subsidy}) = N(c; a_f h + b_f, \sigma_f^2) = \frac{1}{\sigma_f \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{c - (a_f h + b_f)}{\sigma_f} \right)^2}$$

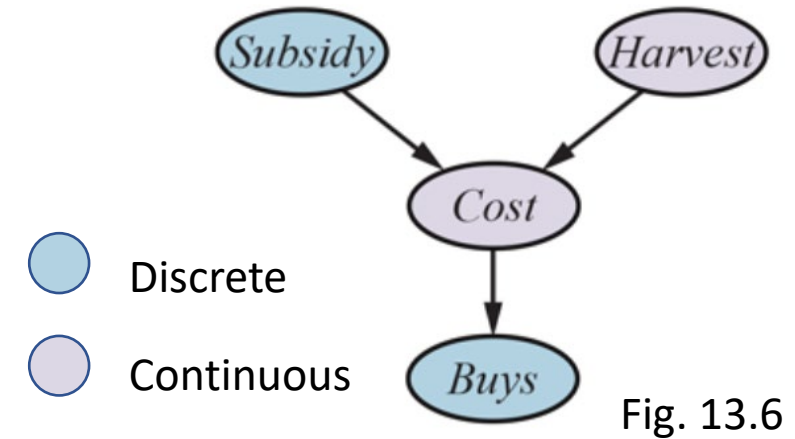


Fig. 13.6

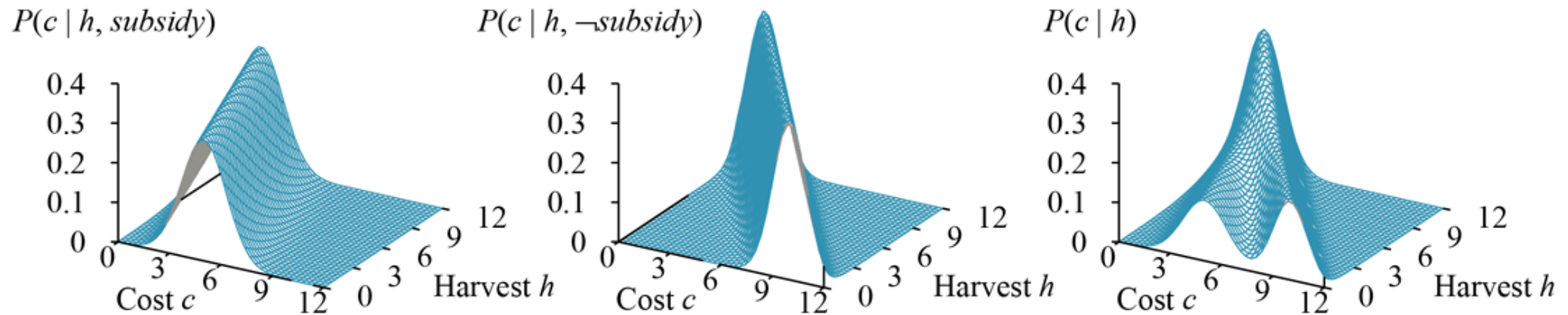


Fig. 13.7

# Discrete with continuous parents

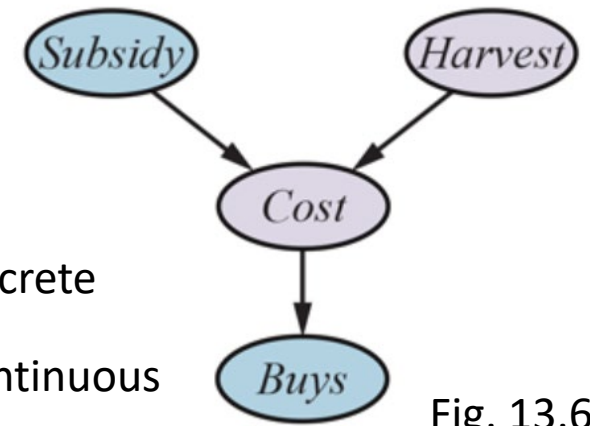


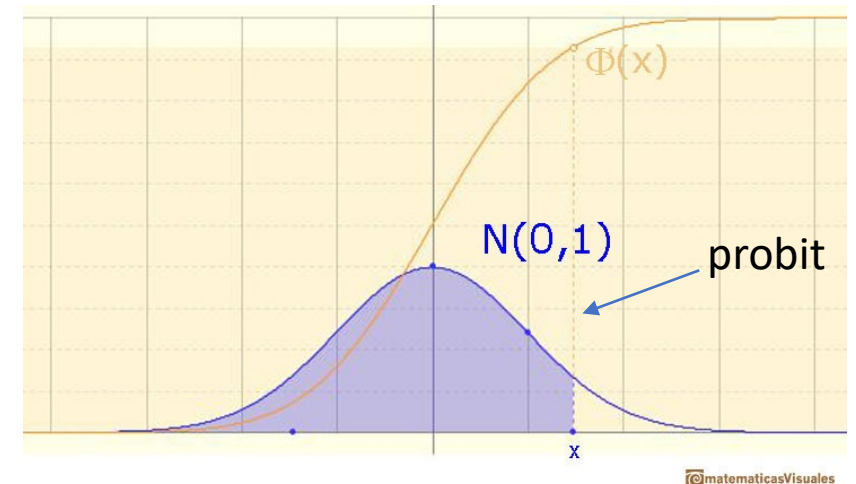
Fig. 13.6

- We need some type of “soft” threshold
- Remember the cumulative density function we introduced with the  $\chi^2$  distribution. In general,  $P(X \leq x)$  (sometimes denoted  $\Phi(x)$ ):

$$P(X \leq x) = \int_{-\infty}^x P(x)dx$$

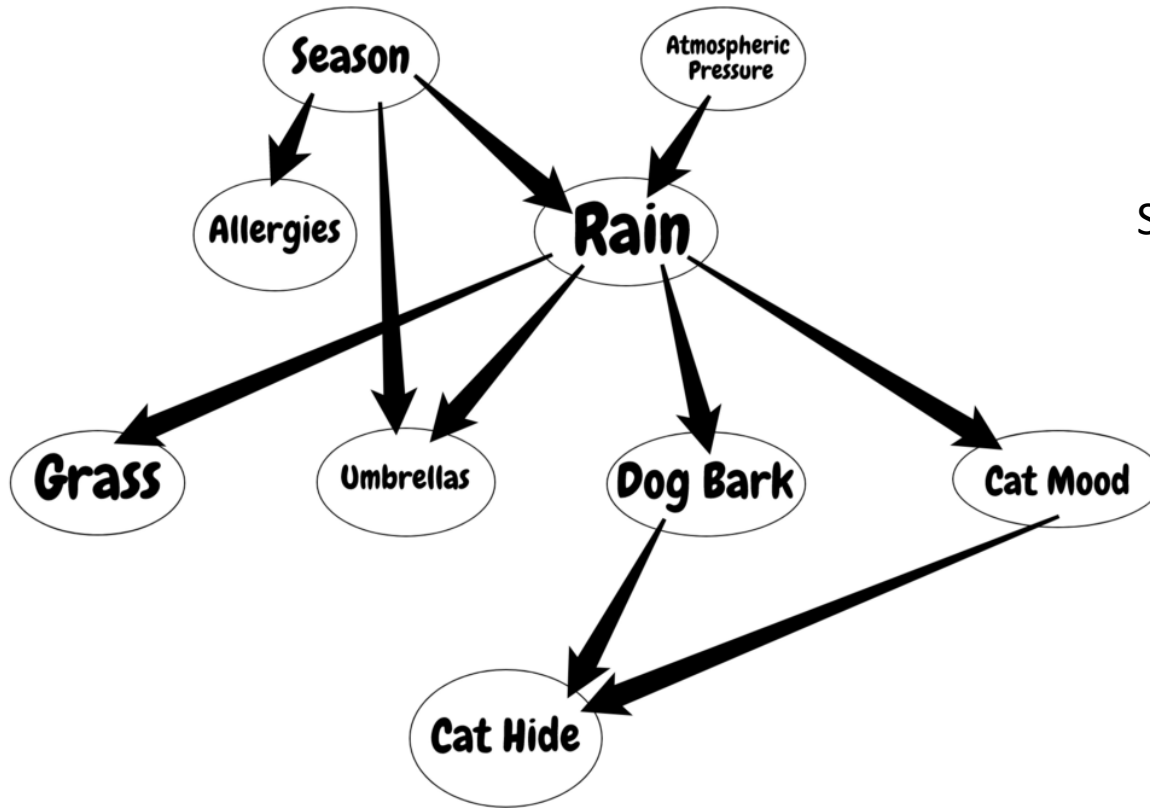
- However, this varies smoothly from 0 to 1 as X increases, which is not exactly what we want
- Invert the probability unit (probit) model:

$$P(B|C = c) = 1 - P\left(C \leq \frac{c - \mu_c}{\sigma_c}\right)$$



Similar model is the inverse logistic (logit) function:  $P(B|C = c) = 1 - \frac{1}{1 + e^{\frac{s \cdot (c - \mu_c)}{\sigma_c}}}$  where  $s$  is the probit's mean

# Another example



See 13.2.4 for a more complex case study

Example by Atakan Güney, [towardsdatascience.com](https://towardsdatascience.com)

# Evaluating probability

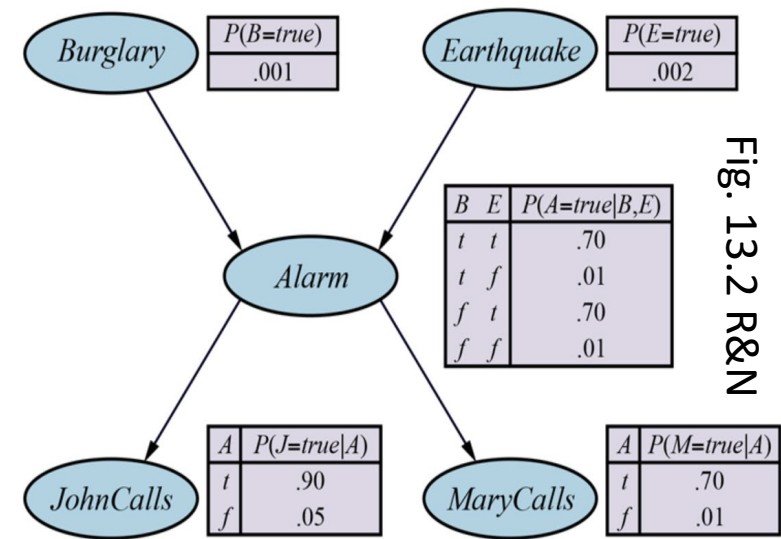
- Let  $X$  represent what we want to know
- Let  $e$  represent one or more evidence values (things we measured)
- Recall that  $P(X|e) = \frac{P(X,e)}{P(e)} = \alpha P(X, e)$  where  $\alpha = 1/P(e)$
- Let  $y$  be variables that are *latent* (hidden or unobservable) are denoted
- Then:

$$P(X|e) = \alpha P(X, e) = \alpha \sum_y P(X, e, y)$$

Note:  $\sum_y$  sums over all combinations of the  $y$  latent variables

# Evaluating probability

- Let us consider the burglar alarm example
- Suppose we want to query:  
 $P(\text{Burglary} | \text{JohnCalls} = \text{True}, \text{MaryCalls} = \text{True})$



$$P(b|j, m) = \alpha \sum_e \sum_a P(b)P(e)P(a|b, e)P(j|a)P(m|a)$$

naïve complexity for Booleans  $O(n2^n)$

- the first two terms do not depend on a. Hence:

$$P(b|j, m) = \alpha \sum_e P(b)P(e) \sum_a P(a|b, e)P(j|a)P(m|a)$$

complexity for Booleans  $O(2^n)$

# Evaluating probability

**function** ENUMERATION-ASK( $X, \mathbf{e}, bn$ ) **returns** a distribution over  $X$

**inputs:**  $X$ , the query variable

$\mathbf{e}$ , observed values for variables  $\mathbf{E}$

$bn$ , a Bayes net with variables  $vars$

$\mathbf{Q}(X) \leftarrow$  a distribution over  $X$ , initially empty

**for each** value  $x_i$  of  $X$  **do**

$\mathbf{Q}(x_i) \leftarrow$  ENUMERATE-ALL( $vars, \mathbf{e}_{x_i}$ )

where  $\mathbf{e}_{x_i}$  is  $\mathbf{e}$  extended with  $X = x_i$

**return** NORMALIZE( $\mathbf{Q}(X)$ )

**function** ENUMERATE-ALL( $vars, \mathbf{e}$ ) **returns** a real number

**if** EMPTY?( $vars$ ) **then return** 1.0

$V \leftarrow$  FIRST( $vars$ )

**if**  $V$  is an evidence variable with value  $v$  in  $\mathbf{e}$

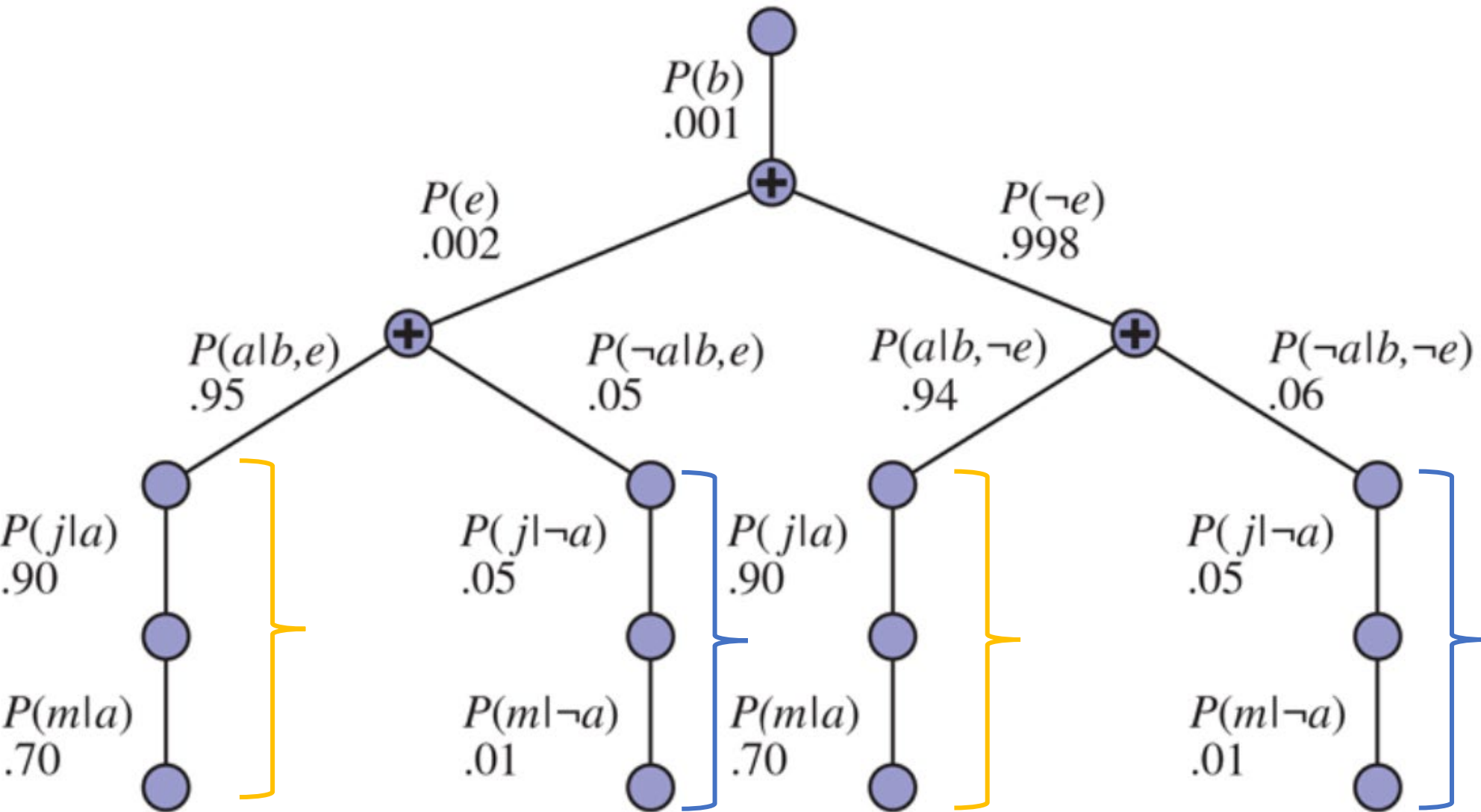
**then return**  $P(v | parents(V)) \times$  ENUMERATE-ALL(REST( $vars$ ),  $\mathbf{e}$ )

**else return**  $\sum_v P(v | parents(V)) \times$  ENUMERATE-ALL(REST( $vars$ ),  $\mathbf{e}_v$ )

where  $\mathbf{e}_v$  is  $\mathbf{e}$  extended with  $V = v$

Fig. 13.11 R&N

# Redundancies



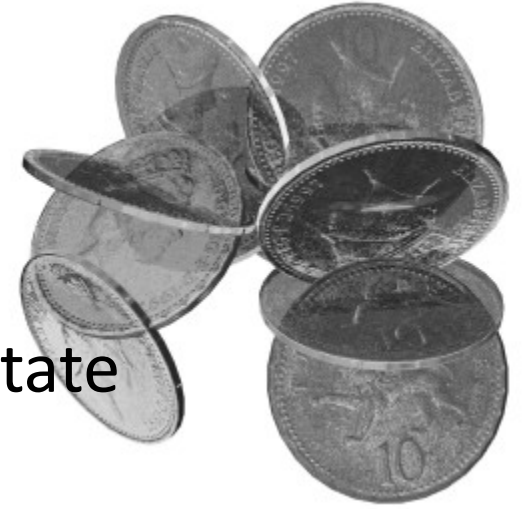


# More efficient computation

- Evaluating common subgraphs once is more efficient and makes a difference in large graphs. Variable elimination algorithm (13.3.2) does this
- There are also approximate evaluation algorithms that are covered later in chapter 13

You are not responsible for these.

# hidden Markov models



- Used for modeling processes that have an unobservable state
- Example
  - 2 coins behind a screen with different odds of heads/tails  
Here we'll let one coin be fair, the other is biased
  - I have a process
    - Flip coin
    - Choose the next coin to flip
  - All we observe are sequences: H, H, H, H, T, H, T, T, H, ...
- hidden Markov models let us model these types of systems

# Markov property

- Let  $q_i$  be the state that we are in at time  $i$ . If  $q_4 = \text{fair}$ , we are using the fair coin for the 4<sup>th</sup> flip in our previous example.
- Chain rule states

$$P(q_1, q_2, \dots, q_T) = P(q_1) \prod_{i=2}^T P(q_i | q_{i-1} q_{i-2} \dots q_1)$$

- Markov property specifies conditional independence after 1 step (can be generalized to N steps)

$$P(q_i | q_{i-1} q_{i-2} \dots q_1) = P(q_i | q_{i-1})$$

# Observed Markov models

- Finite state machine with state transition probabilities:

$$A = \begin{bmatrix} a_{0,0} & a_{0,1} & a_{0,2} & a_{0,3} & a_{0,4} \\ a_{1,0} & a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} \\ a_{2,0} & a_{2,1} & a_{2,2} & a_{2,3} & a_{2,4} \\ a_{3,0} & a_{3,1} & a_{3,2} & a_{3,3} & a_{3,4} \\ a_{4,0} & a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} \end{bmatrix}$$

$$\sum_{j=0}^4 a_{ij} = 1$$

where  $0 \leq i \leq 4$

$a_{2,3}$  is the probability of moving from state 2 to state 3.

So the probability of going from state 2 at time 5 to state 3 at time 6 would be:

$$P(q_6 = 3 | q_5 = 2) = a_{2,3}$$

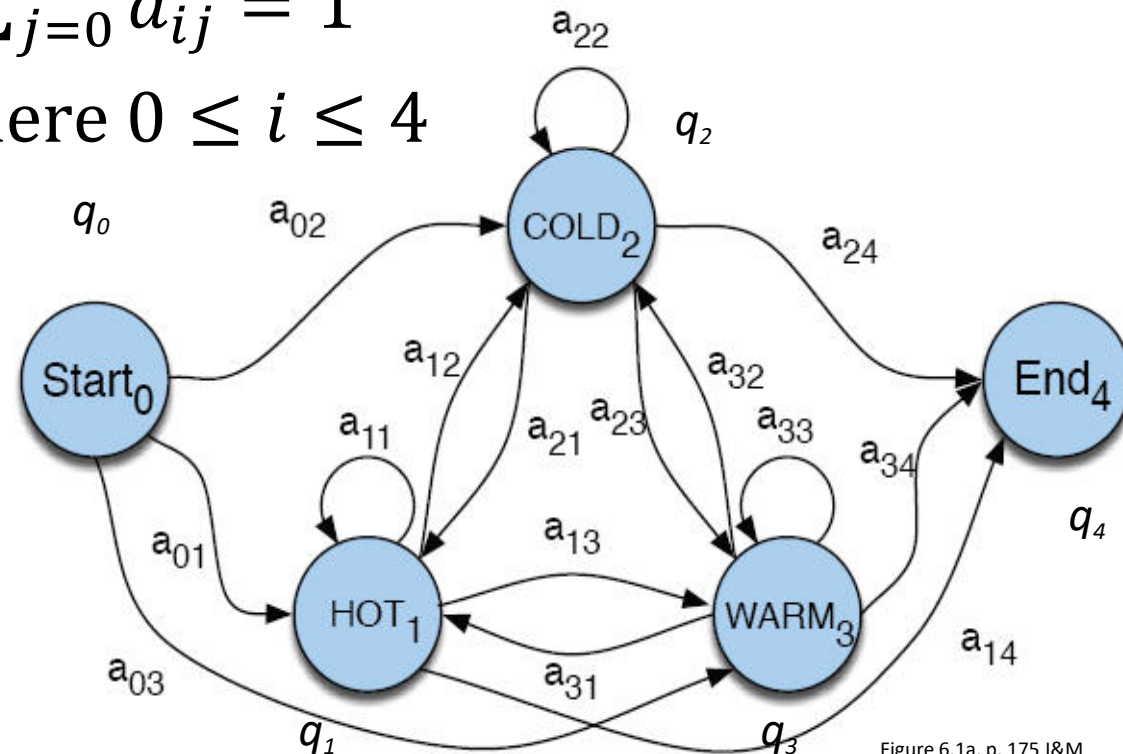


Figure 6.1a, p. 175 J&M

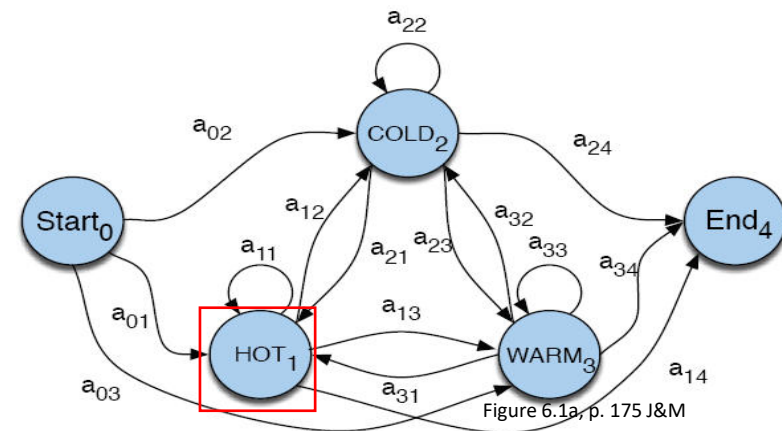
# Observed Markov models

- Markov for a state sequence:

$$P(q_i | q_1, \dots, q_{i-1}) = P(q_i | q_{i-1})$$

e.g.,  $P(\text{warm}_3 | \text{hot}_1) = a_{13}$

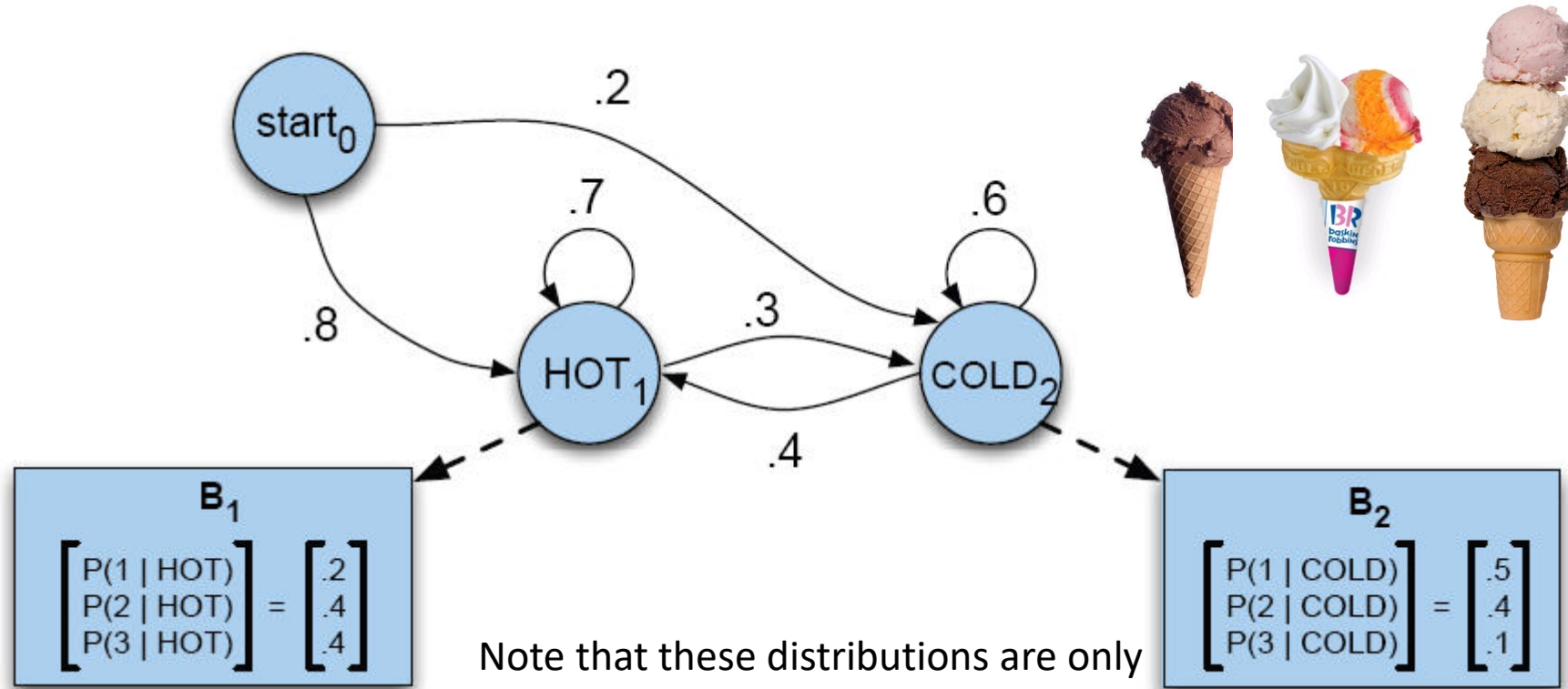
- Chain rule & Markov property



$$P(q_1, q_2, \dots, q_T) = P(q_1) \prod_{i=2}^T P(q_i | q_{i-1})$$

# State dependent distributions

Number of scoops of ice cream eaten daily



Note that these distributions are only dependent upon the state.



# State-dependent transitions

- R&N present observation probabilities differently.
  - Each observation probability is written as an  $S \times S$  *matrix*
- We will use the notation of Rabiner's 1989 tutorial article in *Proc IEEE*

# Notation worth remembering

- Observations (features)

$$O = \{o_1, o_2, \dots, o_T\}$$

- Observations come from a discrete or continuous space and are independent of one another

$$o_i \in \mathbb{R}^D \text{ or } o_i \in \{v_1, v_2, \dots, v_{N_v}\}$$

- States  $S = \{s_1, s_2, \dots, s_{N_s}\}$

- State sequences<sup>1</sup>

$$q_1, q_2, \dots, q_T \text{ where } q_i \in S$$



# Markov chains

- Sequence can be seen as moving from one state to another, dependent only upon the previous state:
  - $P(\text{high} \mid \text{yesterday high, day before changing}) = P(\text{high} \mid \text{yesterday high})$



Low  
pressure



Pressure  
changing



High  
pressure

# State transition distribution

- Matrix A describes the state transition probabilities:

$$A = \begin{matrix} & \begin{matrix} \text{transition to} \\ \text{low} & \text{changing} & \text{high} \end{matrix} \\ \begin{matrix} \text{transition from} \\ \text{low} \\ \text{changing} \\ \text{high} \end{matrix} & \begin{bmatrix} 1/2 & 1/4 & 1/4 \\ 1/4 & 1/2 & 1/4 \\ 1/4 & 1/4 & 1/2 \end{bmatrix} \end{matrix}$$

$$a_{ij} = P(q_t = s_j | q_{t-1} = s_i)$$

Sometimes abbreviated as:  $a_{ij} = P(q_t = j | q_{t-1} = i)$

- $P(\text{high} | \text{low}) = a_{13} = 1/4$

$$\sum_{j=1}^N a_{ij} = 1$$

# Initial state distribution

- The Markov chain has a probability of starting in an initial state, denoted by the vector  $\pi$

$$\pi = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix} \begin{matrix} low \\ changing \\ high \end{matrix}$$

- In this example, the starting state has a uniform (equally probable) distribution.

# State-dependent distributions

- For each state  $s$ , there is a probability of seeing an observation  $o$ :  $b_s(o)$
- For our weather model example:

$$b_1(o) = \begin{cases} \frac{3}{4} & o = \text{rain} \\ \frac{1}{4} & o = \text{sun} \end{cases}$$

$$b_2(o) = \begin{cases} \frac{1}{2} & o = \text{rain} \\ \frac{1}{2} & o = \text{sun} \end{cases}$$

$$b_3(o) = \begin{cases} \frac{1}{4} & o = \text{rain} \\ \frac{3}{4} & o = \text{sun} \end{cases}$$



# What are the odds of that?

$P(O_1 = \text{rain}, Q_1 = \text{changing}, O_2 = \text{sun}, Q_2 = \text{high}, O_3 = \text{sun}, Q_3 = \text{high})$

which we abbreviate as:  $P(o_1, q_1, o_2, q_2, o_3, q_3)$

by chaining Bayes rule  $P(AB) = P(A | B)Pr(B)...$

$$= P(q_1)P(o_1 | q_1)P(q_2 | q_1 o_1)P(o_2 | q_2 q_1 o_1)P(q_3 | o_2 q_2 q_1 o_1)P(o_3 | q_3 o_2 q_2 o_1 q_1)$$

As  $o_i$  only dependent on state  $q_i$

$$= P(q_1)P(o_1 | q_1)P(q_2 | q_1 o_1)P(o_2 | q_2)P(q_3 | o_2 q_2 q_1 o_1)P(o_3 | q_3)$$

and as we have a first order Markov chain

$$= P(q_1)P(o_1 | q_1)P(q_2 | q_1)P(o_2 | q_2)P(q_3 | q_2)P(o_3 | q_3)$$

$$= \pi_{\text{changing}} b_{\text{changing}}(\text{rain}) a_{\text{changing}, \text{high}} b_{\text{high}}(\text{sun}) a_{\text{high}, \text{high}} b_{\text{high}}(\text{sun})$$

or using our state numbers:  $\pi_2 b_2(\text{rain}) a_{2,3} b_3(\text{sun}) a_{3,3} b_3(\text{sun})$



$S_1$   
Low  
pressure



$S_2$   
Pressure  
changing



$S_3$   
High  
pressure



# The *hidden* in hidden Markov model

- Barometer let us observe the state.
- Suppose we cannot observe state.
- Many state sequences are possible, each sequence has a probability of occurrence.

# HMM

- Let  $\Phi(A,B,\pi)$  denote a HMM where:
  - $A$  –  $N \times N$  state transition matrix.  $a_{ij}$  denotes the probability of transitioning from state  $i$  to  $j$ .
  - $B$  –  $\{b_j(k)\}$  – Set of state-dependent probability distributions.  $1 \leq j \leq N_s$ ,  $k$  in  $O$
  - $\pi$  – Initial state distribution.  $\pi_j$  is the probability of starting in state  $j$ .  $1 \leq j \leq N_s$

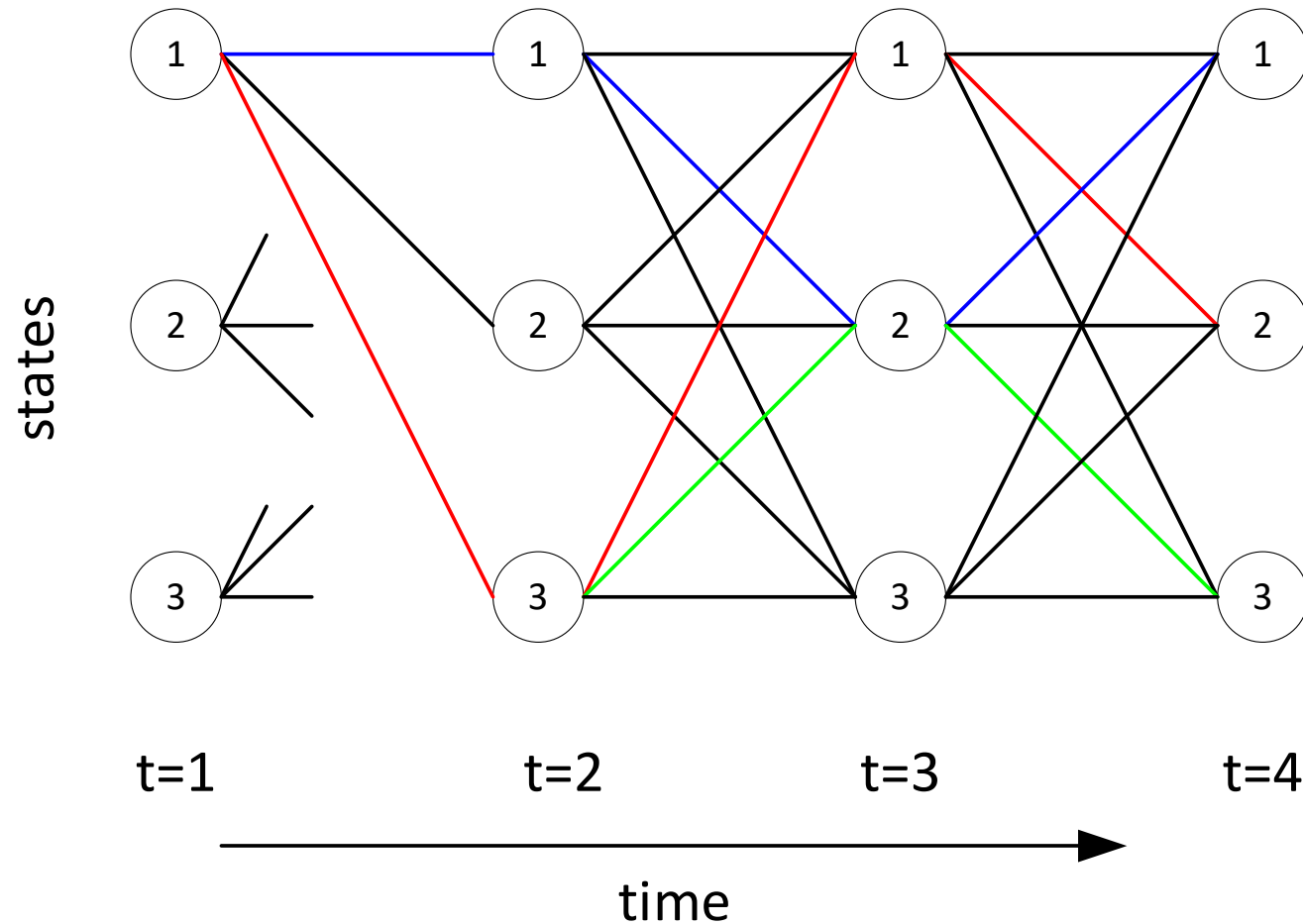
# Top 3 List for HMMs

1. What is the probability of a given sequence  $O$  given model  $\Phi$ ?
2. What state sequence is most likely to account for  $O$  in model  $\Phi$ ?
3. How can we improve the parameters of  $\Phi$  to better account for  $O$ ?



# Probability Evaluation

- Must evaluate all paths through model



Naive approach is exponential!

# Probability evaluation

- Dynamic programming can be used
- Two algorithms
  - Forward procedure

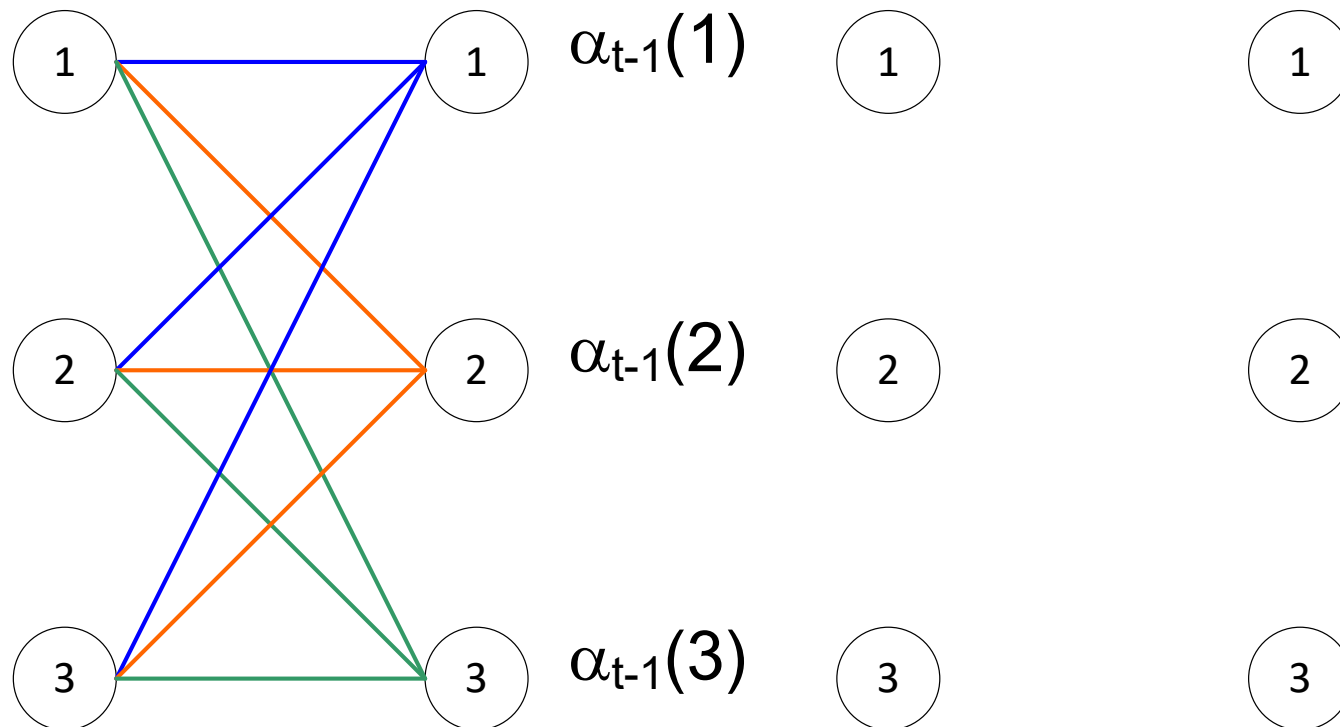
$$\alpha_t(i) = P(o_1, o_2, \dots, o_t, q_t = s_i | \Phi)$$

- Backward procedure

$$\beta_t(i) = P(o_{t+1}, o_{t+2}, \dots, o_T, q_t = s_i | \Phi)$$

# Forward algorithm

- Suppose we know the sum of all paths leading into each state  $j$  at time  $t-1$ :

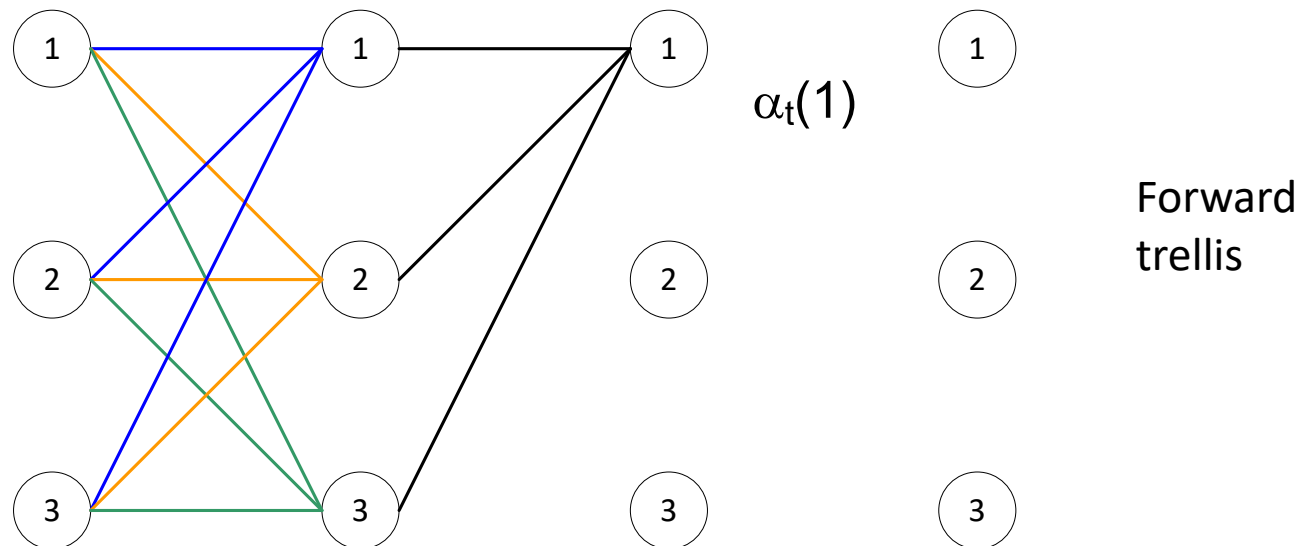


$$\alpha_{t-1}(j) = P(o_1, o_2, \dots, o_{t-1}, q_{t-1} = s_j | \Phi)$$

# Forward algorithm

- Then we can compute the probability of all paths leading into time t.

$$\alpha_t(i) = P(o_1, o_2, \dots, o_t, q_t = s_i | \Phi)$$



# Forward algorithm - $O(N^2T)$

- Initialization

$$\alpha_1(i) = \pi_i b_i(o_1) \quad 1 \leq i \leq N$$

- Induction

$$\alpha_t(j) = \left[ \sum_{i=1}^N \alpha_{t-1}(i) a_{ij} \right] b_j(o_t) \quad \begin{array}{l} 1 \leq j \leq N \\ 2 \leq t \leq T \end{array}$$

- Termination

$$P(O|\Phi) = \sum_{i=1}^N \alpha_T(i)$$

# Optimal State Sequence

- The forward algorithm finds all paths through a model.
- Sometimes, we are interested in the best path through the model:
  - Perhaps we are interested in determining which states are associated with which observations.
  - Frequently, most of the paths contribute very little to the overall probability.

# Viterbi algorithm:

## Determine optimal state sequence

- Another example of dynamic programming
- Finds the most likely path through the model
- Similar to the forward algorithm
  - Uses max instead of sum
  - Keeps extra information about the best path

# Viterbi Algorithm

- Initialization

- Probability to start in state  $i$

$$V_1(i) = \pi_i b_i(o_1) \quad 1 \leq i \leq N$$

- $B_t(i)$  – The previous state which transitioned into state  $i$  at time  $t-1$ . (0 indicates no previous state.)

$$B_1(i) = 0 \quad 1 \leq i \leq N$$



# Viterbi Algorithm $O(N^2T)$

- Induction

- Find the best path leading into the current state and account for the probability of the observation

- Record the previous node on the best path

$$V_t(j) = \max_{1 \leq i \leq N} [V_{t-1}(i) a_{ij}] b_j(o_t) \quad \begin{array}{l} 1 \leq j \leq N \\ 2 \leq t \leq T \end{array}$$

$$B_t(j) = S_{\operatorname{argmax}_{1 \leq i \leq N} [V_{t-1}(i) a_{ij}]} \quad \begin{array}{l} 1 \leq j \leq N \\ 2 \leq t \leq T \end{array}$$

# Viterbi Algorithm

- Termination

$$p^{BestPath}(X) = \max_{1 \leq i \leq N} [V_T(i)]$$

$$q_T^{BestPath} = S_{\arg \max_{1 \leq i \leq N} [V_T(i)]}$$

- Extracting the best path:

for  $t = T-1, T-2, \dots, 1$

$$q_t^{BestPath} = S_{B_{t+1}(q_{t+1}^{BestPath})}$$

# Log Domain Viterbi Implementation

- Operates in the log probability domain
- Multiplications are replaced with addition
- Not covered in text.

# Parameter Estimation

- Application of the Expectation-Maximization (EM) algorithm
  - If we had all the information
    - true state sequence
    - observations
  - then techniques such as maximum-likelihood estimation could be used to improve our parameter set

# EM algorithm

- Problem: Some information is unknown
- Solution:
  - Use the expectation operator to determine the expected values of missing parameters
  - Determine new parameters
  - Repeat

# EM Algorithm

- Guaranteed to converge to a local maximum
- For speech applications, no more than 5-15 iterations are typically required
- For HMMs, the resulting formulae are known as the Baum-Welch reestimation equations.

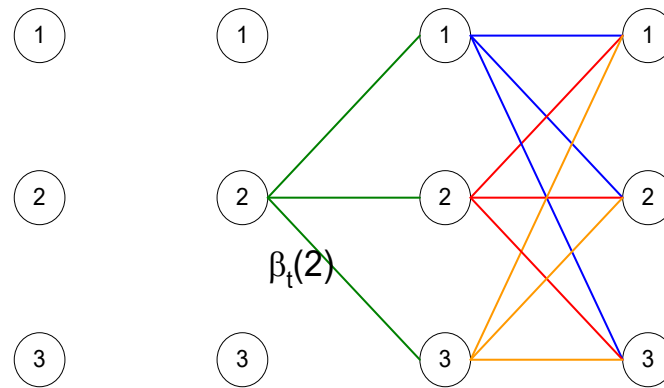
# Some needed concepts

- The backward algorithm – similar to the forward algorithm, but works from  $T$  down towards 1.
- $\gamma_t(i,j)$  – Given a model and observation sequence, the probability of transitioning from state  $i$  at  $t-1$  to  $j$  at  $t$  given the model  $\Phi$  and acoustic evidence  $O$

$$P(q_{t-1} = s_i, q_t = s_j | O, \Phi)$$

# Backward algorithm

$$\beta_t(i) = P(o_{t+1}, o_{t+2}, \dots, o_T \mid q_t = s_i, \Phi)$$



Note:  $\beta_t(i)$  does not include the probability of observing  $o_t$ .





# Backward Algorithm $O(N^2T)$

- Initialization

$$\beta_T(i) = 1 \quad 1 \leq i \leq N$$

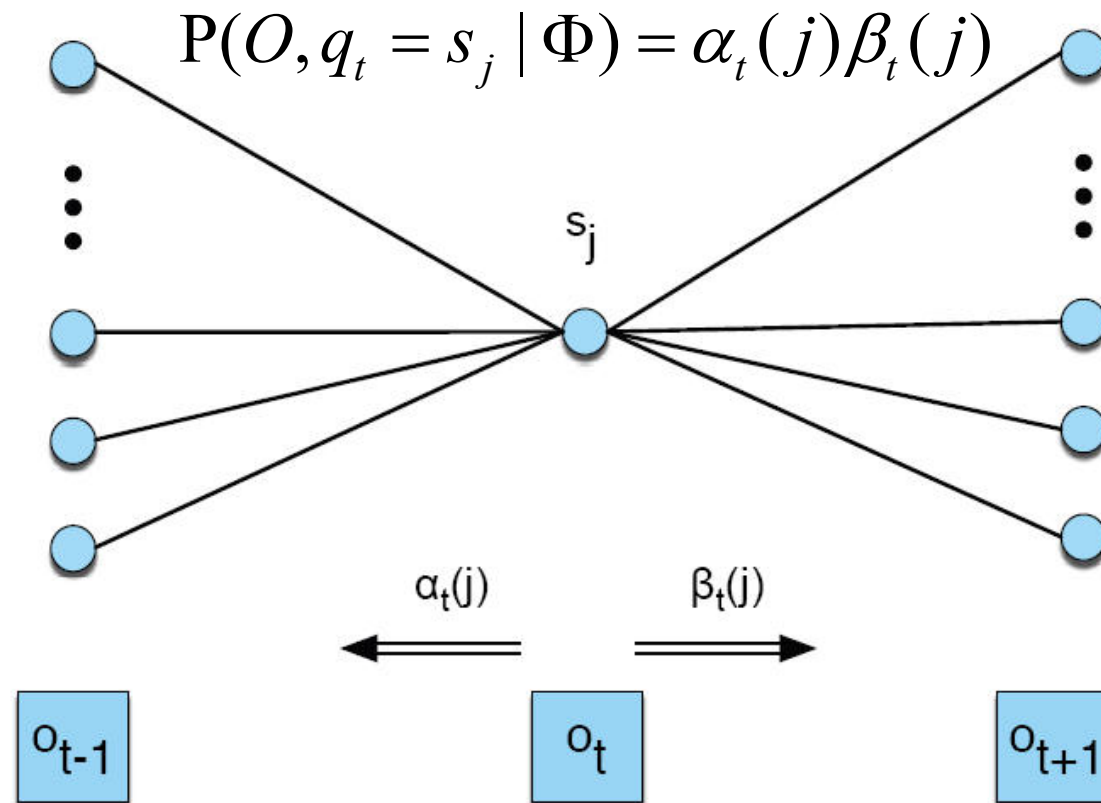
- Induction

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j) \quad \begin{array}{l} 1 \leq j \leq N \\ 1 \leq t \leq T - 1 \end{array}$$

- Termination not needed, but possible

# Forward-Backward relationship

- $\alpha_t(i)$  = all paths into  $q_t=s_i$  and observing  $o_1, o_2, \dots, o_t$ .
- $\beta_t(i)$  = all paths out of  $q_t=s_j$  and observing  $o_{t+1}, o_{t+2}, \dots, o_T$ .

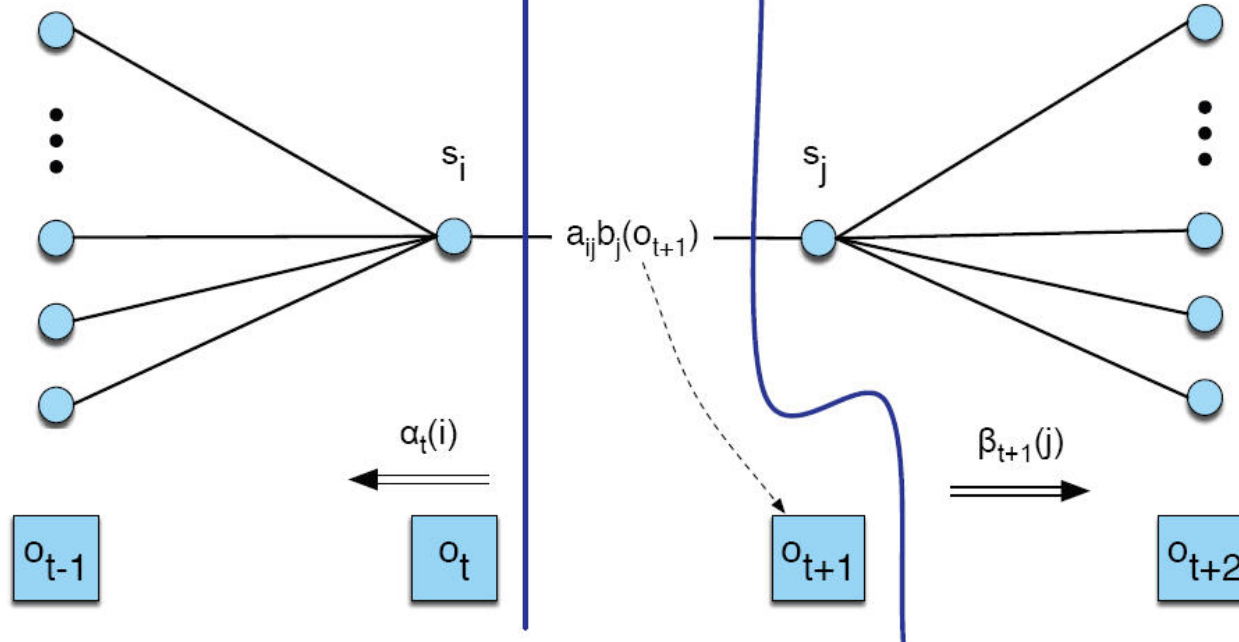


J&M p. 191

# Constrained path probability

Suppose we want the probability of all paths through a specific transition from  $t-1$  to  $t$ :

$$P(O, q_t = s_i, q_{t+1} = s_j | \Phi) \\ = \alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)$$



J&M p. 190

$\gamma_t(i, j)$  Probability state  $i$  to  $j$  at time  $t$

$$\begin{aligned}\gamma_t(i, j) &\stackrel{\Delta}{=} P(q_{t-1} = s_i, q_t = s_j | O, \Phi) \\ &= \frac{P(q_{t-1} = s_i, s_t = s_j, O | \Phi)}{P(O | \Phi)} \\ &= \frac{\alpha_{t-1}(i) a_{ij} b_j(o_t) \beta_t(j)}{\sum_{k=1}^N \alpha_T(k)}\end{aligned}$$

Bayes rule

previous slide, and

$$P(O | \Phi) = \sum_{k=1}^N \alpha_T(k)$$

## Special case: $\gamma_1(i, j)$

- Calls for non-existent transition between  $q_0$  and  $q_1$ . We define  $\alpha_0(i)=1$  and  $a_{ij}$  at time 0 as  $\pi_j$ :

$$\begin{aligned}\gamma_1(i, j) &= \frac{\alpha_0(i)a_{ij}b_j(o_1)\beta_1(j)}{\sum_{k=1}^N \alpha_T(k)} \\ &= \frac{\pi_j b_j(o_1)\beta_1(j)}{\sum_{k=1}^N \alpha_T(k)}\end{aligned}$$

# Baum-Welch equations

- The EM algorithm can be used to derive the Baum-Welch reestimation equations.
- Computation of the expectations is done with the  $\gamma$  function.
- Once the expectation has been computed, a maximum likelihood estimate can be computed for the model.

# Initial state distribution

- Percentage of time that we will be in state  $i$  at time  $t=1$

$$\hat{\pi}_i = \sum_{k=1}^N \gamma_1(i, k)$$

- For many speech applications, we desire a starting state. In this case  $\pi_{\text{start}}=1$ . The reestimation formulas will not change this.

# State transition

- We can think of this as the expected number of transitions from state  $i$  to  $j$  divided by the expected number of all transitions:

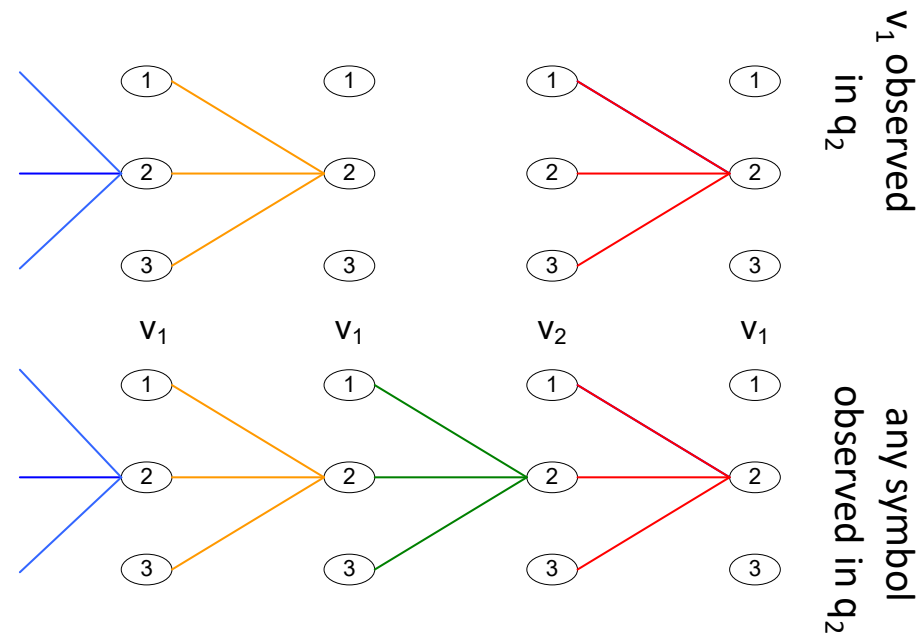
$$\hat{a}_{ij} = \frac{\sum_{t=2}^T \gamma_t(i, j)}{\sum_{t=2}^T \sum_{k=1}^N \gamma_t(i, k)}$$

Note: Huang, Acero, and Hon sum from 1 to  $N$  (eq. 8.40, p 392) which includes the initial state probability  $\pi_i$ . Most authors do not do this.



# State-dependent pdfs

- For each time where symbol  $v_k$  is seen, we need to determine the probability of seeing  $v_k$  given that we are in state  $s_j$ .
- We divide this expectation by the expected probability of any symbol given state  $s_j$ .



# State-dependent pdfs

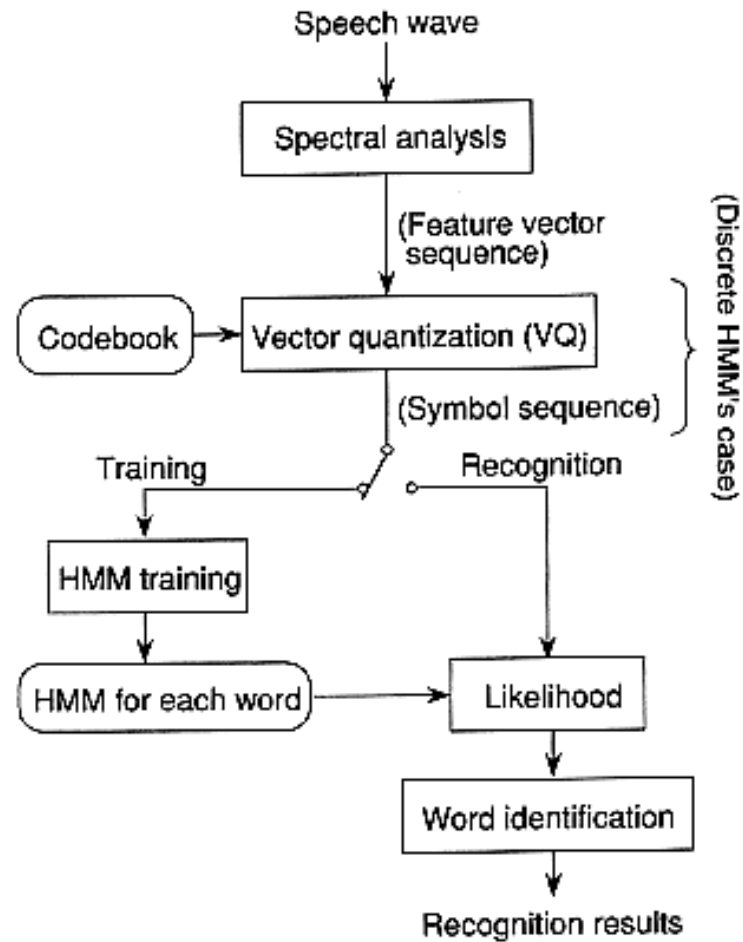
- Divide the expected number of transitions into state  $j$  where symbol  $o_k$  occurs by all transitions into state  $j$ :

$$\hat{b}_j(k) = \frac{\sum_{t=1}^T \sum_{i=1}^N \gamma_t(i, j) \delta(o_t, v_k)}{\sum_{t=1}^T \sum_{i=1}^N \gamma_t(i, j)}$$

where  $\delta(o, v) \triangleq \begin{cases} 1 & o = v \\ 0 & o \neq v \end{cases}$

$$= \frac{\sum_{t \text{ such that } o_t = v_k} \sum_{i=1}^N \gamma_t(i, j)}{\sum_{t=1}^T \sum_{i=1}^N \gamma_t(i, j)}$$

# Isolated Word Recognizer



Vector quantization is the same think as k-means. We map a continuous vector to a discrete symbol.