


Artificial Intelligence


Graphical Models



Mountains & Minds 1

Making Simple Decisions


- So far, we have discussed reasoning under uncertainty, i.e., methods for representing and reasoning about state information that we do not know with certainty.
- We applied a probabilistic framework for our discussion.
- For our agent-based system, we now need to use this information within a process that allows the agent to act rationally.
- In other words, we need to determine how to make good decisions over actions based on our beliefs.



Mountains & Minds 2

Utility Theory

- The framework we will use to guide decision making will be based on *utility theory*.
- An agent's preferences between world states will be represented via a *utility function*.
- The utility function will assign a number to each state to express desirability of that state.



Mountains & Minds 3

Some Notation

- Let $U(S)$ denote the utility of state S .
- Given an action A , let $Result(A)$ be the outcome state resulting from applying action A . Note: In general, the outcome state will be conditioned on both the action *and* the current state.
- Let $P(Result(A) | Do(A), E)$ denote the probability that the result of A given performing A under evidence E is $Result(A)$.
- Then the expected utility, denoted $EU(A | E)$ is defined as

$$EU(A | E) = \sum_i P(Result_i(A) | Do(A), E)U(Result_i(A))$$

- The principle of *maximum expected utility* (MEU) states that a rational agent should choose an action that maximizes the agent's expected utility.

Constraints on Preference

- We represent rational preferences between states for an agent as follows:

$A \succ B$ A is preferred to B

$A \sim B$ the agent is indifferent between A and B

$A \succeq B$ the agent prefers A to B or is indifferent between them

- A *lottery* is a probability distribution over a set of actual outcomes.
- Lottery L with possible outcomes C_1, \dots, C_n that can occur with probabilities p_1, \dots, p_n is written,
 - $L = [p_1, C_1; p_2, C_2; \dots; p_n, C_n]$

Constraints on Preference

- Orderability:** Given any two states, a rational agent must either prefer one to the other or be indifferent.

$$(A \succ B) \vee (B \succ A) \vee (A \sim B)$$

- Transitivity:** Given any three states, if A is preferred over B and B is preferred over C , then A is preferred over C .

$$(A \succ B) \wedge (B \succ C) \Rightarrow (A \succ C)$$

- Continuity:** If some state B is between states A and C in preference, then there is some probability p for which one is indifferent towards definitely getting B or having a lottery yield A or C according to probabilities p and $(1 - p)$.

$$(A \succ B \succ C) \Rightarrow \exists p [p, A; 1 - p, C] \sim B$$

Constraints on Preference

- **Substitutability:** If an agent is indifferent on two lotteries A and B , then it is indifferent between two more complex lotteries that are the same except B is substituted for A in one of them.

$$(A \sim B) \Rightarrow [p, A; 1-p, C] \sim [p, B; 1-p, C]$$

- **Monotonicity:** Given two lotteries, each with outcomes A and B , if an agent prefers A to B , then the agent must prefer the lottery with higher probability for A .

$$(A \succ B) \Rightarrow (p \geq q \Leftrightarrow [p, A; 1-p, B] \succeq [q, A; 1-q, B])$$

- **Decomposability:** Compound lotteries can be reduced to simpler ones via the laws of probability.

$$[p, A; 1-p, [q, B; 1-q, C]] \sim [p, A; (1-p)q, B; (1-p)(1-q), C]$$

Utility Principle

- If an agent's preferences obey the axioms of utility, then there exists a real-valued function U that operates on states such that

- $U(A) > U(B)$ iff A is preferred to B and
- $U(A) = U(B)$ iff the agent is indifferent between A and B .

$$U(A) > U(B) \Leftrightarrow A \succ B$$

$$U(A) = U(B) \Leftrightarrow A \sim B$$

- The utility of a lottery is the sum of the probability of each outcome times the utility of that outcome.

$$U([p_1, S_1; \dots; p_n, S_n]) = \sum_i p_i U(S_i)$$

Expected Monetary Value

- Suppose we base utility on monetary value of some event.
- For example, suppose we have a lottery where you must make a choice:
 - You are guaranteed \$1,000,000, or
 - You can gamble on the flip of a coin where heads gets you nothing but tails gets you \$3,000,000.
- Then the expected monetary value of this game is (assuming a fair coin) $0.5 \times \$0 + 0.5 \times \$3,000,000 = \$1,500,000$.
- The EMV of doing nothing is \$1,000,000.
- What would you do?

Expected Monetary Value

- This is not a simple decision.
- Specifically, utility is not directly proportional to EMV but is based on change in lifestyle.
- In terms of expected utility,
 - $EU(\text{Accept}) = 0.5U(S_k) + 0.5U(S_{k+3,000,000})$
 - $EU(\text{Decline}) = U(S_{k+1,000,000})$
- One study (Grayson 1960) demonstrated that utility of money was almost exactly proportional to the *logarithm* of the amount.

Multiattribute Utility Theory

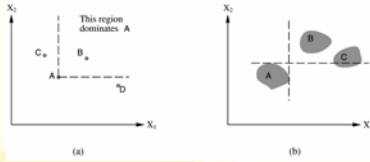
- Suppose we are attempting to determine where to locate a new airport.
- Several factors must be considered in this analysis:
 - disruption due to construction
 - cost of land
 - distance from population centers
 - noise pollution
 - safety conditions
 - etc.
- Characterizing outcomes by two or more attributes leads to *multiattribute utility theory*.

Multiattribute Utility Theory

- Let $\mathbf{X} = X_1, \dots, X_n$ be a set of attributes to consider.
- Let $\mathbf{x} = \langle x_1, \dots, x_n \rangle$ be a complete vector of assignments for the attributes.
- Then multiattribute utility theory considers complex utility functions of the form $U(\mathbf{X})$.
- We will consider various forms and features of different utility functions.

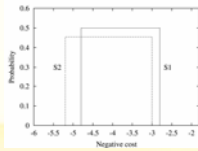
Dominance

- We say there exists *strict dominance* of $U(\mathbf{x}_j)$ over $U(\mathbf{x}_i)$ if an option is of lower value on x_j than x_i on all attributes.
- When strict dominance applies, the dominated alternative need not be considered further.



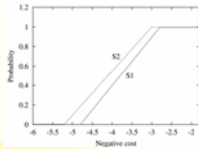
Stochastic Dominance

- Consider a situation with one attribute.
- Suppose we are considering cost to site an airport at S_1 or S_2 .
 - $C(S_1) - U(\$2.8B, \$4.8B)$
 - $C(S_2) - U(\$3B, \$5.2B)$
 - In this case, U represents uniform distribution.



Stochastic Dominance

- Given that utility decreases with cost, S_1 should dominate S_2 .
- This does not come from comparing expected cost; knowing a value absolutely for one makes it harder to decide.
- Consider the *cumulative* distributions.
- If S_1 is always "to the right" of S_2 , then S_1 is stochastically cheaper than S_2 .



Stochastic Dominance

- If two actions A_1 and A_2 lead to distributions $p_1(x)$ and $p_2(x)$ on attribute X , then A_1 stochastically dominates A_2 on X if

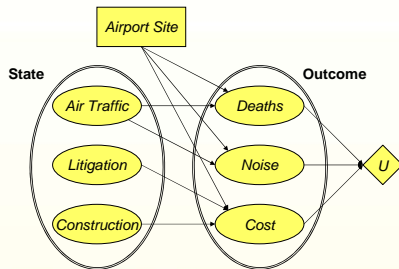
$$\forall x \int_{-\infty}^x p_1(x') dx' \leq \int_{-\infty}^x p_2(x') dx'$$

- If A_1 stochastically dominates A_2 , then for any monotonically nondecreasing utility function $U(x)$, the expected utility of A_1 is at least as high as the expected utility of A_2

Decision Networks

- Decision networks (also called influence diagrams) combine Bayesian networks with additional node types for actions and utilities
- **Chance nodes:** Represented as ovals and correspond to random variables.
- **Decision nodes:** Represented as rectangles and correspond to choices of actions.
- **Utility nodes:** Represented as diamonds and correspond to computation of utility based on parents.

Example Decision Network

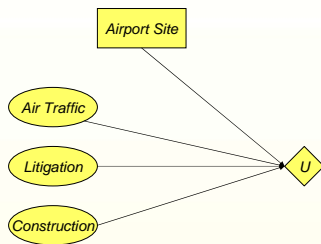


Simplified Form

- Frequently, a simplified form of decision network is used.
- Chance nodes describing “outcome states” are omitted.
- Utility nodes are connected directly to current state nodes and decision node.
- Then the utility function is defined over actions rather than states.
- Specifically,

$$EU(A|E) = \sum_i P(\text{Result}_i(A) | \text{Do}(A), E) U(\text{Result}_i(A))$$

Action-Utility Diagrams



Evaluating Decision Networks

- Actions are selected by evaluating the decision network for each possible decision.
- Given the decision node is set to a particular condition, it can be treated like a chance node that has been set as an evidence variable.
- The evaluation is like a Bayesian network:
 - Set the evidence variable(s) for the current state.
 - For each possible value of decision node:
 - Set the decision node to that value.
 - Calculate the posterior probabilities for the parent nodes of the utility node.
 - Calculate the utility for the action.
 - Return the action with the highest utility.

The Value of Information

- So far, we have assumed that when we make a decision, all relevant information has been provided.
- In general, this is rarely the case.
- Given we need to seek information, how do we determine what information to get?
- *Information value theory* enables to evaluate the quality of information to help an agent decide what information to seek.

Example

- An oil company wants to buy one of n indistinguishable blocks of ocean drilling rights.
- Exactly one block contains oil worth C dollars.
- The price of one block is C/n dollars.
- A survey has been performed on block 3, indicating definitely whether or not the block has oil.
- How much should the company pay for that survey?
- To answer this, we will consider what the company would do with the information.

Example

- With probability $1/n$, the survey will indicate oil in block 3. In this case, the company will buy block 3 for C/n dollars and make a profit of $C - C/n = (n-1)C/n$ dollars.
- With probability $(n-1)/n$, the survey will show that the block contains no oil (since one and only one block has oil). In that case, the company will buy a different block.
- If block 3 does not have oil, then the probability of finding oil in another block changes from $1/n$ to $1/(n-1)$, so the company makes an expected profit of $C/(n-1) - C/n = C/n(n-1)$.

Example

- We are now ready to compute expected profit, given the survey information.

$$\begin{aligned} \text{ExpectedProfit}(B_3) &= P(B_3)\text{Profit}(B_3) + P(-B_3)\text{Profit}(-B_3) \\ &= \frac{1}{n} \times \frac{(n-1)C}{n} + \frac{n-1}{n} \times \frac{C}{n(n-1)} \\ &= C/n \end{aligned}$$

- Note this is the same as the purchase price of a block.
- Thus the company would be willing to pay up to the price of a block for the survey.

Value of Information

- We assume that exact evidence is obtained about the value of some random variable E_j .
- We, therefore, refer to the *value of perfect information (VPI)*.
- To compute this, suppose the agent's current knowledge is given as E .
- Then the value of the current best action α is

$$EU(\alpha | E) = \max_A \sum_i U(\text{Result}_i(A))P(\text{Result}_i(A) | Do(A), E)$$

- We can extend this for after E_j is obtained.

$$EU(\alpha_{E_j} | E, E_j) = \max_A \sum_i U(\text{Result}_i(A))P(\text{Result}_i(A) | Do(A), E, E_j)$$

Value of Information

- Note that, in the latter case, E_j is *currently* unknown.
- Consequently, we must "average" over all possible values e_{jk} that might be discovered for E_j using current beliefs about its value.
- Thus, the value of discovering E_j , given the current information E , is given as

$$VPI_E(E_j) = \left(\sum_k P(E_j = e_{jk} | E) EU(\alpha_{e_{jk}} | E, E_j = e_{jk}) \right) - EU(\alpha | E)$$

- The subtraction corresponds to a "change" in utility due to obtaining information E_j .

Example

- Deciding whether to collect the evidence still may not be so easy.
- Suppose we are deciding between two actions, A_1 and A_2 .
- Suppose their expected utilities are U_1 and U_2 respectively.
- Suppose further that E_j will yield some new expected utilities U'_1 and U'_2 for the actions.
- Prior to collecting E_j , we have probability distributions over values of U_1 and U_2 .

Example

- We will assume our two actions A_1 and A_2 are deciding between two routes through a mountain range in the winter with avalanches possible.
 - A_1 is a straight highway through a low pass.
 - A_2 is a winding dirt road over the top.
- With no other information, it is safe to say that A_1 should be preferred over A_2 (i.e., $U_1 > U_2$).
- Suppose satellite imagery is available to give the actual conditions of the two roads, and this yields the new utilities, U'_1 and U'_2 .
- The report is unlikely to change the utilities significantly.

Example

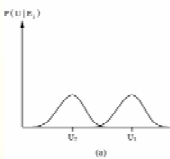
- Suppose now we consider two winding dirt roads through the mountains.
- Further, suppose the distances are approximately the same, and we are taking an injured passenger for help.
- Here, the distributions for the two utilities are very close, and getting additional information could help choose the better road.
- In this case, the value of the additional information could be very high.

Example

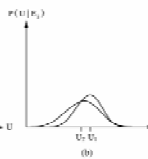
- Finally, suppose it is summer rather than winter.
- During this season, blockage due to avalanche is unlikely.
- Now, the satellite imagery might indicate one road is more scenic than another.
- In this case, we “might” change our choice given additional information, but the change in value of the information is sufficiently low that we would not get the reports.

Example

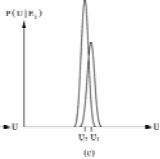
Clear separation of utilities, thus low change in utility



Poor separation of utilities, high change in utility



Poor separation of utilities, low change in utility



Properties of VPI

- First, the value of information is non-negative.

$$\forall j, E \quad VPI_E(E_j) \geq 0$$

- Second, VPI depends on the current state of information, so VPI is not additive.

$$VPI_E(E_j, E_k) \neq VPI_E(E_j) + VPI_E(E_k)$$

- Third, VPI is order-independent.

$$\begin{aligned} VPI_E(E_j, E_k) &= VPI_E(E_j) + VPI_{E, E_j}(E_k) \\ &= VPI_E(E_k) + VPI_{E, E_k}(E_j) \end{aligned}$$

Information-Gathering Agent

- Suppose an observable evidence variable E_j has an associated cost $Cost(E_j)$.
- This cost reflects the cost for obtaining the evidence.
- An information-gathering agent would attempt to select the piece of evidence with the highest value.
- Assume the result of the action $Request(E_j)$ is that the next percept provides E_j .

Information-Gathering Agent

- One type of agent can select evidence in a “myopic” way, looking ahead only one percept at a time.
- This approach is analogous to greedy search.

```
function INFORMATION-GATHERING-AGENT(percept)
  returns an action
  static: D, a decision network

  integrate percept into D
  j ← the value that maximizes  $VPI(E_j) - Cost(E_j)$ 
  if  $VPI(E_j) > Cost(E_j)$  then return REQUEST( $E_j$ )
  else return the best action from D
```

Bayes' Decision Theory

- The following is a *brief* introduction to Bayes' decision theory—more will be provided during the class on Bayesian learning.
- Suppose we are attempting to decide between two choices, A_1 and A_2 .
- Assume for the sake of discussion that only one choice is the “right” choice.
- Finally, assume that we must make our choice using only some set of observations e .

Bayesian Risk

- In Bayesian decision theory, we want to decide between choices based on some set of beliefs.
- Associated with each choice is the possibility that it is the wrong choice (or at least not the best choice).
- The wrong choice can then lead to some negative impact relative to the best action.
- We will denote the "loss" due to making some choice α_i for some "correct choice" A_j as $\lambda(\alpha_i | A_j)$.
- Then we can define the "expected loss" (also called *risk*) as

$$R(\alpha_i) = \sum_{j=1} \lambda(\alpha_i | A_j) P(A_j | \mathbf{e})$$

Minimizing Risk

- We find the last term using Bayes' Rule
- $$P(A_j | \mathbf{e}) = \frac{P(\mathbf{e} | A_j) P(A_j)}{P(\mathbf{e})}$$
- Now we want to make a choice that minimizes risk.
- $$\alpha(\mathbf{e}) = \arg \min R(\alpha_i | \mathbf{e})$$
- Let λ_{ij} be the loss of choosing A_j when A_i is better.
 - So, in choosing between A_1 and A_2 ,

$$R(\alpha_1 | \mathbf{e}) = \lambda_{11} P(A_1 | \mathbf{e}) + \lambda_{12} P(A_2 | \mathbf{e})$$

$$R(\alpha_2 | \mathbf{e}) = \lambda_{21} P(A_1 | \mathbf{e}) + \lambda_{22} P(A_2 | \mathbf{e})$$

Decision Rule

- Using risk, choose A_1 if $R(\alpha_1 | \mathbf{e}) < R(\alpha_2 | \mathbf{e})$.
- Substituting the definitions for risk,

$$\lambda_{11} P(A_1 | \mathbf{e}) + \lambda_{12} P(A_2 | \mathbf{e}) < \lambda_{21} P(A_1 | \mathbf{e}) + \lambda_{22} P(A_2 | \mathbf{e})$$
- Rearranging terms,

$$(\lambda_{21} - \lambda_{11}) P(A_1 | \mathbf{e}) > (\lambda_{12} - \lambda_{22}) P(A_2 | \mathbf{e})$$
- Using Bayes' Rule

$$(\lambda_{21} - \lambda_{11}) P(\mathbf{e} | A_1) P(A_1) / P(\mathbf{e}) > (\lambda_{12} - \lambda_{22}) P(\mathbf{e} | A_2) P(A_2) / P(\mathbf{e})$$
- That is,

$$(\lambda_{21} - \lambda_{11}) P(\mathbf{e} | A_1) P(A_1) > (\lambda_{12} - \lambda_{22}) P(\mathbf{e} | A_2) P(A_2)$$
- Assuming no loss for "correct choice" and equal risk, choose A_1 if $P(\mathbf{e} | A_1) P(A_1) > P(\mathbf{e} | A_2) P(A_2)$.

Reasoning Over Time

- A changing world is modeling using random variable(s) at a point (i.e., slice) in time.
- Relating these random variables across time frames describes how the state evolves through time.
- These state transitions can be probabilistic.
- Example Models:
 - Hidden Markov Models
 - Kalman Filters
 - Dynamic Bayesian Networks

Markov Processes

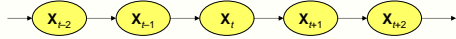
- A process is considered to be *stationary* if the rules governing the process do not change over time.
- Suppose we are interested in the state of a set of random variables, \mathbf{X} at some time slice t (i.e., \mathbf{X}_t) and want to determine the probability distribution over the values of \mathbf{X} based on how the process evolves, $\mathbf{P}(\mathbf{X}_t | \mathbf{X}_{0:t-1})$.
- The *Markov Assumption* states that the current state of a process depends only on a *finite* history of previous states.

First-Order Markov Processes

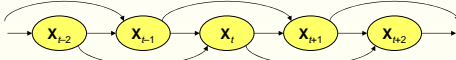
- The simplest case is when the current state depends only on the previous state.
- The corresponding *transition model* is given by $\mathbf{P}(\mathbf{X}_t | \mathbf{X}_{0:t-1}) = \mathbf{P}(\mathbf{X}_t | \mathbf{X}_{t-1})$.
- Sometimes, we are not able to observe state directly but can only infer state based on a set of evidence variables, \mathbf{E}_t .
- The corresponding *sensor model* is given by $\mathbf{P}(\mathbf{E}_t | \mathbf{X}_{0:t}, \mathbf{E}_{0:t-1}) = \mathbf{P}(\mathbf{E}_t | \mathbf{X}_t)$.
- In addition to the transition and sensor models, we need to specify the prior probability $\mathbf{P}(\mathbf{X}_0)$.

Example

First-Order Markov Process



Second-Order Markov Process



Inference Problems

- **Filtering:** Computing a belief state, i.e., the posterior distribution over the current state, given the evidence to date.
 - $P(X_t | e_{1:t})$
- **Prediction:** Computing the posterior distribution over *future* states, given the evidence to date.
 - $P(X_{t+k} | e_{1:t})$
- **Smoothing:** Computing the posterior distribution over a *past* state, given all evidence up to the present.
 - $P(X_k | e_{1:t})$ for $0 \leq k < t$.
- **Explanation:** Given a sequence of observations, find a sequence of states most likely to have generated those observations.
 - $\text{argmax } P(x_{1:t} | e_{1:t})$

Dynamic Bayesian Networks

- A DBN is a Bayesian network that represents a temporal probability model that can be used to perform the previous temporal inference problems.
- Hidden Markov Models and Kalman Filters are two examples of DBNs.
- In general, each time slice of a DBN can have any number of state variables X_t and evidence variables E_t .

Specifying DBNs

- Three kinds of information are required to specify a DBN:
 - The prior distribution over the state variables, $P(\mathbf{X}_0)$.
 - The transition model, $P(\mathbf{X}_{t+1} | \mathbf{X}_t)$.
 - The sensor model, $P(E_t | \mathbf{X}_t)$.
- Recall the assumption that the transition and sensor models are stationary, i.e., they are the same for all t .
- Then only two slices need to be specified, $t = 0$ (base case), and $t = 1$ (for all subsequent cases).
- Inference involves “unrolling” the DBN.

46

Signal Processing

- We are interested in processing audio signals to extract and interpret speech.
- Two forms of signal models:
 - Deterministic models—parameterized combination of basis functions (e.g., sine waves or exponentials)
 - Stochastic models—statistical models as combinations of Gaussian processes, Poisson processes, Markov processes.
- We will focus on a DBN, the Hidden Markov Model, for interpreting the signals.

47

Discrete Markov Processes

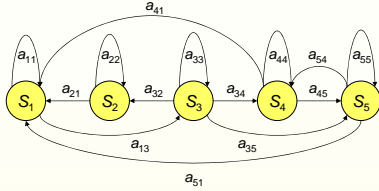
- Consider a system described at any time as being in one of N distinct states, S_1, S_2, \dots, S_N .
- At regularly spaced, discrete times, the system undergoes a change of state.
- Transitions follow some set of probabilities associated with each state.

$$P[q_t = S_j | q_{t-1} = S_i, q_{t-2} = S_k, \dots]$$

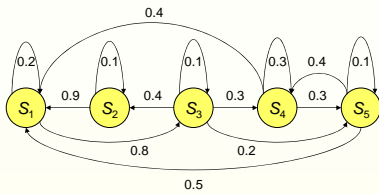
$$= P[q_t = S_j | q_{t-1} = S_i]$$

48

Discrete Markov Processes



Discrete Markov Processes



Discrete Markov Processes

$$A = \{a_{ij}\} = \begin{bmatrix} 0.2 & 0.0 & 0.8 & 0.0 & 0.0 \\ 0.9 & 0.1 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.4 & 0.1 & 0.3 & 0.2 \\ 0.4 & 0.0 & 0.0 & 0.3 & 0.3 \\ 0.5 & 0.0 & 0.0 & 0.4 & 0.1 \end{bmatrix}$$

Questions to Ask (Given Model)

- Let $O = S_1 S_2 \dots S_k$ be some sequence of states.
- What is the probability of following a particular sequence of state transitions of length k assuming it starts in state S_1 ?

$$P(O|M) = 1 * \prod_{i=2}^k P(O_i | O_{i-1})$$

- Given that the model is in a known state, what is the probability it stays in that state for exactly k time steps?

$$P(O|M, q_1 = S_i) = (a_{ii})^{k-1} (1 - a_{ii})$$

- Given that the model starts in a known state, what is the expected number of time steps it will remain in that state?

$$E[O_i] = \sum_{d=1}^{\infty} d * (a_{ii})^{d-1} (1 - a_{ii}) = \frac{1}{1 - a_{ii}}$$

Hidden Markov Models

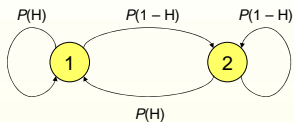
- So far, we have assumed we can directly observe all of the states.
- Now we consider the situation where the observations are probability functions of the states.
- We will consider two classic problems:
 - The Coin Toss problem.
 - The Urn and Ball problem.

Coin Toss Models

- You are in a room with a barrier.
- Another person is on the other side of the barrier.
- The other person is tossing one or more coins (you don't know how many).
- The only information you get is the sequence of coin tosses.
- Your task is to construct an HMM to explain the observations.

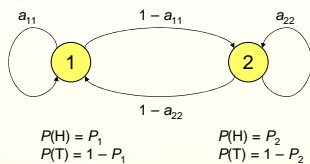
Possible HMMs

- One-coin model
 - Directly observable (i.e., not really hidden).
 - Only need to determine probability for Heads.



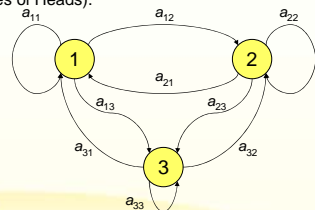
Possible HMMs

- Two-coin model
 - States represent choice of two coins.
 - Each state includes model of underlying coin (i.e., individual probabilities of Heads).



Possible HMMs

- Three-coin model
 - States represent choice of three coins.
 - Each state includes model of underlying coin (i.e., individual probabilities of Heads).



Urn and Ball Models

- Suppose we have N large glass urns.
- Within each urn we have a large number of colored balls.
- Assume there are M distinct colors available.
- A genie in the room randomly chooses an urn and pulls out a ball.
- The color is reported, and the ball is replaced.
- The process repeats for some number steps.
- We want the resultant HMM.
 - Each state is an urn.
 - Color probabilities are associated with each urn.

Hidden Markov Models

- **Def:** $HMM = \langle N, M, A, B, \pi \rangle$ where
 - N = the number of states in the model, where the states are denoted $S = \{S_1, \dots, S_N\}$.
 - M = the number of distinct observation symbols per state (i.e., the discrete alphabet size), where the symbols are denoted $V = \{V_1, \dots, V_M\}$.
 - A = the state transition probability distribution, $A = \{a_{ij}\} = P(q_{t+1} = S_j | q_t = S_i)$.
 - B_j = the observation symbol probability distribution in state j , $B_j = b_j(k) = P(v_k \text{ at } t | q_t = S_j)$
 - π = the initial state distribution, $\pi = \{\pi_i\} = P(q_1 = S_i)$

Observation Sequences

- Given an HMM, an observation sequence, $O = O_1 \dots O_T$, can be built as follows:
 - Choose an initial state q_1 according to π .
 - Set $t = 1$.
 - Choose $O_t = v_k$ for state S_j according to $b_j(k)$.
 - Transition to successor state $q_{t+1} = S_j$ according to a_{ij} .
 - Set $t = t + 1$ and repeat until full sequence generated.

Three Problems for HMMs

- **Model Evaluation:** Given observation sequence $O = \{O_1 \dots O_T\}$, and model $M = (A, B, \pi)$, how do we efficiently compute $P(O|M)$, i.e., the probability of the observation sequence given the model?
- **Sequence Determination:** Given observation sequence O and model M , how do we choose a corresponding state sequence $Q = q_1 \dots q_T$ that best explains the observations?
- **Model Discovery:** How do we adjust the model parameters M to maximize $P(O|M)$?

Model Evaluation

- We want to calculate $P(O|M)$.
- Consider sequence $Q = q_1 \dots q_T$.

$$P(O|M) = \sum_Q P(O|Q, M)P(Q|M)$$

$$P(O|Q, M) = \prod_{t=1}^T P(O_t | q_t, M)$$

$$= b_{q_t}(O_t) \cdot b_{q_2}(O_2) \cdots b_{q_T}(O_T)$$

$$P(Q|M) = \pi_{q_1} a_{q_1 q_2} a_{q_2 q_3} \cdots a_{q_{T-1} q_T}$$

- This approach is not efficient!

Forward-Backward Procedure

- Define a "forward" variable to be $\alpha_t(i) = P(O_1 O_2 \dots O_t, q_t = S_i | M)$, i.e., the probability of a partial observation sequence O_1, O_2, \dots, O_t and state S_i at time t , given our model.

- Initialization: $\alpha_1(i) = \pi_i b_i(O_1), 1 \leq i \leq N$

- Induction:

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(O_{t+1})$$

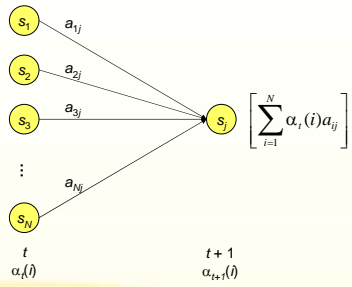
$$1 \leq t \leq T-1$$

$$1 \leq j \leq N$$

- Termination:

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

Forward-Backward Procedure



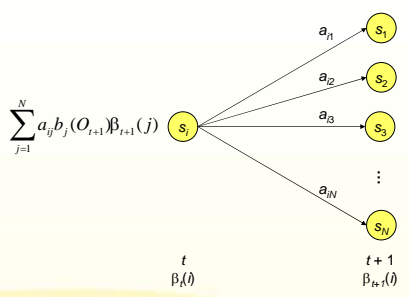
Forward-Backward Procedure

- Define a "backward" variable to be $\beta_t(i) = P(O_{t+1} O_{t+2} \dots O_T | q_t = S_i, M)$, i.e., the probability of a partial observation sequence from $t + 1$ to the end, given state S_i at time t and our model.
 - Initialization: $\beta_T(i) = 1, 1 \leq i \leq N$
 - Induction:
$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)$$

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

$$1 \leq i \leq N$$
- This part is required for the model discovery problem and state sequence determination, not here.

Forward-Backward Procedure



State Sequence Determination

- Solving this problem depends on the definition of “best.”
- Define $\gamma_t(i) = P(q_t = S_i | O, M)$, i.e., the probability of being in state S_i at time t given the observation sequence O and model M .
- This can be expressed entirely in terms of forward and backward variables.

$$\gamma_t(i) = \frac{\alpha_t(i)\beta_t(i)}{P(O|M)} = \frac{\alpha_t(i)\beta_t(i)}{\sum_{j=1}^N \alpha_t(j)\beta_t(j)}$$

State Sequence Determination

- Individually, the most likely state q_t at time t can be determined as

$$q_t = \arg \max_{1 \leq i \leq N} [\gamma_t(i)], \forall t \in [1, T]$$
- Unfortunately, this says nothing about best sequence. In fact, impossible sequences may be found.
- One solution is to maximize expected number of correct pairs of states (q_t, q_{t+1}) .
- Most common approach is to find the single best sequence, i.e., maximize $P(Q|O, M)$, which is equivalent to maximizing $P(Q, O|M)$.
- The *Viterbi Algorithm*, a dynamic programming algorithm, finds this.

Viterbi Algorithm

- Define the “best” score along a single path at time t that accounts for the first t observations and ends in state S_i to be



$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P[q_1, q_2, \dots, q_t = i, O_1, O_2, \dots, O_t | M]$$

- By induction, we have

$$\delta_{t+1}(j) = [\max_i \delta_t(i) a_{ij}] \cdot b_j(O_{t+1})$$



Viterbi Algorithm

- Maintain array $\psi_t(j)$ to keep arguments that maximize δ .
- Initialization: $\delta_1(i) = \pi_i b_i(O_1), \forall i \in [1, N]$
 $\psi_1(i) = 0$
- Recursion: $\delta_t(j) = \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] b_j(O_t), j \in [1, N], t \in [2, T]$
 $\psi_{t1}(j) = \arg \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}]$
- Termination: $P^* = \max_{1 \leq i \leq N} [\delta_T(i)]$
 $q_T^* = \arg \max_{1 \leq i \leq N} [\delta_T(i)]$
- Path Backtracking: $q_t^* = \psi_{t+1}(q_{t+1}^*), t = T-1, T-2, \dots, 1$


70


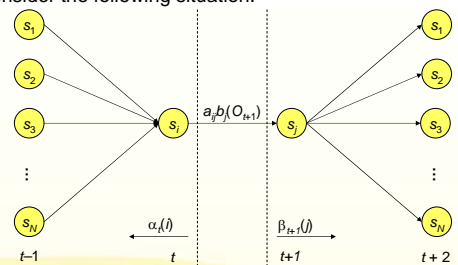
Model Discovery



- This problem involves adjusting the model parameters (A, B, π) to maximize the probability of the observation sequence given the model.
- There is no known analytical solution to this problem.
- Local optimization approach based on “expectation maximization”, called the *Baum-Welch method*.


71


Baum-Welch Method

- Consider the following situation.




72


Baum-Welch Method

- Define $\xi_t(i, j) = P(q_t = S_i, q_{t+1} = S_j \mid O, M)$ to be the probability of being in state S_i at time t and state S_j at time $t+1$, given the observation sequence O and the model M .
- As shown in the figure, this can be defined in terms of forward and backward variables.

$$\begin{aligned} \xi_t(i, j) &= \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{P(O \mid M)} \\ &= \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)} \end{aligned}$$

Baum-Welch Method

- Now recall that we defined $\gamma_t(i)$ to be the probability of being in state S_i at time t , given the observation sequence O and model M .
- We can relate γ to ξ as follows:

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$$
- Then summing $\gamma_t(i)$ over the time index t yields the expected number of times that state S_i is visited (i.e., the expected number of transitions from S_i).
- Similarly, summing $\xi_t(i, j)$ over t is the expected number of transitions from S_i to S_j .

Baum-Welch Method

- Consider now the following:

$$\bar{\pi}_i = \text{expected \# times in state } S_i \text{ at time } (t = 1) = \gamma_1(i)$$

$$\bar{a}_{ij} = \frac{\text{expected \# transitions from } S_i \text{ to } S_j}{\text{expected \# transitions from } S_i}$$

$$= \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

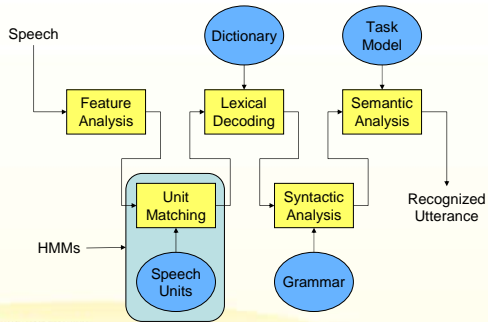
$$\bar{b}_j(k) = \frac{\text{expected \# times in } S_j \text{ observing symbol } v_k}{\text{expected \# times in } S_j}$$

$$= \frac{\sum_{t=1, O_t=v_k}^T \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)}$$

Baum-Welch Method

- Assume we start with model $M = (A, B, \pi)$
- Use this model to compute the right-hand sides of $\bar{\pi}_i, \bar{a}_{ij}, \bar{b}_j(k)$
- This yields re-estimated model $\bar{M} = (\bar{A}, \bar{B}, \bar{\pi})$
- Repeat until convergence.
- Baum and his colleagues proved this process will converge on a fixed point that maximizes $P(O|M)$ for a particular set of training data (i.e., observation sequences).

Continuous Speech Recognition



Word Recognition

- For each word v in the vocabulary, an HMM, M^v is constructed using the Baum-Welch method.
- A word to be recognized is decomposed into a sequence of signals (i.e., observations)—vector quantization.
- Each sequence of observations is processed by each M^v .
- For each M^v , $P(O|M^v)$ is computed using the Viterbi algorithm.
- The word corresponding to $\arg \max P(O|M^v)$ is returned as the most probable word given the observation sequence.
