

Hidden Markov Models

CISC 5800
Professor Daniel Leeds

Representing sequence data



- Spoken language
- DNA sequences
- Daily stock values

Example: spoken language

F?r plu? fi?e is nine

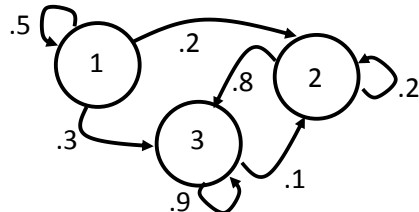
- Between F and r expect a vowel: "aw", "ee", "ah"; NOT "oh", "uh"
- At end of "plu" expect consonant: "g", "m", "s"; NOT "d", "p"

2

Markov Models

Start with:

- n states: s_1, \dots, s_n
- Probability of initial start states: Π_1, \dots, Π_n
- Probability of transition between states: $A_{i,j} = P(q_t=s_j | q_{t-1}=s_i)$

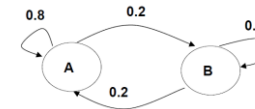


3

A dice-y example

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

- Two colored die



- What is the probability we start at s_A ?
- What is the probability we have the sequence of die choices:

$$s_A, s_A?$$

- What is the probability we have the sequence of die choices:

$$s_B, s_A, s_B, s_A?$$

4

A dice-y example

- What is the probability we have the sequence of die choices: s_B, s_A, s_B, s_A ?
- $\Pi_A = 0.3, \Pi_B = 0.7$
- Dynamic programming: find answer for q_t , then compute q_{t+1}

State\Time	t_1	t_2	t_3
s_A	0.3		
s_B	0.7		

$$p_t(i) = \sum_j p(q_t = s_i | q_{t-1} = s_j) p_{t-1}(j)$$

$p_t(i) = P(q_t = s_i)$ -- Probability state i at time t

Hidden Markov Models

- Actual state q "hidden"
- State produces visible data o : $\phi_{i,j} = P(o_t = x_i | q_t = s_j)$
- Compute

$$P(O, Q | \theta) = p(q_1 | \pi) \prod_{t=2}^T p(q_t | q_{t-1}, A) \prod_{t=1}^T p(o_t | q_t, \phi)$$

Deducing die based on observed "emissions"

Each color is biased

o	$P(o s_A)$	$P(o s_B)$
1	.3	.1
2	.2	.1
3	.2	.1
4	.1	.2
5	.1	.2
6	.1	.3

Intuition – balance transition and emission probabilities

Observed numbers: 554565254556 – the 2 is probably from s_B

Observed numbers: 554565213321 – the 2 is probably from s_A

Deducing die based on observed "emissions"

Each color is biased

o	$P(o s_R)$	$P(o s_B)$
1	.3	.1
2	.2	.1
3	.2	.1
4	.1	.2
5	.1	.2
6	.1	.3

- We see: 5 What is probability of $o=5, q=B$ (blue)
 $\Pi_B \phi_{5,B} = 0.7 \times 0.2 = 0.14$
- We see: 5, 3 What is probability of $o=5,3 | q=B, B$?
 $\Pi_B \phi_{5,B} A_{B,B} \phi_{3,B} = 0.7 \times 0.2 \times 0.8 \times 0.1 = 0.0112$

Goal: calculate most likely states given observable data

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

Define and use $\delta_t(i)$

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

$\delta_t(i)$: max possible value of $P(q_1, \dots, q_t, o_1, \dots, o_t)$ given we insist $q_t = s_i$

Find the most likely path from q_1 to q_t that

- $q_t = s_i$
- Outputs are o_1, \dots, o_t

9

Viterbi algorithm: $\delta_t(i)$

$$\delta_1(i) = \pi_i P(o_1 | q_1 = s_i) = \pi_i \phi_{1,i}$$

$$\delta_t(i) = \max_j \delta_{t-1}(j) P(q_t = s_i | q_{t-1} = s_j) P(o_t | q_t = s_i) = \max_j \delta_{t-1}(j) \phi_{t,i} A_{t,j}$$

$$P(Q^* | O) = \arg \max_Q P(Q | O) = \arg \max_i \delta_t(i)$$

10

Viterbi algorithm: bigger picture

Compute all $\delta_t(i)$'s

- At time $t=1$ compute $\delta_1(i)$ for every state i
- At time $t=2$ compute $\delta_2(i)$ for every state i (based on $\delta_1(i)$ values)
- ...
- At time $t=T$ compute $\delta_T(i)$ for every state i (based on $\delta_{T-1}(i)$ values)

Find states going from $t=T$ back to $t=1$ to lead to max $\delta_T(i)$

- Now find state j that gives maximum value for $\delta_T(j)$
- Find state k at time $T-1$ used to maximize $\delta_T(j)$
- ...
- Find state z at time 1 used to maximize $\delta_2(y)$

11

Parameters in HMM

Initial probabilities: π_i

Transition probabilities $A_{i,j}$

Emission probabilities $\phi_{i,j}$

How do we learn these values?

12

First, assume we know the states

Learning HMM parameters: π_i

Compute MLE for each parameter

\mathbf{x}^1 : A B, A, A, B
 \mathbf{x}^2 : B, B, B, A, A
 \mathbf{x}^3 : A, A, B, A, B
 \vdots

$$\pi^* = \operatorname{argmax}_{\pi} \prod_k \pi(q_1) \prod_{t=2}^T p(q_t | q_{t-1}) \prod_{t=1}^T p(o_t | q_t, \phi)$$

$$\pi_A = \frac{\#D(q_1 = s_A)}{\#D}$$

13

First, assume we know the states

Learning HMM parameters: $A_{i,j}$

Compute MLE for each parameter

\mathbf{x}^1 : A, B, A, A, B
 \mathbf{x}^2 : B, B, B, A, A
 \mathbf{x}^3 : A, A, B, A, B
 \vdots

$$A^* = \operatorname{argmax}_A \prod_k \pi(q_1) \prod_{t=2}^T p(q_t | q_{t-1}) \prod_{t=1}^T p(o_t | q_t, \phi)$$

$$A_{i,j} = \frac{\#D(q_t = s_i, q_{t-1} = s_j)}{\#D(q_{t-1} = s_j)}$$

14

First, assume we know the states

Learning HMM parameters: $\phi_{i,j}$

Compute MLE for each parameter

\mathbf{x}^1 : A, B, A, A, B
 \mathbf{o}^1 : 2, 5, 3, 3, 6
 \mathbf{x}^2 : B, B, B, A, A
 \mathbf{o}^2 : 4, 5, 1, 3, 2
 \mathbf{x}^3 : A, A, B, A, B
 \mathbf{o}^3 : 1, 4, 5, 2, 6
 \vdots

$$\phi^* = \operatorname{argmax}_{\phi} \prod_k \pi(q_1) \prod_{t=2}^T p(q_t | q_{t-1}) \prod_{t=1}^T p(o_t | q_t, \phi)$$

$$\phi_{i,j} = \frac{\#D(o_t = i, q_t = s_j)}{\#D(q_t = s_j)}$$


15

Challenges in HMM learning

Learning parameters (π, A, ϕ) with known states is not too hard

BUT usually states are unknown

If we had the parameters and the observations, we could figure out the states: Viterbi $P(Q^* | O) = \operatorname{argmax}_Q P(Q | O)$



16

Expectation-Maximization, or “EM”

Problem: Uncertain of y^i (class), uncertain of θ^i (parameters)

Solution: Guess y^i , deduce θ^i , re-compute y^i , re-compute θ^i ... etc.

OR: Guess θ^i , deduce y^i , re-compute θ^i , re-compute y^i

Will converge to a solution

E step: Fill in expected values for missing labels y

M step: Regular MLE for θ given known and filled-in variables

Also useful when there are holes in your data

17

Computing states q_t

Instead of picking one state: $q_t = s_i$, find $P(q_t = s_i | \mathbf{o})$

$$P(q_t = s_i | o_1, \dots, o_T) = \frac{\alpha_t(i)\beta_t(i)}{\sum_j \alpha_t(j)\beta_t(j)}$$

Forward probability: $\alpha_t(i) = P(o_1 \dots o_t \wedge q_t = s_i)$

Backward probability: $\beta_t(i) = P(o_{t+1} \dots o_T | q_t = s_i)$

18

Details of forward probability

Forward probability: $\alpha_t(i) = P(o_1 \dots o_t \wedge q_t = s_i)$

$$\alpha_1(i) = \phi_{o_1,i} \pi_i = P(o_1 | q_1 = s_i) P(q_1 = s_i)$$

$$\alpha_t(i) = \phi_{o_t,i} \sum_j A_{i,j} \alpha_{t-1}(j)$$

$$\alpha_t(i) = P(o_t | q_t = s_i) \sum_j P(q_t = s_i | q_{t-1} = s_j) \alpha_{t-1}(j)$$

20

Details of backward probability

Backward probability: $\beta_t(i) = P(o_{t+1} \dots o_T | q_t = s_i)$

$$\beta_t(i) = \sum_j A_{j,i} \phi_{o_{t+1},j} \beta_{t+1}(j)$$

$$\beta_t(i) = \sum_j P(q_{t+1} = s_j | q_t = s_i) P(o_{t+1} | q_{t+1} = s_j) \beta_{t+1}(j)$$

Final β : $\beta_{T-1}(i)$

$$\beta_{T-1}(i) = \sum_j A_{j,i} \phi_{o_T,j}$$

$$= P(q_T = s_j | q_{T-1} = s_i) P(o_T | q_T = s_j)$$

21

E-step: State probabilities

One state:

$$P(q_t = s_i | o_1, \dots, o_T) = \frac{\alpha_t(i)\beta_t(i)}{\sum_j \alpha_t(j)\beta_t(j)} = S_t(i)$$

Two states in a row:

$$P(q_t = s_j, q_{t+1} = s_i | o_1, \dots, o_T) = \frac{\alpha_t(j)A_{i,j}\phi_{o_{t+1},i}\beta_{t+1}(i)}{\sum_i \sum_j \alpha_t(j)A_{i,j}\phi_{o_{t+1},i}\beta_{t+1}(i)} = S_t(i, j)$$

22

Recall: when states known

$$\pi_A = \frac{\#D(q_1=s_A)}{\#D}$$

$$A_{i,j} = \frac{\#D(q_t=s_i, q_{t-1}=s_j)}{\#D(q_{t-1}=s_j)}$$

$$\phi_{i,j} = \frac{\#D(o_t=i)}{\#D(q_t=s_j)}$$

23

M-step

$$A_{i,j} = \frac{\sum_t S_t(i,j)}{\sum_t S_t(i)}$$

$$\phi_{obs,i} = \frac{\sum_t \mathbb{1}_{o_t=obs} S_t(i)}{\sum_t S_t(i)}$$

$$\pi_i = S_1(i)$$

Known states:

$$\bullet \pi_A = \frac{\#D(q_1=s_A)}{\#D}$$

$$\bullet A_{i,j} = \frac{\#D(q_t=s_i, q_{t-1}=s_j)}{\#D(q_{t-1}=s_j)}$$

$$\bullet \phi_{i,j} = \frac{\#D(o_t=i)}{\#D(q_t=s_j)}$$

24

Review of HMMs in action

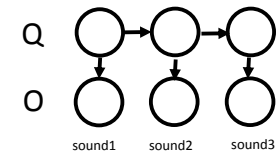
For classification, find highest probability class given features

Features for one sound:

- $[q_1, o_1, q_2, o_2, \dots, q_T, o_T]$

Conclude word:

Generates states:



26