

## Hidden Markov Models I

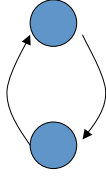
Data Mining

## Learning Bayes Net Structure

Many of these slides are derived from Seyoung Kim, Tom Mitchell, Ziv Bar-Joseph. Thanks!

### What's wrong with Bayesian networks

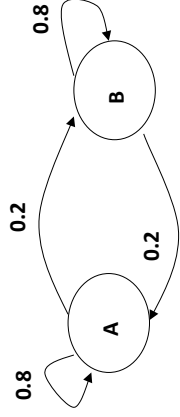
- Bayesian networks are very useful for modeling joint distributions
- But they have their limitations:
  - Cannot account for temporal / sequence models
  - DAG's (no self or any other loops)



This is not a valid Bayesian network!

### Example: Gambling on dice outcome

- Two dices, both skewed (output model).
- Can either stay with the same dice or switch to the second dice (transition mode).



### Hidden Markov models

- Model a set of observation with a set of hidden states
  - Robot movement
    - Observations:** range sensor, visual sensor
    - Hidden states:** location (on a map)
  - Speech processing
    - Observations:** sound signals
    - Hidden states:** parts of speech, words
  - Biology
    - Observations:** amino acid sequence
    - Hidden states:** 3d structure of protein

### Problem setup

- Dice A
  - 1. 0.3
  - 2. 0.2
  - 3. 0.2
  - 4. 0.1
  - 5. 0.1
  - 6. 0.1
- Dice B
  - 1. 0.1
  - 2. 0.1
  - 3. 0.1
  - 4. 0.2
  - 5. 0.2
  - 6. 0.3
- Stay on Dice A: 0.8
- Switch to Dice B: 0.2
- Stay on Dice B: 0.8
- Switch to Dice A: 0.2

Start on dice A with probability 0.3, B with probability 0.7

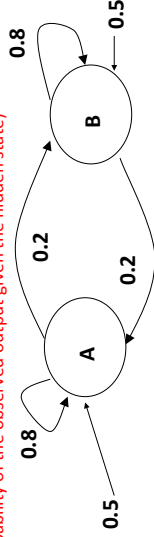
Observed rolls:

1 1 3 4 1 4 4 1 4 5 6 1 2 5 3

Prediction problem: Which rolls were created by which die?

## A Hidden Markov model

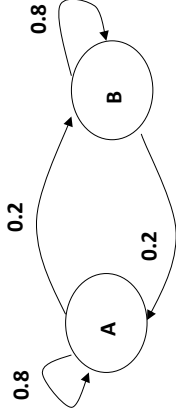
- A set of states  $\{s_1 \dots s_n\}$ 
  - In each time point we are in exactly one of these states denoted by  $q_t$
- $\Pi_t$ , the probability that we *start* at state  $s_i$  (**start die**)
- A transition probability model,  $P(q_t = s_i \mid q_{t-1} = s_j)$  (**switch die**)
- A set of possible outputs  $\Sigma$  (**die roll outcomes**)
  - At time  $t$  we emit a symbol  $\alpha \in \Sigma$  (e.g. at time  $t$  we emit a “1”)
- An emission probability model,  $p(\alpha_t = \alpha \mid s_i)$ 
  - (**probability of the observed output given the hidden state**)



## What can we ask when using a HMM?

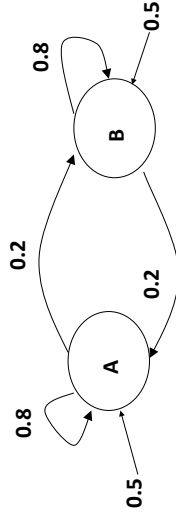
A few examples:

- “What dice is currently being used?”
- “What is the probability of a 6 in the next role?”
- “What is the probability of 6 in any of the next 3 roles?”



## Which die is currently being used?

- We played  $t$  rounds so far
- We want to determine  $P(q_t = A)$
- Lets assume for now that we cannot observe any outputs (we are blind folded)
- How can we compute this?



## HMMs have the Markov Property

An important aspect of this definitions is the Markov property:  $q_{t+1}$  is conditionally independent of  $q_{t-1}$  (and any earlier time points) given  $q_t$

More formally  $P(q_{t+1} = s_i \mid q_t = s_j) = P(q_{t+1} = s_i \mid q_t = s_j, q_{t-1} = s_l)$

## Inference in HMMs

- Computing  $P(Q)$  and  $P(q_t = s_i)$ 
  - If we cannot look at observations
- Computing  $P(Q \mid O)$  and  $P(q_t = s_i \mid O)$ 
  - When we have observation and care about the last state only
- Computing  $\text{argmax}_Q P(Q \mid O)$ 
  - When we care about the entire path

## $P(q_t = A)$ ?

- Simple answer:

Lets determine  $P(Q)$  where  $Q$  is any path that ends in A

$$Q = q_{1'} \dots q_{t-1'}, A$$

$$P(Q) = P(q_{1'}, \dots, q_{t-1'}, A)$$

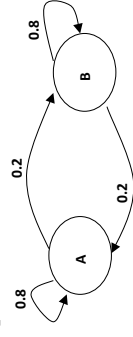
$$= P(A \mid q_{1'}, \dots, q_{t-1'}) P(q_{1'}, \dots, q_{t-1'})$$

$$= P(A \mid q_{t-1'}) P(q_{1'}, \dots, q_{t-1'})$$

$$= \dots$$

$$= P(A \mid q_{t-1'}) \dots P(q_2 \mid q_1) P(q_1)$$

Markov property!



## $P(q_t = A)$ ?

- Simple answer:
  - Lets determine  $P(Q)$  where  $Q$  is any path that ends in  $A$   
 $Q = q_1, \dots, q_{t-1}, A$   
 $P(Q) = P(q_1, \dots, q_{t-1}, A)$   
 $= P(A \mid q_1, \dots, q_{t-1}) P(q_1, \dots, q_{t-1})$   
 $= P(A \mid q_{t-1}) P(q_1, \dots, q_{t-1})$   
 $= \dots$   
 $= P(A \mid q_{t-1}) \dots P(q_2 \mid q_1) P(q_1)$

- $P(q_t = A) = \sum P(Q)$  where the sum is over all sets of  $t$  states that end in  $A$

## $P(q_t = A)$ , the smart way

- Lets define  $p_t(i)$  as the probability of being in state  $i$  at time  $t$ :  $p_t(i) = p(q_t = s_i)$
- We can determine  $p_t(i)$  by induction
  - $p_1(i) = \Pi_i$  (the probability that we *start* at state  $s_i$ )
  - $p_t(i) = ?$

## $P(q_t = A)$ ?

- Simple answer:
  - Lets determine  $P(Q)$  where  $Q$  is any path that ends in  $A$   
 $Q = q_1, \dots, q_{t-1}, A$   
 $P(Q) = P(q_1, \dots, q_{t-1}, A) = P(A \mid q_1, \dots, q_{t-1}) P(q_1, \dots, q_{t-1}) = P(A \mid q_{t-1}) P(q_1, \dots, q_{t-1}) = \dots = P(A \mid q_{t-1}) \dots P(q_2 \mid q_1) P(q_1)$

- $P(q_t = A) = \sum P(Q)$  where the sum is over all sets of  $t$  states that end in  $A$

Q: How many sets  $Q$  are there?  
 A: A lot! ( $2^{t-1}$ )  
 Not a feasible solution

## $P(q_t = A)$ , the smart way

- Lets define  $p_t(i)$  = probability state  $i$  at time  $t = p(q_t = s_i)$
- We can determine  $p_t(i)$  by induction
  - $p_1(i) = \Pi_i$
  - $p_t(i) = \sum_j p(q_t = s_i \mid q_{t-1} = s_j) p_{t-1}(j)$

## $P(q_t = A)$ , the smart way

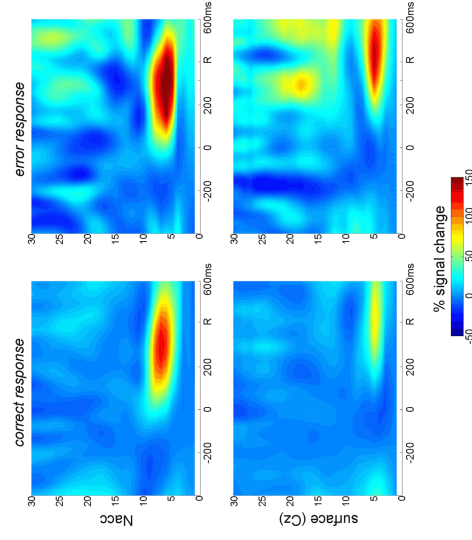
- Lets define  $p_t(i)$  = probability state  $i$  at time  $t = p(q_t = s_i)$
- We can determine  $p_t(i)$  by induction
  - $p_1(i) = \Pi_i$
  - $p_t(i) = \sum_j p(q_t = s_i \mid q_{t-1} = s_j) p_{t-1}(j)$

This type of computation is called dynamic programming  
 Complexity:  $O(n^2 * t)$

Number of states in our HMM

Time / state	t1	t2	t3
s1	.3		
s2	.7		

krigolson@uvic.ca



## Assignment Announcements

- Assignment 1
  - I will accept assignments up until Monday night.
  - **20% off per day on the WHOLE MARK**
  - We can only see the final assignment you hand in on connex
- Assignment 2
  - Out later today
  - No programming
  - **NO LATE DAYS!** Because we want to release the key to help you study for the mid term

## Inference in HMMs

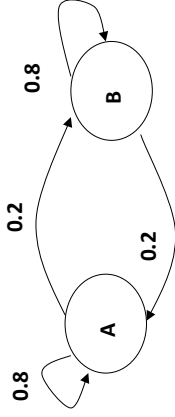
- Computing  $P(Q)$  and  $P(q_t = s_i)$  ✓
- Computing  $P(Q | O)$  and  $P(q_t = s_i | O)$
- Computing  $\text{argmax}_Q P(Q)$

## But what if we observe outputs?

- So far, we assumed that we could not observe the outputs
- In reality, we almost always can.



v	$P(v   A)$	$P(v   B)$
1	.3	.1
2	.2	.1
3	.2	.1
4	.1	.2
5	.1	.2
6	.1	.3

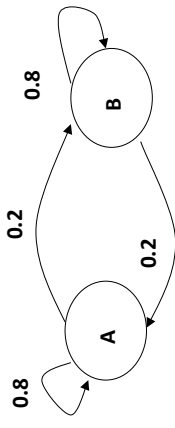


v	$P(v   A)$	$P(v   B)$
1	.3	.1
2	.2	.1
3	.2	.1
4	.1	.2
5	.1	.2
6	.1	.3

Does observing the sequence

5, 6, 4, 5, 6, 6

Change our belief about the state?



## $P(q_t = A)$ when outputs are observed

- We want to compute  $P(q_t = A | O_1 \dots O_t)$
- For ease of writing we will use the following notations (commonly used in the literature)

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

$$b_i(o_t) = P(o_t | s_i)$$

Transition probability

Emission probability

## $P(q_t = A)$ when outputs are observed

- We want to compute  $P(q_t = A | O_1 \dots O_t)$
- Lets start with a simpler question. Given a sequence of states  $Q$ , what is  $P(Q | O_1 \dots O_t) = P(Q | O)$ ?
  - It is pretty simple to move from  $P(Q)$  to  $P(q_t = A)$
  - In some cases  $P(Q)$  is the more important question
    - Speech processing
    - NLP

## P(Q | O)

- We can use Bayes rule:

$$P(Q|O) = \frac{P(O|Q)P(Q)}{P(O)}$$

Easy,  $P(O|Q) = P(o_1 | q_1) P(o_2 | q_2) \dots P(o_t | q_t)$

## P(Q | O)

- We can use Bayes rule:

$$P(Q|O) = \frac{P(O|Q)P(Q)}{P(O)}$$

Easy,  $P(Q) = P(q_1) P(q_2 | q_1) \dots P(q_t | q_{t-1})$

## P(Q | O)

- We can use Bayes rule:

$$P(Q|O) = \frac{P(O|Q)P(Q)}{P(O)}$$

Hard!

## P(O)

- What is the probability of seeing a set of observations:
  - An important question in it own rights, for example classification using two HMMs
- Define  $\alpha_t(i) = P(o_1, o_2, \dots, o_t \wedge q_t = s_i)$
- $\alpha_t(i)$  is the probability that we:
  - Observe  $o_1, o_2, \dots, o_t$
  - End up at state  $i$

How do we compute  $\alpha_t(i)$ ?

## Computing $\alpha_t(i)$

- $\alpha_t(i) = P(o_1 \wedge q_1 = i) = P(o_1 | q_1 = s_i) \Pi_i$

$$\alpha_{t+1}(i) = P(o_1 \dots o_{t+1} \wedge q_{t+1} = s_i) = \sum_j P(o_1 \dots o_t \wedge q_t = s_j \wedge o_{t+1} \wedge q_{t+1} = s_i) = \sum_j P(o_{t+1} \wedge q_{t+1} = s_i | o_1 \dots o_t \wedge q_t = s_j) P(o_1 \dots o_t \wedge q_t = s_j) = \sum_j P(o_{t+1} \wedge q_{t+1} = s_i | o_1 \dots o_t \wedge q_t = s_j) \alpha_t(j) = \sum_j P(o_{t+1} | q_{t+1} = s_i) P(q_{t+1} = s_i | q_t = s_j) \alpha_t(j) =$$

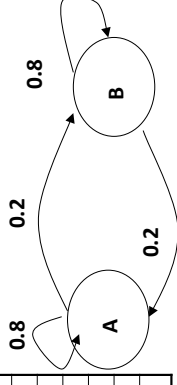
Markov property

## Example: Computing $\alpha_3(B)$

- We observed 2,3,6
- $$\alpha_1(A) = P(2 \wedge q_1 = A) = P(2 | q_1 = A) \Pi_A = 2 * .7 = .14, \alpha_1(B) = .1 * .3 = .03$$
- $$\alpha_2(A) = \sum_{j=A,B} a_{j,A}(3) \alpha_1(j) = 2 * .8 * .14 + 2 * .2 * .03 = 0.0236, \alpha_2(B) = 0.0052$$
- $$\alpha_3(B) = \sum_{j=A,B} b_{j,B}(6) \alpha_2(j) = 3 * .2 * 0.0236 + 3 * .8 * 0.0052 = 0.00264$$

$\Pi_A = 0.7$   
 $\Pi_B = 0.3$

v	P(v   A)	P(v   B)
1	.3	.1
2	.2	.1
3	.2	.1
4	.1	.2
5	.1	.2
6	.1	.3



## Where we are

- We want to compute  $P(Q | O)$
- For this, we only need to compute  $P(O)$
- We know how to compute  $\alpha_t(i)$

From now its easy

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

so

$$P(O) = P(o_1, o_2, \dots, o_t) = \sum_i P(o_1, o_2, \dots, o_t \wedge q_t = s_i) = \sum_i \alpha_t(i)$$

note that

$$p(q_t = s_i | o_1, o_2, \dots, o_t) = \frac{\alpha_t(i)}{\sum_j \alpha_t(j)}$$

$$P(A | B) = P(A \wedge B) / P(B)$$

## Most probable path

- We are almost done ...
- One final question remains  
How do we find the most probable path, that is  $Q^*$  such that  
 $P(Q^* | O) = \arg \max_Q P(Q | O)$ ?

- This is an important path
  - The words in speech processing
  - The set of genes in the genome
  - etc.

## Most probable path

$$\begin{aligned} \arg \max_Q P(Q | O) &= \arg \max_Q \frac{P(O | Q) P(Q)}{P(O)} \\ &= \arg \max_Q P(O | Q) P(Q) \end{aligned}$$

We will use the following definition:

$$\delta_i(i) = \max_{q_1 \dots q_{i-1}} p(q_1 \dots q_{i-1} \wedge q_i = s_i \wedge O_1 \dots O_i)$$

In other words we are interested in the most likely path from 1 to  $t$  that:

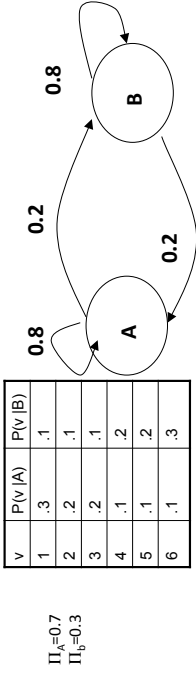
1. Ends in  $S_i$
2. Produces outputs  $O_1 \dots O_i$

## Inference in HMMs

- Computing  $P(Q)$  and  $P(q_t = s_i)$  ✓
- Computing  $P(Q | O)$  and  $P(q_t = s_i | O)$  ✓
- Computing  $\arg \max_Q P(Q)$

## Example

- What is the most probable set of states leading to the sequence:  
1, 2, 2, 5, 6, 5, 1, 2, 3 ?



## Computing $\delta_t(i)$

$$\begin{aligned} \delta_t(i) &= p(q_t = s_i \wedge O_t) \\ &= p(q_t = s_i) p(O_t | q_t = s_i) \\ &= \pi_i b_i(O_t) \end{aligned}$$

$$\delta_t(i) = \max_{q_1 \dots q_{t-1}} p(q_1 \dots q_{t-1} \wedge q_t = s_i \wedge O_1 \dots O_t)$$

Q: Given  $\delta_t(i)$ , how can we compute  $\delta_{t+1}(i)$ ?

A: To get from  $\delta_t(i)$  to  $\delta_{t+1}(i)$  we need to

1. Add an emission for time  $t+1$  ( $O_{t+1}$ )
2. Transition to state  $s_i$

$$\begin{aligned} \delta_{t+1}(i) &= \max_{q_1 \dots q_t} p(q_1 \dots q_t \wedge q_{t+1} = s_i \wedge O_1 \dots O_{t+1}) \\ &= \max_j \delta_t(j) p(q_{t+1} = s_i | q_t = s_j) p(O_{t+1} | q_{t+1} = s_i) \\ &= \max_j \delta_t(j) a_{ji} b_i(O_{t+1}) \end{aligned}$$

## The Viterbi algorithm

$$\begin{aligned}\delta_{t+1}(i) &= \max_{q_1 \in \mathcal{Q}} p(q_1 | \mathbb{K} \ q_t \wedge q_{t+1} = s_i \wedge O_1 \dots O_{t+1}) \\ &= \max_j \delta_t(j) p(q_{t+1} = s_i | q_t = s_j) p(O_{t+1} | q_{t+1} = s_i) \\ &= \max_j \delta_t(j) a_{j,i} b_i(O_{t+1})\end{aligned}$$

- Once again we use dynamic programming for solving  $\delta_t(i)$

- Once we have  $\delta_t(i)$ , we can solve for our  $P(Q^* | O)$

By:

$$P(Q^* | O) = \operatorname{argmax}_Q P(Q | O) =$$

path defined by  $\operatorname{argmax}_i \delta_t(i)$ ,

## Wikipedia Page for Viterbi is Great

- [https://en.wikipedia.org/wiki/Viterbi\\_algorithm](https://en.wikipedia.org/wiki/Viterbi_algorithm)
- A nice animation
  - [https://en.wikipedia.org/wiki/Viterbi\\_algorithm#/media/File:Viterbi\\_animated\\_demo.gif](https://en.wikipedia.org/wiki/Viterbi_algorithm#/media/File:Viterbi_animated_demo.gif)

## Inference in HMMs

- Computing  $P(Q)$  and  $P(q_t = s_i)$  ✓
- Computing  $P(Q | O)$  and  $P(q_t = s_i | O)$  ✓
- Computing  $\operatorname{argmax}_Q P(Q)$  ✓

## What you should know

- Why HMMs? Which applications are suitable?
- Inference in HMMs
  - No observations
  - Probability of next state w. observations
  - Maximum scoring path (Viterbi)