

Lecture Notes on
Pattern Recognition

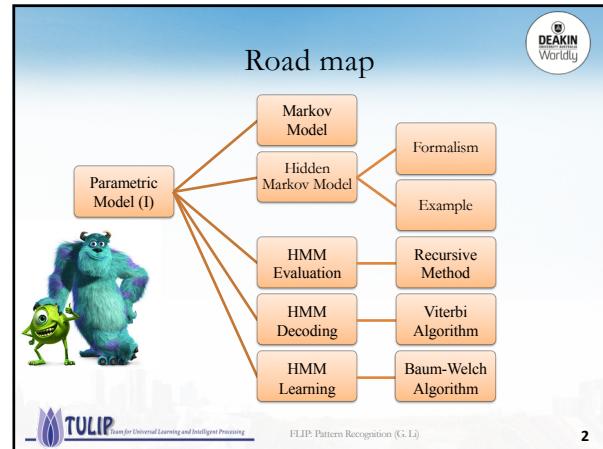
Session 04(A): Parametric Model (I)

Gang Li
School of Information Technology
Deakin University, VIC 3125, Australia

TULIP Team for Universal Learning and Intelligent Processing

DEAKIN Worldly

2



Markov Model

- First Order Markov Model
- Limitations

TULIP Team for Universal Learning and Intelligent Processing

DEAKIN Worldly

3

Sequential and Temporal Patterns

- Some applications have an inherent temporality.
 - Speech recognition
 - Gesture recognition
 - Human activity recognition

TULIP Team for Universal Learning and Intelligent Processing

DEAKIN Worldly

4

Sequential and Temporal Patterns

- **Temporal** patterns:
 - The order of data points is important
 - (i.e., time series).
 - Data can be represented by a number of states.
 - States at time t are influenced directly by states in previous time steps
 - (i.e., correlated).
- **Sequential** patterns:
 - The order of data points is irrelevant.

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

5

First Order Markov Model

Markov Model

- States: $s_1, s_2 \dots s_N$
- Time-steps: $t=0, 1, \dots$

TULIP Team for Universal Learning and Intelligent Processing

DEAKIN Worldly

6

First Order Markov Model

Markov Model

- States: $s_1, s_2 \dots s_N$
- Time-steps: $t=0, 1, \dots$
- On the t -th time-step the system is in exactly one of the available states q_t

FLIP: Pattern Recognition (G. Li) 7

First Order Markov Model

Markov Model

- States: $s_1, s_2 \dots s_N$
- Time-steps: $t=0, 1, \dots$
- On the t -th time-step the system is in exactly one of the available states q_t
- The next state is chosen randomly by a distribution based on current status.

FLIP: Pattern Recognition (G. Li) 8

First Order Markov Model

Markov Model

- States: $s_1, s_2 \dots s_N$
- Time-steps: $t=0, 1, \dots$
- On the t -th time-step the system is in exactly one of the available states q_t
- The next state is chosen randomly by a distribution based on current status.

FLIP: Pattern Recognition (G. Li) 9

First Order Markov Model

Markov Model

- States: $s_1, s_2 \dots s_N$
- Time-steps: $t=0, 1, \dots$
- On the t -th time-step the system is in exactly one of the available states q_t
- The next state is chosen randomly by a distribution based on current status.

FLIP: Pattern Recognition (G. Li) 10

First Order Markov Model

Markov Property

- q_{t+1} is conditionally independent of $\{q_{t-1}, q_{t-2}, \dots, q_1, q_0\}$ given q_t .
- In other words:
 - $P(q_{t+1} = s_i | q_t = s_i) = P(q_{t+1} = s_j | q_t = s_i, \text{any earlier history})$
 - $a_{ij} = P(q_{t+1} = s_j | q_t = s_i)$

FLIP: Pattern Recognition (G. Li) 11

First Order Markov Model

Markov Property

- q_{t+1} is conditionally independent of $\{q_{t-1}, q_{t-2}, \dots, q_1, q_0\}$ given q_t .
- In other words:
 - $P(q_{t+1} = s_j | q_t = s_i) = P(q_{t+1} = s_j | q_t = s_i, \text{any earlier history})$
 - $a_{ij} = P(q_{t+1} = s_j | q_t = s_i)$

FLIP: Pattern Recognition (G. Li) 12

First Order Markov Model (Weather Prediction Model)

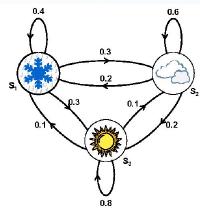
- Assume three weather states:

- S_1 : Precipitation (rain, snow, hail, etc.)
- S_2 : Cloudy
- S_3 : Sunny

- Transition Matrix

$$\begin{matrix} & S_1 & S_2 & S_3 \\ S_1 & 0.4 & 0.3 & 0.3 \\ S_2 & 0.2 & 0.6 & 0.2 \\ S_3 & 0.1 & 0.1 & 0.8 \end{matrix}$$

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



FLIP: Pattern Recognition (G. Li)

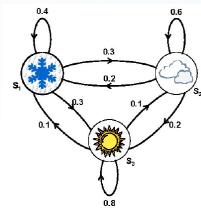
13

First Order Markov Model (Weather Prediction Model)

- The probability of a sequence of states $q^T = (q_1, q_2, \dots, q_T)$ is the product of the corresponding transition probabilities

- “sunny-sunny-sunny-rainy-rainy-sunny-cloudy-sunny”?

$$= p(S_3)a_{33}a_{33}a_{31}a_{11}a_{13}a_{32}a_{23}$$



FLIP: Pattern Recognition (G. Li)

14

Limitations of Markov models

- In Markov models, each state is uniquely associated with an observable event.
 - Once an observation is made, the state of the system is trivially retrieved.
 - Such systems are not of practical use for most applications.

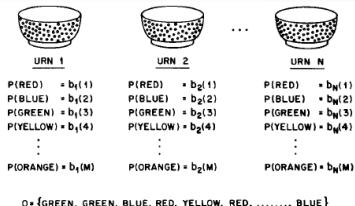


FLIP: Pattern Recognition (G. Li)

15

Limitations of Markov models

- From the sequence of balls you observed, figure out the sequence of URNs the balls are taken from?

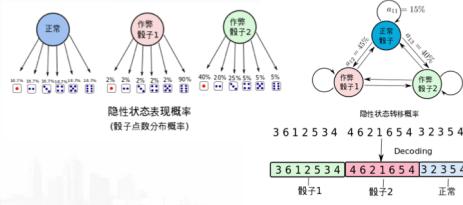


FLIP: Pattern Recognition (G. Li)

16

Limitations of Markov models

- Given a sequence of dice rolls, assuming that we know that the casino player may switch between a fair and two loaded die:



FLIP: Pattern Recognition (G. Li)

17

Hidden Markov Model

- HMM Formalism
- HMM Example



FLIP: Pattern Recognition (G. Li)

18

Hidden Markov Models

First Order HMM

- Augment Markov model such that when it is in a **hidden state** $\omega(t)$ it also emits some **visible symbol** $v(t)$ among a set of possible symbols.
- We have access to the **visible states** only, while **hidden states** are unobservable.

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 19

Hidden Markov Models

- Hidden states**: the (TRUE) states of a system that may be described by a Markov process
 - e.g., the weather
- Observable states**: the states of the process that are 'visible'
 - e.g., seaweed dampness

Hidden State	Soggy	Damp	Dryish	Dry
Sun	0.60	0.20	0.15	0.05
Cloud	0.25	0.25	0.35	0.25
Rain	0.05	0.10	0.35	0.50

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 20

Hidden Markov Models

- Initial Distribution**: contains the probability of the (hidden) model being in a particular initial hidden state.
- State transition matrix**: holding the probability of a hidden state given the previous hidden state.
- Output matrix**: containing the probability of observing a particular observable state given that the hidden model is in a particular hidden state.

Hidden State	Soggy	Damp	Dryish	Dry
Sun	0.60	0.20	0.15	0.05
Cloud	0.25	0.25	0.35	0.25
Rain	0.05	0.10	0.35	0.50

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 21

HMM Formalism

- An HMM is defined by $\lambda = \{\Omega, V, P, A, B\}$:
 - $\Omega = \{S_1, \dots, S_n\}$ are n possible **states**
 - $V = \{O_1, \dots, O_m\}$ are m possible **observations**
 - $P = \{\pi_1, \dots, \pi_n\}$ are the **prior state probabilities**
 - $A = \{a_{ij}\}$ is the $(n \times n)$ **state transition matrix**
 - $B = \{b_{ik}\}$ is the $(n \times m)$ **observation state matrix**

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 22

An HMM Example

- Start randomly in state S_1 or S_2
- Choose one of the output symbols in each state at random.

$n = 3, m = 3$
 $\pi_1 = 1/2 \quad \pi_2 = 1/2 \quad \pi_3 = 0$
 $a_{11} = 0 \quad a_{12} = 1/3 \quad a_{13} = 2/3$
 $a_{21} = 1/3 \quad a_{22} = 0 \quad a_{23} = 2/3$
 $a_{31} = 1/3 \quad a_{32} = 1/3 \quad a_{33} = 1/3$
 $b_1(X) = 1/2 \quad b_1(Y) = 1/2 \quad b_1(Z) = 0$
 $b_2(X) = 0 \quad b_2(Y) = 1/2 \quad b_2(Z) = 1/2$
 $b_3(X) = 1/2 \quad b_3(Y) = 0 \quad b_3(Z) = 1/2$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 23

An HMM Example

- Start randomly in state S_1 or S_2
- Choose one of the output symbols in each state at random.
- Let us generate a sequence of observations

$n = 3, m = 3$
 $\pi_1 = 1/2 \quad \pi_2 = 1/2 \quad \pi_3 = 0$
 $a_{11} = 0 \quad a_{12} = 1/3 \quad a_{13} = 2/3$
 $a_{21} = 1/3 \quad a_{22} = 0 \quad a_{23} = 2/3$
 $a_{31} = 1/3 \quad a_{32} = 1/3 \quad a_{33} = 1/3$
 $b_1(X) = 1/2 \quad b_1(Y) = 1/2 \quad b_1(Z) = 0$
 $b_2(X) = 0 \quad b_2(Y) = 1/2 \quad b_2(Z) = 1/2$
 $b_3(X) = 1/2 \quad b_3(Y) = 0 \quad b_3(Z) = 1/2$

$q_0 =$	•	$O_0 =$	
$q_1 =$		$O_1 =$	
$q_2 =$		$O_2 =$	

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 24

An HMM Example

- Start randomly in state S_1 or S_2
- Choose one of the output symbols in each state at random.
- Let us generate a sequence of observations

$n = 3, m = 3$
 $\pi_1 = 1/2 \quad \pi_2 = 1/2 \quad \pi_3 = 0$

$a_{11} = 0$	$a_{12} = 1/3$	$a_{13} = 2/3$
$a_{21} = 1/3$	$a_{22} = 0$	$a_{23} = 2/3$
$a_{31} = 1/3$	$a_{32} = 1/3$	$a_{33} = 1/3$

$b_1(X) = 1/2$	$b_1(Y) = 1/2$	$b_1(Z) = 0$
$b_2(X) = 0$	$b_2(Y) = 1/2$	$b_2(Z) = 1/2$
$b_3(X) = 1/2$	$b_3(Y) = 0$	$b_3(Z) = 1/2$

50-50 choice between X and Y

$q_0 =$	S_1	$O_0 =$	---
$q_1 =$	---	$O_1 =$	---
$q_2 =$	---	$O_2 =$	---

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

25

An HMM Example

- Start randomly in state S_1 or S_2
- Choose one of the output symbols in each state at random.
- Let us generate a sequence of observations

$n = 3, m = 3$
 $\pi_1 = 1/2 \quad \pi_2 = 1/2 \quad \pi_3 = 0$

$a_{11} = 0$	$a_{12} = 1/3$	$a_{13} = 2/3$
$a_{21} = 1/3$	$a_{22} = 0$	$a_{23} = 2/3$
$a_{31} = 1/3$	$a_{32} = 1/3$	$a_{33} = 1/3$

$b_1(X) = 1/2$	$b_1(Y) = 1/2$	$b_1(Z) = 0$
$b_2(X) = 0$	$b_2(Y) = 1/2$	$b_2(Z) = 1/2$
$b_3(X) = 1/2$	$b_3(Y) = 0$	$b_3(Z) = 1/2$

Goto S_1 with probability 2/3 or S_2 with prob. 1/3

$q_0 =$	S_1	$O_0 =$	X
$q_1 =$	---	$O_1 =$	---
$q_2 =$	---	$O_2 =$	---

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

26

An HMM Example

- Start randomly in state S_1 or S_2
- Choose one of the output symbols in each state at random.
- Let us generate a sequence of observations

$n = 3, m = 3$
 $\pi_1 = 1/2 \quad \pi_2 = 1/2 \quad \pi_3 = 0$

$a_{11} = 0$	$a_{12} = 1/3$	$a_{13} = 2/3$
$a_{21} = 1/3$	$a_{22} = 0$	$a_{23} = 2/3$
$a_{31} = 1/3$	$a_{32} = 1/3$	$a_{33} = 1/3$

$b_1(X) = 1/2$	$b_1(Y) = 1/2$	$b_1(Z) = 0$
$b_2(X) = 0$	$b_2(Y) = 1/2$	$b_2(Z) = 1/2$
$b_3(X) = 1/2$	$b_3(Y) = 0$	$b_3(Z) = 1/2$

50-50 choice between X and Z

$q_0 =$	S_1	$O_0 =$	X
$q_1 =$	S_3	$O_1 =$	---
$q_2 =$	---	$O_2 =$	---

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

27

An HMM Example

- Start randomly in state S_1 or S_2
- Choose one of the output symbols in each state at random.
- Let us generate a sequence of observations

$n = 3, m = 3$
 $\pi_1 = 1/2 \quad \pi_2 = 1/2 \quad \pi_3 = 0$

$a_{11} = 0$	$a_{12} = 1/3$	$a_{13} = 2/3$
$a_{21} = 1/3$	$a_{22} = 0$	$a_{23} = 2/3$
$a_{31} = 1/3$	$a_{32} = 1/3$	$a_{33} = 1/3$

$b_1(X) = 1/2$	$b_1(Y) = 1/2$	$b_1(Z) = 0$
$b_2(X) = 0$	$b_2(Y) = 1/2$	$b_2(Z) = 1/2$
$b_3(X) = 1/2$	$b_3(Y) = 0$	$b_3(Z) = 1/2$

Each of the three next states is equally likely

$q_0 =$	S_1	$O_0 =$	X
$q_1 =$	S_3	$O_1 =$	X
$q_2 =$	---	$O_2 =$	---

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

28

An HMM Example

- Start randomly in state S_1 or S_2
- Choose one of the output symbols in each state at random.
- Let us generate a sequence of observations

$n = 3, m = 3$
 $\pi_1 = 1/2 \quad \pi_2 = 1/2 \quad \pi_3 = 0$

$a_{11} = 0$	$a_{12} = 1/3$	$a_{13} = 2/3$
$a_{21} = 1/3$	$a_{22} = 0$	$a_{23} = 2/3$
$a_{31} = 1/3$	$a_{32} = 1/3$	$a_{33} = 1/3$

$b_1(X) = 1/2$	$b_1(Y) = 1/2$	$b_1(Z) = 0$
$b_2(X) = 0$	$b_2(Y) = 1/2$	$b_2(Z) = 1/2$
$b_3(X) = 1/2$	$b_3(Y) = 0$	$b_3(Z) = 1/2$

50-50 choice between X and Z

$q_0 =$	S_1	$O_0 =$	---
$q_1 =$	S_3	$O_1 =$	---
$q_2 =$	S_3	$O_2 =$	---

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

29

An HMM Example

- Start randomly in state S_1 or S_2
- Choose one of the output symbols in each state at random.
- Let us generate a sequence of observations

$n = 3, m = 3$
 $\pi_1 = 1/2 \quad \pi_2 = 1/2 \quad \pi_3 = 0$

$a_{11} = 0$	$a_{12} = 1/3$	$a_{13} = 2/3$
$a_{21} = 1/3$	$a_{22} = 0$	$a_{23} = 2/3$
$a_{31} = 1/3$	$a_{32} = 1/3$	$a_{33} = 1/3$

$b_1(X) = 1/2$	$b_1(Y) = 1/2$	$b_1(Z) = 0$
$b_2(X) = 0$	$b_2(Y) = 1/2$	$b_2(Z) = 1/2$
$b_3(X) = 1/2$	$b_3(Y) = 0$	$b_3(Z) = 1/2$

This is what the observer has to work with...

$q_0 =$	S_1	$O_0 =$	X
$q_1 =$	S_3	$O_1 =$	X
$q_2 =$	S_3	$O_2 =$	Z

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

30

Another HMM Example

• Given a sequence of dice rolls, assuming that we know that the casino player may switch between a fair and two loaded die:

FLIP: Pattern Recognition (G. Li)

31

Another HMM Example

• Given a sequence of dice rolls, assuming that we know that the casino player may switch between a fair and two loaded die:

FLIP: Pattern Recognition (G. Li)

32

Another HMM Example

• Given a sequence of dice rolls, assuming that we know that the casino player may switch between a fair and two loaded die:

FLIP: Pattern Recognition (G. Li)

33

Three basic HMM problems

- Evaluation**
 - Determine the probability that a particular sequence of **visible states** was generated by a given model
 - i.e., Forward/Backward algorithm
- Decoding**
 - Given a sequence of **visible states**, determine the most likely sequence of **hidden states** that led to those observations
 - i.e., using Viterbi algorithm
- Learning**
 - Given a set of visible observations, determine a_{ij} and b_{jk}
 - i.e., using EM algorithm - Baum-Welch algorithm

FLIP: Pattern Recognition (G. Li)

34

HMM Evaluation

- HMM Evaluation
- Recursive Method

FLIP: Pattern Recognition (G. Li)

35

An HMM Example

(Evaluation: Prob. of a sequence of observations)

FLIP: Pattern Recognition (G. Li)

36

An HMM Example

(Evaluation: Prob. of a sequence of observations)

- What is $P(\mathbf{O}) = P(O_1 O_2 O_3)$
= $P(O_1 = X \wedge O_2 = X \wedge O_3 = Z)$?
- Slow, stupid way:

$$P(\mathbf{O}) = \sum_{Q \in \text{Paths of length 3}} P(\mathbf{O} \wedge Q)$$

$$= \sum_{Q \in \text{Paths of length 3}} P(\mathbf{O} | Q) P(Q)$$

How do we compute $P(Q)$ for an arbitrary path Q^2

- How do we compute $P(O_i | Q)$ for an arbitrary path Q^2

Example in the case $Q = S_1 S_3 S_2$:
 $P(X|S_1) P(Z|S_3) P(Y|S_2) = 1/2 * 2/3 * 1/3 = 1/9$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li)

37

An HMM Example

(Evaluation: Prob. of a sequence of observations)

- What is $P(\mathbf{O}) = P(O_1 O_2 O_3)$
= $P(O_1 = X \wedge O_2 = X \wedge O_3 = Z)$?
- Slow, stupid way:

$$P(\mathbf{O}) = \sum_{Q \in \text{Paths of length 3}} P(\mathbf{O} \wedge Q)$$

$$= \sum_{Q \in \text{Paths of length 3}} P(\mathbf{O} | Q) P(Q)$$

How do we compute $P(O_i | Q)$ for an arbitrary path Q^2

Example in the case $Q = S_1 S_3 S_2$:
 $P(X|S_1) P(X|S_3) P(Z|S_2) = 1/2 * 1/2 * 1/2 = 1/8$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li)

38

An HMM Example

(Evaluation: Prob. of a sequence of observations)

- What is $P(\mathbf{O}) = P(O_1 O_2 O_3)$
= $P(O_1 = X \wedge O_2 = X \wedge O_3 = Z)$?

$P(\mathbf{O})$ would need 27 $P(Q)$ computations and 27 $P(O_i | Q)$ computations

A sequence of 20 observations would need $3^{20} = 3.5$ billion computations and 3.5 billion $P(O_i | Q)$ computations

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li)

39

An HMM Example

(Evaluation: Prob. of a sequence of observations)

- Given observations $O_1 O_2 \dots O_T$, define $\alpha_t(i) = P(O_1 O_2 \dots O_t \wedge q_t = S_i | \lambda)$

where $1 \leq t \leq T$

$n = 3, m = 3$
 $\pi_1 = 1/2 \quad \pi_2 = 1/3 \quad \pi_3 = 0$

a ₁₁ = 0	a ₁₂ = 1/3	a ₁₃ = 2/3
a ₂₁ = 1/3	a ₂₂ = 0	a ₂₃ = 2/3
a ₃₁ = 1/3	a ₃₂ = 1/3	a ₃₃ = 1/3

b ₁ (X) = 1/2	b ₁ (Y) = 1/2	b ₁ (Z) = 0
b ₂ (X) = 0	b ₂ (Y) = 1/2	b ₂ (Z) = 1/2
b ₃ (X) = 1/2	b ₃ (Y) = 0	b ₃ (Z) = 1/2

$\alpha_t(i) = P(O_1 \wedge q_1 = S_i)$
 $= P(q_1 = S_i) P(O_1 | q_1 = S_i)$

$\alpha_{t+1}(j) = P(O_1, \dots, O_t \wedge q_t = S_i \wedge O_{t+1} \wedge q_{t+1} = S_j)$
 $= \sum_{i=1}^n P(O_1, \dots, O_t \wedge q_t = S_i \wedge O_{t+1}, \dots, O_n \wedge q_{t+1} = S_j)$
 $= \sum_{i=1}^n P(O_{t+1}, \dots, O_n | O_1, \dots, O_t, q_t = S_i) P(O_1, \dots, O_t \wedge q_t = S_i)$
 $= \sum_{i=1}^n P(q_{t+1} = S_j | O_{t+1}, \dots, O_n) P(O_{t+1}, \dots, O_n | O_1, \dots, O_t, q_t = S_i) \alpha_t(i)$
 $= \sum_i a_{ij} b_j(O_{t+1}, \dots, O_n) \alpha_t(i)$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li)

40

An HMM Example

(Evaluation: Prob. of a sequence of observations)

- Given observations $O_1 O_2 \dots O_T$, define $\alpha_t(i) = P(O_1 O_2 \dots O_t \wedge q_t = S_i | \lambda)$

where $1 \leq t \leq T$

$\alpha_t(i) = P(O_1 \wedge q_1 = S_i)$
 $= P(q_1 = S_i) P(O_1 | q_1 = S_i)$

$\alpha_{t+1}(j) = P(O_1, \dots, O_t \wedge q_t = S_i \wedge O_{t+1}, \dots, O_n \wedge q_{t+1} = S_j)$
 $= \sum_{i=1}^n P(O_1, \dots, O_t \wedge q_t = S_i \wedge O_{t+1}, \dots, O_n \wedge q_{t+1} = S_j)$
 $= \sum_{i=1}^n P(O_{t+1}, \dots, O_n | O_1, \dots, O_t, q_t = S_i) P(O_1, \dots, O_t \wedge q_t = S_i)$
 $= \sum_{i=1}^n P(q_{t+1} = S_j | O_{t+1}, \dots, O_n) P(O_{t+1}, \dots, O_n | O_1, \dots, O_t, q_t = S_i) \alpha_t(i)$
 $= \sum_i a_{ij} b_j(O_{t+1}, \dots, O_n) \alpha_t(i)$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li)

41

An HMM Example

(Evaluation: Prob. of a sequence of observations)

- Given observations $O_1 O_2 O_3 = XXZ$
 $\alpha_t(i) = P(O_1 O_2 \dots O_t \wedge q_t = S_i | \lambda)$

$\alpha_1(i) = b_i(O_1) \pi_i$
 $\alpha_{t+1}(j) = \sum_i a_{ij} b_j(O_{t+1}) \alpha_t(i)$

$n = 3, m = 3$
 $\pi_1 = 1/2 \quad \pi_2 = 1/3 \quad \pi_3 = 0$

a ₁₁ = 0	a ₁₂ = 1/3	a ₁₃ = 2/3
a ₂₁ = 1/3	a ₂₂ = 0	a ₂₃ = 2/3
a ₃₁ = 1/3	a ₃₂ = 1/3	a ₃₃ = 1/3

b ₁ (X) = 1/2	b ₁ (Y) = 1/2	b ₁ (Z) = 0
b ₂ (X) = 0	b ₂ (Y) = 1/2	b ₂ (Z) = 1/2
b ₃ (X) = 1/2	b ₃ (Y) = 0	b ₃ (Z) = 1/2

$\alpha_1(1) = \frac{1}{4}$ $\alpha_1(2) = 0$ $\alpha_1(3) = 0$
 $\alpha_2(1) = 0$ $\alpha_2(2) = 0$ $\alpha_2(3) = \frac{1}{12}$
 $\alpha_3(1) = 0$ $\alpha_3(2) = \frac{1}{72}$ $\alpha_3(3) = \frac{1}{72}$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li)

42

An HMM Example

(Evaluation: Prob. of a sequence of observations)

$n = 3, m = 3$
 $\pi_1 = 1/2 \quad \pi_2 = 1/2 \quad \pi_3 = 0$

$a_{11} = 0$	$a_{12} = 1/3$	$a_{13} = 2/3$
$a_{21} = 1/3$	$a_{22} = 0$	$a_{23} = 2/3$
$a_{31} = 1/3$	$a_{32} = 1/3$	$a_{33} = 1/3$

$b_1(X) = 1/2$	$b_2(Y) = 1/2$	$b_3(Z) = 0$
$b_1(X) = 0$	$b_2(Y) = 1/2$	$b_3(Z) = 1/2$
$b_1(X) = 1/2$	$b_2(Y) = 0$	$b_3(Z) = 1/2$

$\alpha_1(1) = \frac{1}{4}$ $\alpha_1(2) = 0$ $\alpha_1(3) = 0$
 $\alpha_2(1) = 0$ $\alpha_2(2) = 0$ $\alpha_2(3) = \frac{1}{12}$
 $\alpha_3(1) = 0$ $\alpha_3(2) = \frac{1}{72}$ $\alpha_3(3) = \frac{1}{72}$

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

43

An HMM Example

(Evaluation: Prob. of a sequence of observations)

$n = 3, m = 3$
 $\pi_1 = 1/2 \quad \pi_2 = 1/2 \quad \pi_3 = 0$

$a_{11} = 0$	$a_{12} = 1/3$	$a_{13} = 2/3$
$a_{21} = 1/3$	$a_{22} = 0$	$a_{23} = 2/3$
$a_{31} = 1/3$	$a_{32} = 1/3$	$a_{33} = 1/3$

$b_1(X) = 1/2$	$b_2(Y) = 1/2$	$b_3(Z) = 0$
$b_1(X) = 0$	$b_2(Y) = 1/2$	$b_3(Z) = 1/2$
$b_1(X) = 1/2$	$b_2(Y) = 0$	$b_3(Z) = 1/2$

$\alpha_1(i) = P(O_1 O_2 \dots O_n | S_i)$

We can cheaply compute $\alpha_1(i) = P(O_1 O_2 \dots O_n | q_i = S_i)$

How can we cheaply compute $P(O_1 O_2 \dots O_n)$?

$$\sum_{i=1}^N \alpha_i(i)$$

How can we cheaply compute $P(q_i = S_i | O_1 O_2 \dots O_n)$?

$$\frac{\alpha_i(i)}{\sum_{j=1}^N \alpha_j(j)}$$

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

44

HMM Decoding

- HMM Decoding
- Viterbi Algorithm

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

45

Viterbi Algorithm

(Decoding: Most probable path given observations)

$n = 3, m = 3$
 $\pi_1 = 1/2 \quad \pi_2 = 1/2 \quad \pi_3 = 0$

$a_{11} = 0$	$a_{12} = 1/3$	$a_{13} = 2/3$
$a_{21} = 1/3$	$a_{22} = 0$	$a_{23} = 2/3$
$a_{31} = 1/3$	$a_{32} = 1/3$	$a_{33} = 1/3$

$b_1(X) = 1/2$	$b_2(Y) = 1/2$	$b_3(Z) = 0$
$b_1(X) = 0$	$b_2(Y) = 1/2$	$b_3(Z) = 1/2$
$b_1(X) = 1/2$	$b_2(Y) = 0$	$b_3(Z) = 1/2$

$\arg\max_Q P(Q | O_1 O_2 \dots O_r)$

This is what the observer has to work with...

$q_0 =$	S_1	$O_0 =$	X
$q_1 =$	S_3	$O_1 =$	X
$q_2 =$	S_3	$O_2 =$	Z

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

46

Viterbi Algorithm

(Decoding: Most probable path given observations)

$n = 3, m = 3$
 $\pi_1 = 1/2 \quad \pi_2 = 1/2 \quad \pi_3 = 0$

$a_{11} = 0$	$a_{12} = 1/3$	$a_{13} = 2/3$
$a_{21} = 1/3$	$a_{22} = 0$	$a_{23} = 2/3$
$a_{31} = 1/3$	$a_{32} = 1/3$	$a_{33} = 1/3$

$b_1(X) = 1/2$	$b_2(Y) = 1/2$	$b_3(Z) = 0$
$b_1(X) = 0$	$b_2(Y) = 1/2$	$b_3(Z) = 1/2$
$b_1(X) = 1/2$	$b_2(Y) = 0$	$b_3(Z) = 1/2$

$\arg\max_Q P(Q | O_1 O_2 \dots O_r)$

Slow, stupid answer:

$$\arg\max_Q P(Q | O_1 O_2 \dots O_r)$$

$$= \arg\max_Q \frac{P(O_1 O_2 \dots O_r | Q) P(Q)}{P(O_1 O_2 \dots O_r)}$$

$$= \arg\max_Q P(O_1 O_2 \dots O_r | Q) P(Q)$$

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

47

Viterbi Algorithm

(Decoding: Most probable path given observations)

- We're going to compute the following variables:

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P(q_1 q_2 \dots q_{t-1} \wedge q_t = S_i \wedge O_1 \dots O_t)$$

= Probability of Length **t-1** path with the maximum chance of:
...OCCURRING
and
...ENDING UP IN STATE S_i
and
...PRODUCING OUTPUT $O_1 \dots O_t$

- DEFINE: $mpp_t(i) =$ that path,
 $\delta_t(i) = P(mpp_t(i))$

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

48

Viterbi Algorithm

(Decoding: Most probable path given observations)

$$\delta_t(i) = \max_{q_1 q_2 \dots q_{t-1}} P(q_1 q_2 \dots q_{t-1} \wedge q_t = S_i \wedge O_1 O_2 \dots O_t)$$

$$mpp_t(i) = \arg \max_{q_1 q_2 \dots q_{t-1}} P(q_1 q_2 \dots q_{t-1} \wedge q_t = S_i \wedge O_1 O_2 \dots O_t)$$

$$\delta_t(i) = \text{one choice } P(q_t = S_i \wedge O_t) = P(q_t = S_i) P(O_t | q_t = S_i) = \pi_i b_i(O_t)$$

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

49

Viterbi Algorithm

(Decoding: Most probable path given observations)

- Suppose we have $\delta_t(i)$ and $mpp_t(i)$ for all i .
- Task:
 - How to get $\delta_{t+1}(j)$ and $mpp_{t+1}(j)$?

What is the prob of that path?
 $\delta_t(j) \times P(S_j \rightarrow S_{t+1} | O_{t+1}) = \delta_t(j) a_{jj} b_j(O_{t+1})$

SO The most probable path to S_t has S^{*} as its penultimate state with
 $i^* = \operatorname{argmax}_i \delta_t(i) a_{ii} b_i(O_{t+1})$

The most prob path with last two states S_1, S_2 is
 • the most prob path to S_t , followed by transition $S_t \rightarrow S_{t+1}$
 Prob= $\delta_t(N)$

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

50

Viterbi Algorithm

(Decoding: Most probable path given observations)

- Suppose we have $\delta_t(i)$ and $mpp_t(i)$ for all i .
- Task:
 - How to get $\delta_{t+1}(j)$ and $mpp_{t+1}(j)$?

$\delta_{t+1}(j) = \delta_t(i^*) a_{ij} b_j(O_{t+1})$
 $mpp_{t+1}(j) = mpp_t(i^*) S_{i^*}$

$i^* = \operatorname{argmax}_i \delta_t(i) a_{ii} b_i(O_{t+1})$

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

51

Andrew Viterbi

"The algorithm was originally created for improving communication from space by being able to operate with a weak signal but today it has a multitude of applications,"
 Viterbi

- Life Fellow Andrew Viterbi created an algorithm used in every cell phone and cofound Qualcomm.
- For the algorithm, which carries his name, he was awarded the Benjamin Franklin Medal in electrical engineering by the Franklin Institute in Philadelphia, one of the United States' oldest centres of science education and development.
- Another successful venture for the company was the creation of code-division multiple access (CDMA), which was introduced commercially in 1995 in cellphones and is still big today.

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

52

HMM Learning

- Learning HMM
- BAUM-WELCH algorithm

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

53

Learning an HMM

- We have been doing things like
 $P(O_1 O_2 \dots O_T | \lambda)$
 λ is the HMM parameters.
- But how do you know this λ ?
 - Hand crafted? OR
 - But usually, especially in Speech or Genetics, it is better to infer it from large amounts of data
 - $O_1 O_2 \dots O_T$ with a big "T".

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

54

Learning an HMM

- Task: we have some observations $O_1 O_2 \dots O_T$ and we want to estimate λ from them.
- Two approaches
 - MLE
 - B.E.

$\lambda = \operatorname{argmax} P(O_1 \dots O_T | \lambda)$

- Work out $P(\lambda | O_1 \dots O_T)$, and then
 - take $E[\lambda]$ or
 - max $P(\lambda | O_1 \dots O_T)$

 **TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 55

MLE HMM from Observations

- Expectation**

$$\gamma_i(i) = P(q_i = S_i | O_1 O_2 \dots O_T, \lambda)$$

$$\varepsilon_i(i, j) = P(q_i = S_i \wedge q_{i+1} = S_j | O_1 O_2 \dots O_T, \lambda)$$

$$\sum_{t=1}^{T-1} \gamma_t(i) = \text{expected number of transitions out of state } i \text{ during path}$$

$$\sum_{t=1}^{T-1} \varepsilon_t(i, j) = \text{expected number of transitions out of } i \text{ and into } j \text{ during path}$$

 **TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 56

MLE HMM from Observations

- Maximization**

Notice that $\frac{\sum_{j=1}^{T-1} \varepsilon_t(i, j)}{\sum_{i=1}^T \gamma_t(i)} = \frac{\left(\begin{array}{c} \text{expected frequency} \\ i \rightarrow j \end{array} \right)}{\left(\begin{array}{c} \text{expected frequency} \\ i \end{array} \right)}$

= Estimate of Prob(Next state $S_j | \text{This state } S_i$)

We can re-estimate

$$a_{ij} \leftarrow \frac{\sum \varepsilon_t(i, j)}{\sum \gamma_t(i)}$$

We can also re-estimate

$$b_j(O_k) \leftarrow \frac{E[\# \text{times it emits symbol } k \text{ while at state } \omega_j]}{E[\# \text{times it emits any other symbol while at state } \omega_j]}$$

 **TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 57

EM for HMMs (BAUM-WELCH algorithm)

Expectation <ul style="list-style-type: none"> If we knew λ we could estimate expectations such as <ul style="list-style-type: none"> Expected number of times in state i Expected number of transitions $i \rightarrow j$ 	Maximization <ul style="list-style-type: none"> If we knew quantities such as <ul style="list-style-type: none"> Expected number of times in state i Expected number of transitions $i \rightarrow j$ We can MLE: $\lambda = \langle \{a_{ij}\}, \{b_j(j)\}, \pi_i \rangle$
---	---

 **TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 58

EM for HMMs (BAUM-WELCH algorithm)

- BAUM-WELCH algorithm**

- Get your observations $O_1 \dots O_T$
- Guess your first λ estimate $\lambda(0)$, $k=0$
- $k = k+1$
- Given $O_1 \dots O_T$, $\lambda(k)$ compute
 $\gamma_t(i), \varepsilon_t(i, j) \quad \forall 1 \leq t \leq T, \quad \forall 1 \leq i \leq N, \quad \forall 1 \leq j \leq N$
- Compute expected freq. of state i , and expected freq. $i \rightarrow j$
- Compute new estimates of $\lambda(k+1)$: $a_{ij}, b_j(k), \pi_i$ accordingly.
- Goto 3, unless converged.

 **TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 59

EM for HMMs (BAUM-WELCH algorithm)

- Summary of BW-Algorithm
 - There are lots of local minima, which are usually adequate models of the data.
 - EM does not estimate the number of states.
 - Trade-off between too few states (inadequately modeling the structure in the data) and too many (fitting the noise).
 - Thus #states is a regularization parameter.
 - Often, HMMs are forced to have some links with zero probability.
 - This is done by setting $a_{ij}=0$ in initial estimate $\lambda(0)$

 **TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 60

Seminar S05

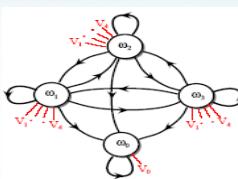


- Topics**
 - Gambler, or Speech Recognition by HMM
 - Viterbi Algorithm Exercise (see next time)
 - Viterbi Algorithm Application
 - Analyze HMM **Forward** or **Backward** estimation **Recursively**. Illustrate them by completing the example in textbook
- Requirements**
 - Prepare a **15 minutes** talk on your chosen topic
 - Make **ppt** to assist your talk
 - Prepare **at least 3 questions** to ask the audience after your talk
 - Get ready to **take questions** from the audience

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 61

Exercise: Find the MPP

Observations: $V^4 = \{v_1, v_3, v_2, v_0\}$
 Initial status: w_1



$$\pi_0 = a_{10} = 0.2, \quad \pi_1 = a_{11} = 0.3,$$

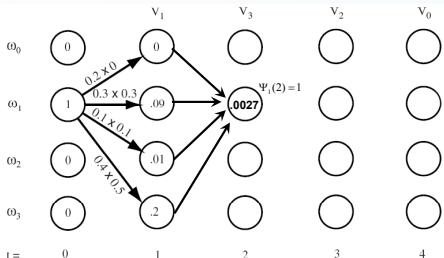
$$\pi_2 = a_{12} = 0.1, \quad \pi_3 = a_{13} = 0.4$$

$$a_{ij} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0.2 & 0.3 & 0.1 & 0.4 \\ 0.2 & 0.5 & 0.2 & 0.1 \\ 0.8 & 0.1 & 0.0 & 0.1 \end{pmatrix}, \quad b_{jk} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0.3 & 0.4 & 0.1 & 0.2 \\ 0 & 0.1 & 0.1 & 0.7 & 0.1 \\ 0 & 0.5 & 0.2 & 0.1 & 0.2 \end{pmatrix}$$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 62

Exercise: Find the MPP

- Complete this Trellis



t = 0 1 2 3 4

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 63

Questions?



TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 64