

HIDDEN MARKOV MODEL

- NIPUN BATRA

SOME APPLICATIONS

① ACTIVITY MONITORING

SOME APPLICATIONS

① ACTIVITY

MONITORING

Accelerometer



Activity

WALKING

RUNNING

SITTING

SOME APPLICATIONS

① ACTIVITY MONITORING

Accelerometer



Activity

WALKING

RUNNING

SITTING

② SPEECH RECOGNITION

RAW SPEECH



SOME APPLICATIONS

① ACTIVITY MONITORING

Accelerometer



Activity

WALKING

RUNNING

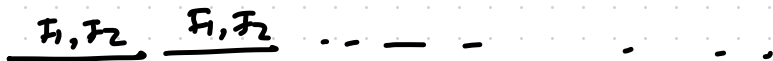
SITTING

② SPEECH RECOGNITION

RAW SPEECH



FEATURES



WORDS

I LIKE MACHINE LEARNING...

SOME APPLICATIONS

- ① ACTIVITY MONITORING
- ② SPEECH RECOGNITION
- ③ PART OF SPEECH TAGGING
- ④ GENE FINDING
- ⋮

SEQUENTIAL DATA

- 1) Amount of rainfall daily
- 2) Price of stock
- 3) Words in a sentence
- 4) Energy consumption of a home

WEATHER EXAMPLE

Let $x_t \in \begin{cases} R: \text{Rainy} \\ S: \text{Sunny} \\ C: \text{cloudy} \end{cases}$

be observed outlook at ' t^{th} ' day

Sample observations: R R R C C R R S S S C S S

Modelling Sequential Data as i.i.d

* Easiest way to treat sequential data \rightarrow ignore sequential aspect and treat observations as i.i.d.

DOWNSIDE?

Modelling Sequential Data as i.i.d

- * Easiest way to treat sequential data \rightarrow ignore sequential aspect and treat observations as i.i.d.

DOWNSIDE?

- * Fails to exploit sequential patterns
 - \rightarrow Nearby observations "correlated"

Modelling Sequential Data as i.i.d

- * Easiest way to treat sequential data \rightarrow ignore sequential aspect and treat observations as i.i.d.

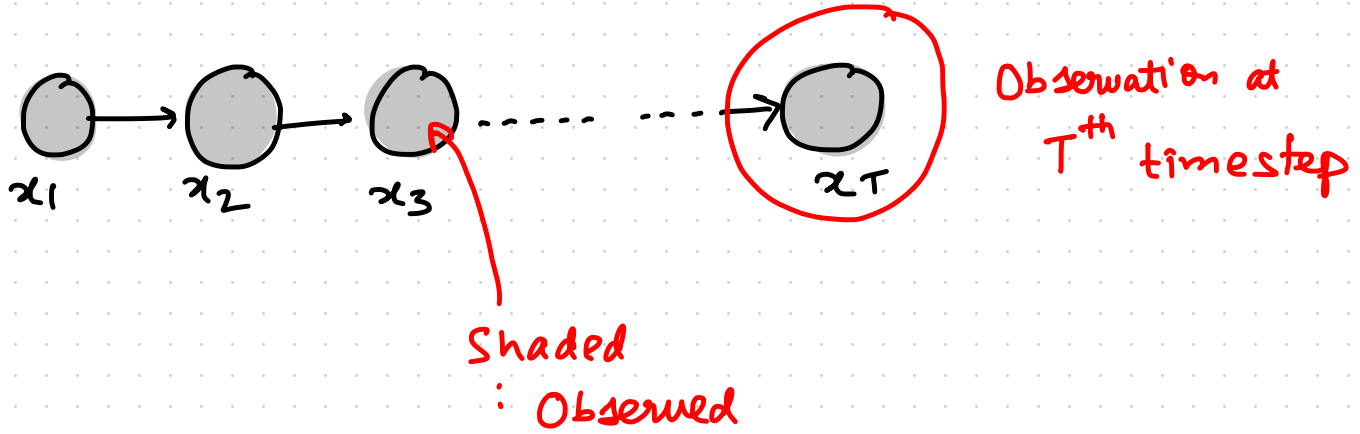
DOWNSIDE?

- * Fails to exploit sequential patterns
 - \rightarrow Nearby observations "correlated"
- * Q: Predict if it will rain today or not?
 - \therefore Rain today $\xrightarrow{\text{w/ high probability}}$ Rain tomorrow

MARKOV MODEL



MARKOV MODEL



MARKOV MODEL



FIRST ORDER MARKOV CHAIN

Future prediction independent of past given present

MARKOV MODEL



FIRST ORDER MARKOV CHAIN

Future prediction independent of past given present

$$\Rightarrow P(x_{t+1} | x_1, x_2, \dots, x_t) = P(x_{t+1} | x_t)$$

MARKOV MODEL



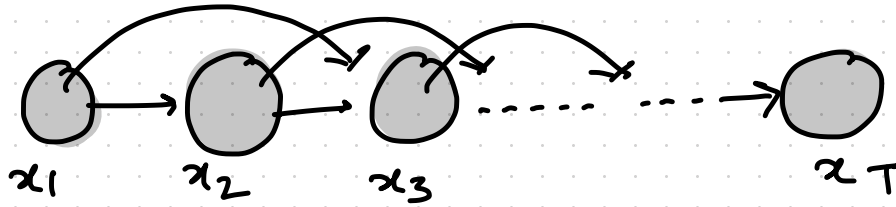
FIRST ORDER MARKOV CHAIN

Future prediction independent of past given present

$$\Rightarrow P(x_{t+1} | x_1, x_2, \dots, x_t) = P(x_{t+1} | x_t)$$

$$\text{JOINT PROB} = P(x_1, \dots, x_T) = P(x_1) \cdot P(x_2 | x_1) \cdot P(x_3 | x_2) \dots P(x_T | x_{T-1})$$

MARKOV MODEL

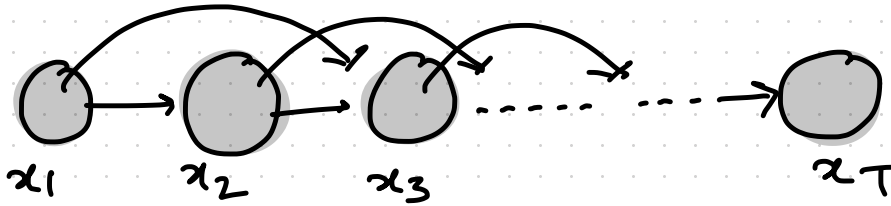


SECOND ORDER MARKOV CHAIN

$$P(x_{t+1} | x_1, x_2, \dots, x_t) = P(x_{t+1} | x_t, x_{t-1})$$

$$\text{JOINT PROB} = P(x_1, \dots, x_T) =$$

MARKOV MODEL



SECOND ORDER MARKOV CHAIN

$$p(x_{t+1} | x_1, x_2, \dots, x_t) = p(x_{t+1} | x_t, x_{t-1})$$

$$\text{JOINT PROB} = p(x_1, \dots, x_T) = p(x_1) \cdot p(x_2 | x_1) \cdot \prod_{n=3}^T p(x_n | x_{n-1}, x_{n-2})$$

PARAMETERS OF FIRST ORDER MARKOV MODEL



FIRST ORDER MARKOV CHAIN

Future prediction independent of past given present

$$\Rightarrow P(x_{t+1} | x_1, x_2, \dots, x_t) = P(x_{t+1} | x_t)$$

$$\text{JOINT PROB} = P(x_1, \dots, x_T) = P(x_1) \cdot P(x_2 | x_1) \cdot P(x_3 | x_2) \dots P(x_T | x_{T-1})$$

PARAMETERS OF FIRST ORDER MARKOV MODEL



① MAIN ASSUMPTION : STATIONARITY

Data evolves over time, but distribution from which data is generated is fixed.

PARAMETERS OF FIRST ORDER MARKOV MODEL



① MAIN ASSUMPTION : STATIONARITY

Data evolves over time, but distribution from which data is generated is fixed.

② we would like to have parameter sharing
(parameters independent of time)

PARAMETERS OF FIRST ORDER MARKOV MODEL



we would like to have parameter sharing

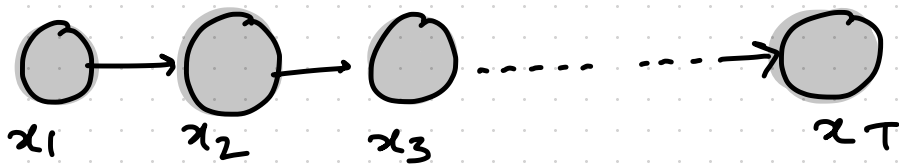
$p(x_t | x_{t-1})$ is same $\forall t$

PARAMETERS OF FIRST ORDER MARKOV MODEL



$$p(x_1, x_2, \dots, x_T) = p(x_1) \prod_{t=2}^T p(x_t | x_{t-1})$$

PARAMETERS OF FIRST ORDER MARKOV MODEL



$$p(x_1, x_2, \dots, x_T) = p(x_1) \prod_{t=2}^T p(x_t | x_{t-1})$$

$$\text{Parameters} = \Theta = \{ \pi, A \}$$

↑
Prior probability

↖ Transition matrix

PARAMETERS OF FIRST ORDER MARKOV MODEL

TRANSITION MATRIX (A)

ASSUME x_t can take 1 of K states

$$A_{jk} \equiv P(x_t = k \mid x_{t-1} = j)$$

$$'j' \text{ and } 'k' \in \{1, \dots, K\}$$

PARAMETERS OF FIRST ORDER MARKOV MODEL

ASSUME x_t can take 1 of K states

$$A_{jk} \equiv P(x_t = k \mid x_{t-1} = j)$$

'j' and 'k' $\in \{1, \dots, K\}$

$j \rightarrow k$

PARAMETERS OF FIRST ORDER MARKOV MODEL

Representing $A_{jk} \equiv P(x_t = k \mid x_{t-1} = j)$

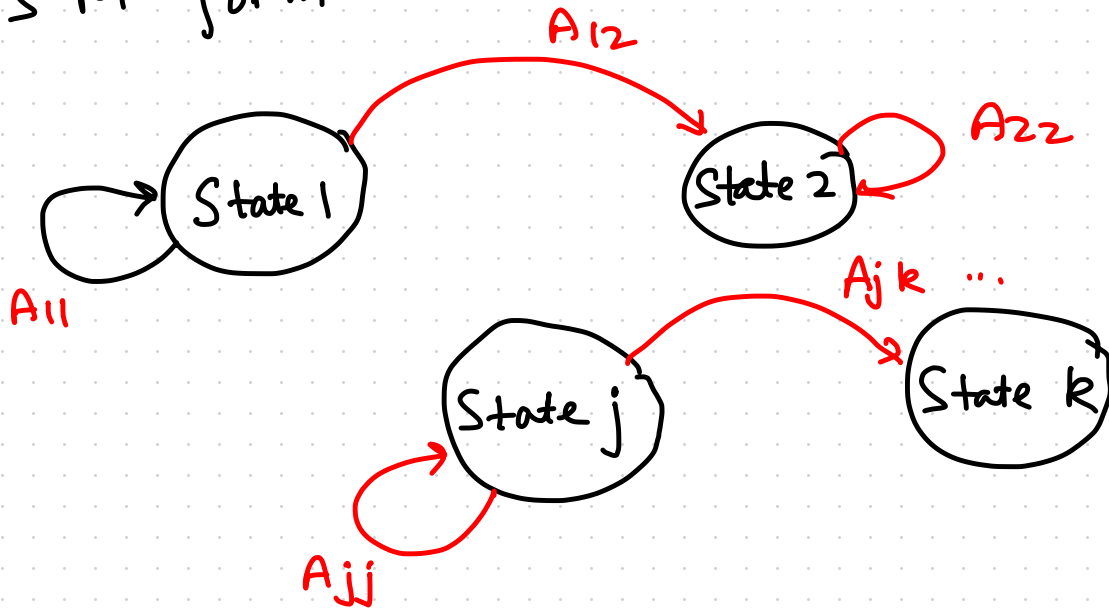
(i) Tabular form

x_{t-1}	x_t	$P(x_t \mid x_{t-1})$
1	1	A_{11}
2	1	A_{21}
\vdots		\cdot
j	k	A_{jk}
k	k	A_{kk}

PARAMETERS OF FIRST ORDER MARKOV MODEL

Representing $A_{jk} \equiv P(x_t = k \mid x_{t-1} = j)$

② FSM form



PARAMETERS OF FIRST ORDER MARKOV MODEL

Representing $A_{jk} \equiv P(x_t = k \mid x_{t-1} = j)$

(3) Adjacency matrix

$x_{t-1} \backslash x_t$	1	2	...	k
1				
2				
...				
j				A_{jk}

PARAMETERS OF FIRST ORDER MARKOV MODEL

Prior Probability (π)

$$\pi_k = P(z_1 = k)$$

Probability of 1st observation

0) Given A & π as:

$$A = \begin{array}{c} R \\ C \\ S \end{array} \begin{array}{ccc} R & C & S \\ \left[\begin{array}{ccc} .4 & .3 & .3 \\ .2 & .6 & .2 \\ .1 & .1 & .8 \end{array} \right] \end{array}$$

$$\pi = \begin{array}{ccc} [.2 & .4 & .4] \\ R & C & S \end{array}$$

WHAT IS PROBABILITY OF SEQUENCE

$$\alpha = \{RRSSC\}$$

0) Given A & π as:

$$A = \begin{matrix} & \begin{matrix} R & C & S \end{matrix} \\ \begin{matrix} R \\ C \\ S \end{matrix} & \begin{bmatrix} .4 & .3 & .3 \\ .2 & .6 & .2 \\ .1 & .1 & .8 \end{bmatrix} \end{matrix}$$

$$\pi = \begin{matrix} & \begin{matrix} R & C & S \end{matrix} \\ \begin{matrix} R \\ C \\ S \end{matrix} & \begin{bmatrix} .2 & .4 & .4 \\ . & . & . \\ . & . & . \end{bmatrix} \end{matrix}$$

WHAT IS PROBABILITY OF SEQUENCE

$$x = \{RRSSC\}$$

$$P(x | \theta = \{A, \pi\}) = P(x_1=R) \cdot P(x_2=R | x_1=R) \cdot P(x_3=S | x_2=R) \\ \cdot P(x_4=S | x_3=S) \cdot P(x_5=C | x_4=S)$$

$$= \pi_R \cdot A_{RR} \cdot A_{RS} \cdot A_{SS} \cdot A_{SC}$$

0) Given A & π as:

$$A = \begin{array}{c} R \\ C \\ S \end{array} \begin{array}{ccc} R & C & S \\ \left[\begin{array}{ccc} .4 & .3 & .3 \\ .2 & .6 & .2 \\ .1 & .1 & .8 \end{array} \right] \end{array}$$

$$\pi = \begin{array}{ccc} [.2 & .4 & .4] \\ R & C & S \end{array}$$

$p(x_1 = R) = 0.2$. What is $p(x_2 = R)$?

0) Given A & π as:

$$A = \begin{array}{c} R \\ C \\ S \end{array} \begin{array}{c} R \\ C \\ S \end{array} \begin{bmatrix} .4 & .3 & .3 \\ .2 & .6 & .2 \\ .1 & .1 & .8 \end{bmatrix}$$

$$\pi = \begin{array}{c} R \\ C \\ S \end{array} \begin{bmatrix} .2 & .4 & .4 \end{bmatrix}$$

$P(x_1 = R) = 0.2$. What is $P(x_2 = R)$?

$$P(x_2 = R) = P(x_1 = R) \cdot P(x_2 = R | x_1 = R) + P(x_1 = S) \cdot P(x_2 = R | x_1 = S) + P(x_1 = C) \cdot P(x_2 = R | x_1 = C)$$

0) Given A & π as:

$$A = \begin{array}{c} R \\ C \\ S \end{array} \begin{array}{c} R \\ C \\ S \end{array} \begin{bmatrix} .4 & .3 & .3 \\ .2 & .6 & .2 \\ .1 & .1 & .8 \end{bmatrix}$$

$$\pi = \begin{array}{c} R \\ C \\ S \end{array} \begin{bmatrix} .2 & .4 & .4 \end{bmatrix}$$

$P(x_1 = R) = 0.2$. What is $P(x_2 = R)$?

$$\begin{aligned} P(x_2 = R) &= P(x_1 = R) \cdot P(x_2 = R | x_1 = R) + \\ &P(x_1 = S) \cdot P(x_2 = R | x_1 = S) + \\ &P(x_1 = C) \cdot P(x_2 = R | x_1 = C) \end{aligned}$$

$$= (.2) * (.4) + (.4) * (.1) + (.4) * (.2) = 0.2 = P(x_1 = R)$$

Hidden Markov Model

↑ what is the hidden component

Hidden Markov Model

↖ what is the hidden component

Example: Coin Toss Model

Assume 2 coins

- Biased ($P(H) = 0.7$)
- Unbiased ($P(H) = 0.5$)
Fair

Hidden Markov Model

↑ what is the hidden component

Example: Coin Toss Model

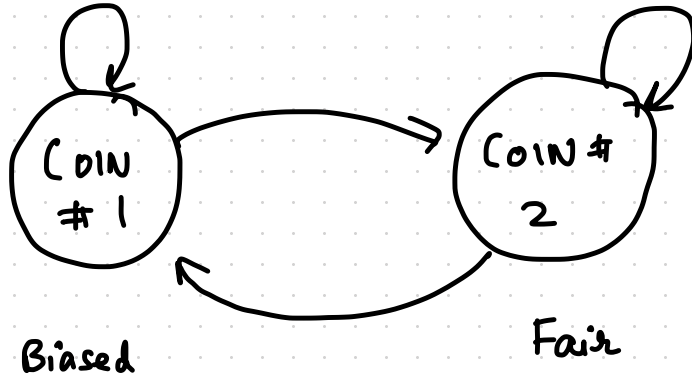
Assume 2 coins

- Biased ($P(H) = 0.7$)
- Unbiased ($P(H) = 0.5$)
Fair

We see only sequence of observations $x = \{H, T, \dots\}$

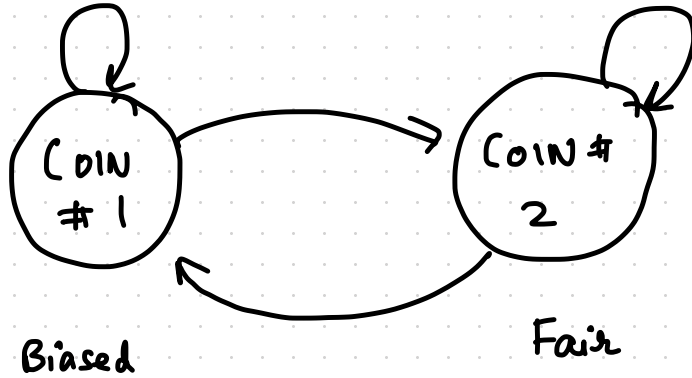
Hidden: which coin was tossed!

HMM (example)



Markov Model

HMM (example)



$$P(H) = 0.7$$
$$P(T) = 0.3$$

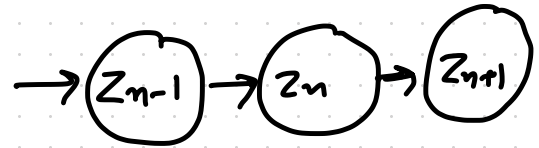
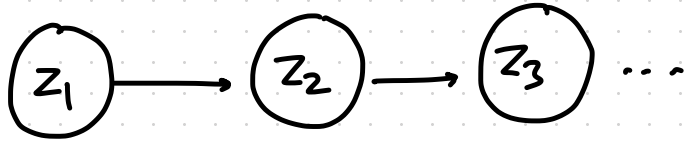
$$P(H) = 0.5$$
$$P(T) = 0.5$$

Markov Model

Emission

HMM (example)

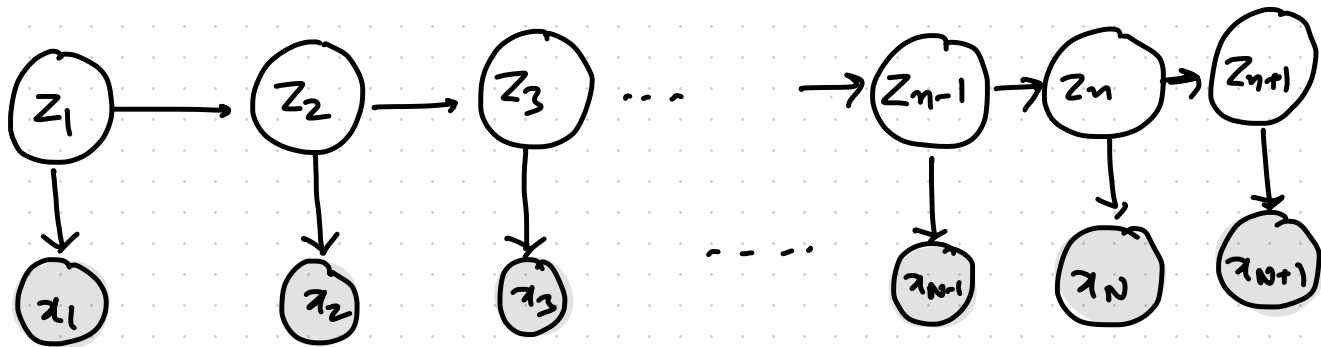
$$Z_i = \left\{ \begin{array}{l} \text{COIN 1,} \\ \text{COIN 2} \end{array} \right\}$$



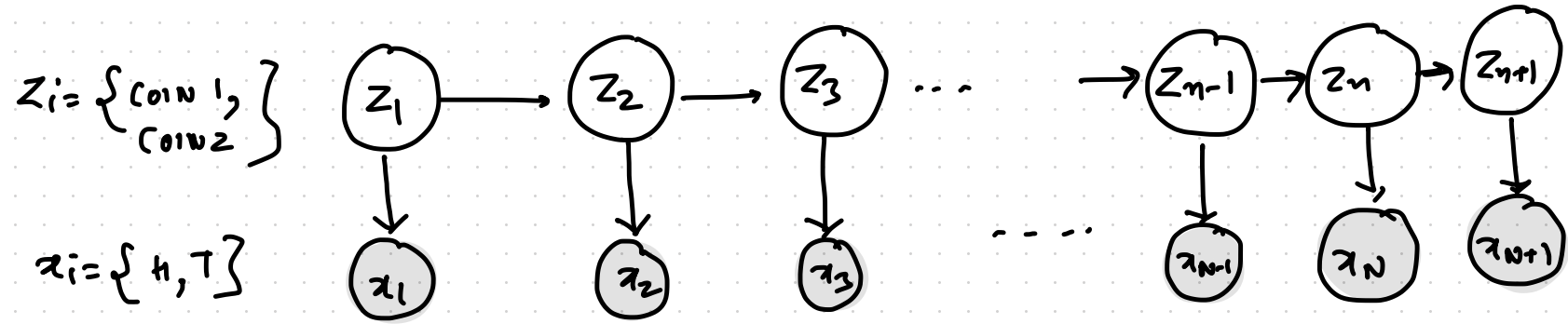
HMM (example)

$$z_i = \left\{ \begin{array}{l} \text{COIN 1,} \\ \text{COIN 2} \end{array} \right\}$$

$$x_i = \{H, T\}$$



HMM (example)



x_i : Observe d_i

z_i : Latent state

HMM Parameters

① Transition matrix A :

$$A_{jk} = P(z_n = k | z_{n-1} = j)$$

HMM Parameters

① Transition matrix A :

$$A_{jk} = P(z_n = k | z_{n-1} = j)$$

② Prior prob. π :

$$\pi_k = P(z_1 = k)$$

HMM Parameters

① Transition matrix A :

$$A_{jk} = p(z_n = k | z_{n-1} = j)$$

② Prior prob. π :

$$\pi_k = p(z_1 = k)$$

③ Emission Probability: ϕ define $p(x_n | z_n, \phi)$

HMM Parameters

③ Emission Probability: ϕ define $p(x_n | z_n, \phi)$

Case I: ϕ is discrete.

e.g. $z_i = \{\text{Fair, Biased}\}$

ϕ defined in terms of conditional probability

Fair

$$P(H) = 0.5$$

$$P(T) = 0.5$$

Biased

$$P(H) = 0.7$$

$$P(T) = 0.3$$

HMM Parameters

③ Emission Probability: ϕ define $p(x_n | z_n, \phi)$

Case II ϕ is continuous

State: Appliance is ON) OFF

ϕ : ON : Power $\sim N(100, 10)$

OFF : Power $\sim N(0, 5)$

SOME APPLICATIONS

① ACTIVITY MONITORING

Accelerometer



Activity

WALKING
 z_1

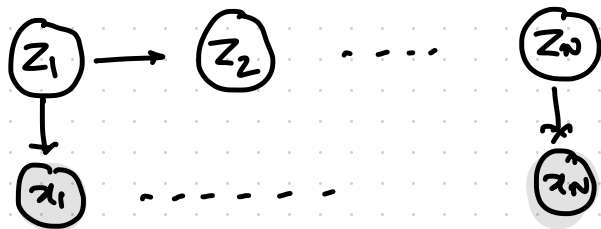
RUNNING
 z_2

SITTING
 z_3

$z_i \in \{ \text{WALKING, RUNNING, SITTING} \}$

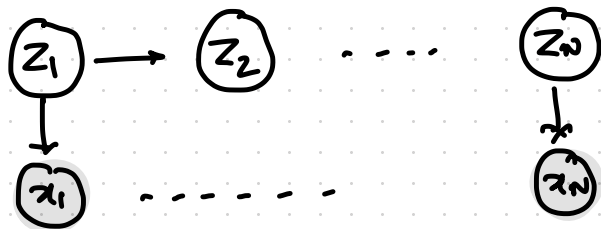
$x_i \in (\text{SOME TRANSFORMATION OF RAW DATA})$

PI: HMM sampling



Given π, A, ϕ generate $x = \{x_1 \dots x_N\}$ & $z = \{z_1 \dots z_N\}$

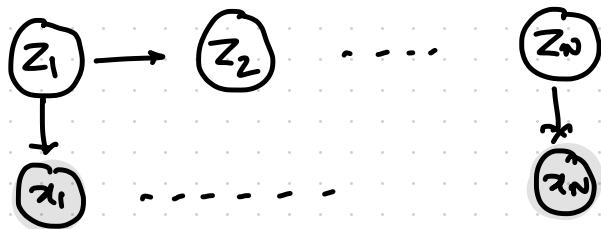
PI: HMM sampling



Given π, A, ϕ generate $x = \{x_1 \dots x_N\}$ & $z = \{z_1 \dots z_N\}$

① CHOOSE z_1 as per π

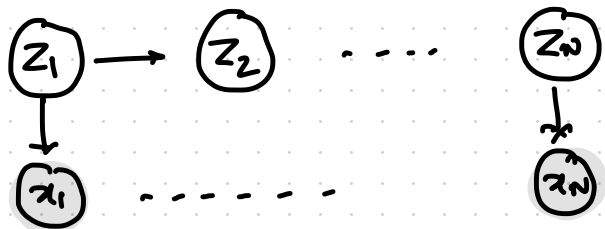
PI: HMM sampling



Given π, A, ϕ generate $x = \{x_1, \dots, x_N\}$ & $z = \{z_1, \dots, z_N\}$

- ① CHOOSE z_1 as per π
- ② Sample x_1 using ϕ and z_1

PI: HMM sampling



Given π, A, ϕ generate $x = \{x_1, \dots, x_n\}$ & $z = \{z_1, \dots, z_n\}$

- ① CHOOSE z_1 as per π
- ② Sample x_1 using ϕ and z_1
- ③ For $n = 2 : N$
 - ③.1 Sample z_n from z_{n-1} using A and z_{n-1}
 - ③.2 Sample x_n from z_n using ϕ and z_n

PI: HMM sampling

Note books: Discrete HMM Simulatⁿ (Unfair Casino)

Continuous HMM Simulatⁿ (Power Appliance)

P II HMM Evidence likelihood

$$X = \{x_1 \dots x_T\}$$

Given $X, \underbrace{\pi, A, \phi}_{\theta}$ what is $L(X|\theta)$?

PII: HMM Evidence likelihood

$$X = \{x_1 \dots x_T\}$$

Given $X, \underbrace{\pi, A, \phi}_{\theta}$ what is $L(X|\theta)$?

$$\text{Likelihood} = p(X|\theta) = \sum_z p(X, z|\theta)$$

(Marginalisat²)

FOR $X = \{H, T\}$

what is $L(X|\theta)$ where $\theta = 2$

$$\pi = \begin{bmatrix} 0.6 & 0.4 \\ \text{Bias} & \text{Fair} \end{bmatrix}$$

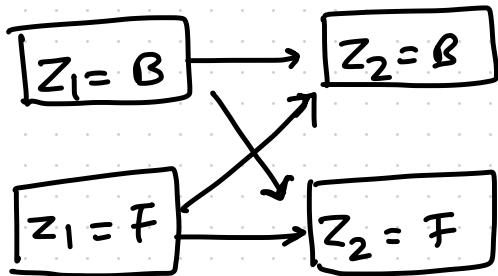
$$A = \begin{array}{cc} & \begin{matrix} B & F \end{matrix} \\ \begin{matrix} \text{Bias} \\ \text{Fair} \end{matrix} & \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix} \end{array}$$

$$\phi = \begin{array}{l} \text{Bias} \\ p(H) = 0.7 \\ p(T) = 0.3 \end{array}$$

$$\begin{array}{l} \text{Fair} \\ p(H) = 0.5 \\ p(T) = 0.5 \end{array}$$

For $X = \{H, H\}$

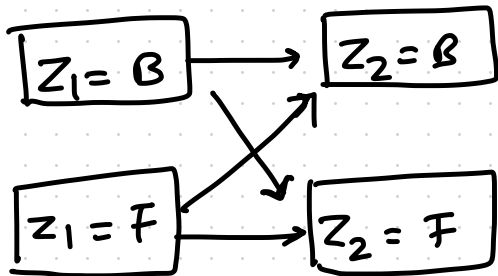
what is $L(X|\theta)$ where $\theta = 2$



Trellis Diagram
for 2 states
and 2 time steps

For $X = \{H, H\}$

what is $L(X|\theta)$ where $\theta = 2$

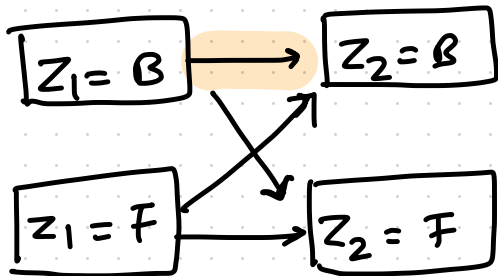


Trellis Diagram
for 2 states
and 2 time steps

$$L(X|\theta) = p(z_1 = B) \cdot p(x_1 = H | z = B) \cdot p(z_2 = B | z_1 = B) \cdot p(x_2 = H | z = B)$$

For $X = \{H, H\}$

what is $L(X|\theta)$ where $\theta = 2$

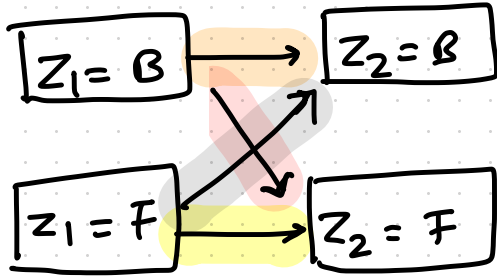


Trellis Diagram
for 2 states
and 2 time steps

$$L(X|\theta) = p(z_1 = B) \cdot p(x_1 = H | z = B) \cdot p(z_2 = B | z_1 = B) \cdot p(x_2 = H | z = B)$$

For $X = \{H, F\}$

what is $L(X|\theta)$ where $\theta = 2$



Trellis Diagram
for 2 states
and 2 time steps

$$\begin{aligned} L(X|\theta) = & p(z_1 = B) \cdot p(x_1 = H | z = B) \cdot p(z_2 = B | z_1 = B) \cdot p(x_2 = H | z = B) \\ & + p(z_1 = F) \cdot p(x_1 = H | z = F) \cdot p(z_2 = B | z_1 = F) \cdot p(x_2 = H | z = B) \\ & + p(z_1 = B) \cdot p(x_1 = H | z = B) \cdot p(z_2 = F | z_1 = B) \cdot p(x_2 = H | z = F) \\ & + p(z_1 = F) \cdot p(x_1 = H | z = F) \cdot p(z_2 = F | z_1 = F) \cdot p(x_2 = H | z = F) \end{aligned}$$

For $X = \{H, H\}$ what is $L(X|\theta)$ where $\theta = 2$

$$L(X|\theta) = p(z_1 = B) \cdot p(x_1 = H | z = B) \cdot p(z_2 = B | z_1 = B) \cdot p(x_2 = H | z = B)$$

$$+ p(z_1 = F) \cdot p(x_1 = H | z = F) \cdot p(z_2 = B | z_1 = F) \cdot p(x_2 = H | z = B)$$

$$+ p(z_1 = B) \cdot p(x_1 = H | z = B) \cdot p(z_2 = F | z_1 = B) \cdot p(x_2 = H | z = F)$$

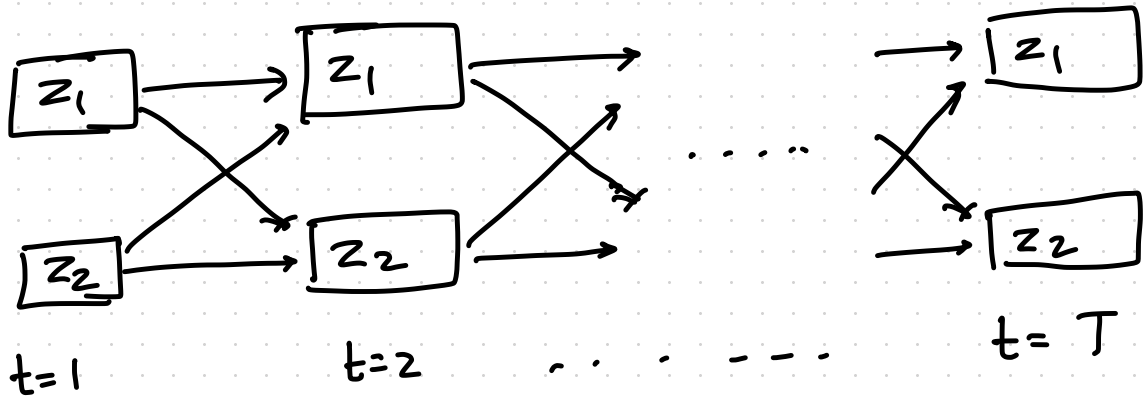
$$+ p(z_1 = F) \cdot p(x_1 = H | z = F) \cdot p(z_2 = F | z_1 = F) \cdot p(x_2 = H | z = F)$$

$$\begin{aligned} &= (0.6)(0.7)(0.9)(0.7) \\ &+ (0.4)(0.5)(0.1)(0.7) \\ &+ (0.6)(0.7)(0.1)(0.5) \\ &+ (0.4)(0.5)(0.9)(0.5) \end{aligned}$$

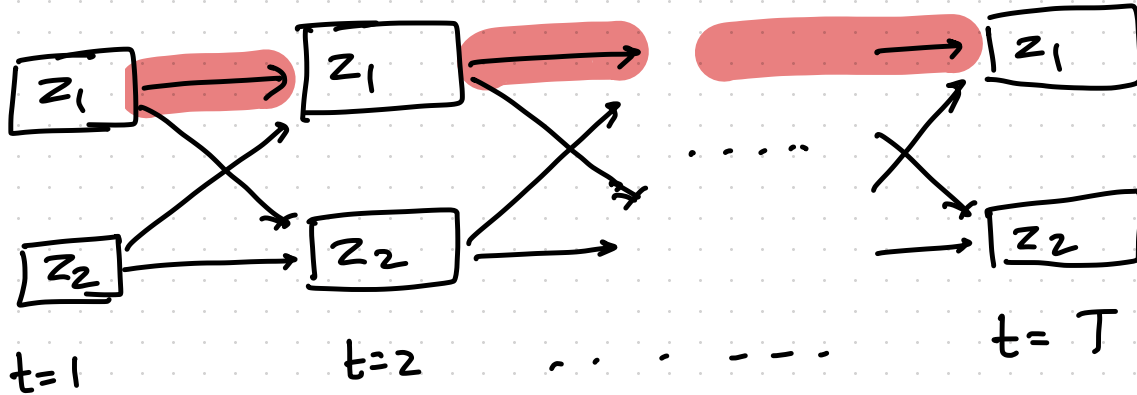
$$\begin{aligned} &= 0.2646 = 0.3896 \\ &+ 0.0140 \\ &+ 0.0210 \\ &+ 0.0900 \end{aligned}$$

For $X = \underbrace{\{H_1, H_2, \dots\}}_{T \text{ timesteps}}$ what is $L(X|\theta)$ where $\theta = \nu$

For $X = \{H_1, H_2, \dots\}$ what is $L(X|\theta)$ where $\theta = \rho$
T timesteps.

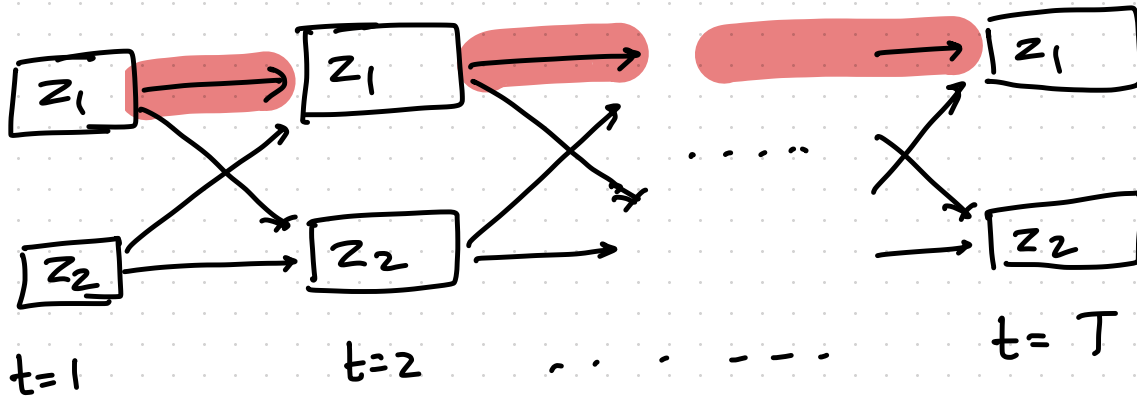


For $X = \{H_1, H_2, \dots\}$ what is $L(X|\theta)$ where $\theta = 2$
T timesteps.



SNO	1	2	...	T
1	z_1	z_1		z_1
...				
K^T

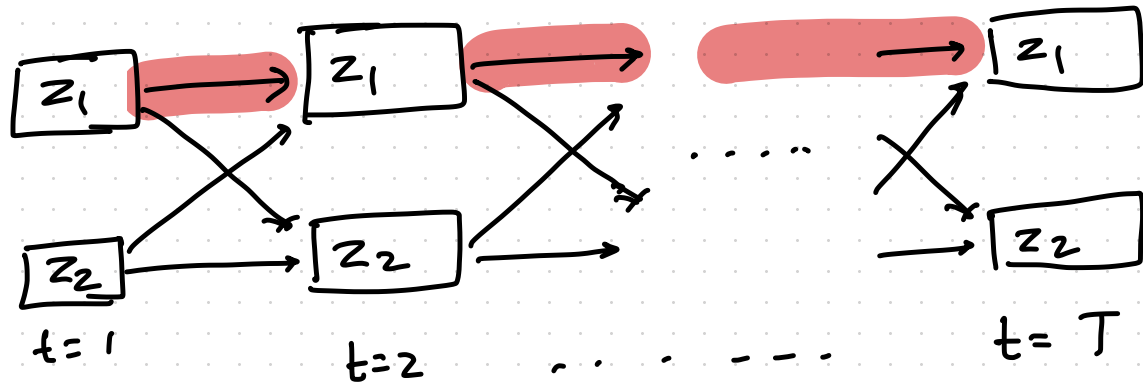
For $X = \{H, H, \dots\}$ what is $L(X|\theta)$ where $\theta = 2$
T timesteps.



SNO	1	2	...	TIME T
1	z_1	z_1	...	z_1
...				
K^T

Total K^T paths required for $L(X|\theta)$

For $X = \{H_1, H_2, \dots\}$ what is $L(X|\theta)$ where $\theta = \downarrow$
 T time steps.



length of each path = T
 $\Rightarrow O(T)$ multiplications

Total k^T paths
 required for
 $L(X|\theta)$

$P(x|\theta)$ computatⁿ is $O(TK^T)$

Can we do better? Hint: Dynamic Programming

FORWARD PROCEDURE

Define $\alpha_t(i) = P(x_{1:t}, z_t = i)$

FORWARD PROCEDURE

Define $\alpha_t(i) = P(x_{1:t}, z_t = i)$

i.e. $\alpha_t(i) =$ Probability of being in state
'i' at time t and
observatⁿ $x_{1:t}$

FORWARD PROCEDURE (Base Step $t=1$)

Define $\alpha_t(i) = P(x_{1:t}, z_t = i)$

For our example for $t=1$

$$\begin{aligned}\alpha_1(B) &= P(x_1, z_1 = \text{Biased}) = \pi_B \cdot P(x_1 = H | z_1 = B) \\ &= 0.6 * 0.7 = 0.42\end{aligned}$$

FORWARD PROCEDURE (Base step $t=1$)

Define $\alpha_t(i) = P(x_{1:t}, z_t = i)$

For our example for $t=1$

$$\begin{aligned}\alpha_1(\text{Bias}) &= P(x_1, z_1 = \text{Biased}) = \pi_B \cdot P(x_1 = H | z_1 = B) \\ &= 0.6 \times 0.7 = 0.42\end{aligned}$$

$$\begin{aligned}\alpha_1(\text{Fair}) &= P(x_1 = H, z_1 = F) = \pi_F \cdot P(x_1 = H | z_1 = F) \\ &= 0.4 \times 0.5 = 0.20\end{aligned}$$

FORWARD PROCEDURE (INDUCTION STEP)

Define $\alpha_t(i) = P(x_{1:t}, z_t = i)$

$$\alpha_1 = 0.42$$

$$\boxed{z_B}$$

$$\boxed{z_F}$$

$$\alpha_1 = 0.2$$

FORWARD PROCEDURE (INDUCTION STEP)

Define $\alpha_t(i) = P(x_{1:t}, z_t = i)$

$$\alpha_1 = 0.42$$

Z_B

$$\alpha_2(B) = ?$$

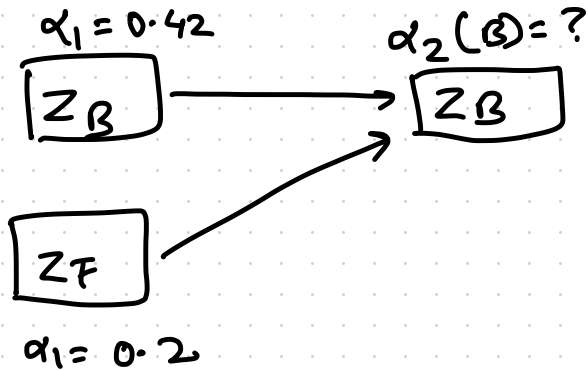
Z_B

Z_F

$$\alpha_1 = 0.2$$

FORWARD PROCEDURE (INDUCTION STEP)

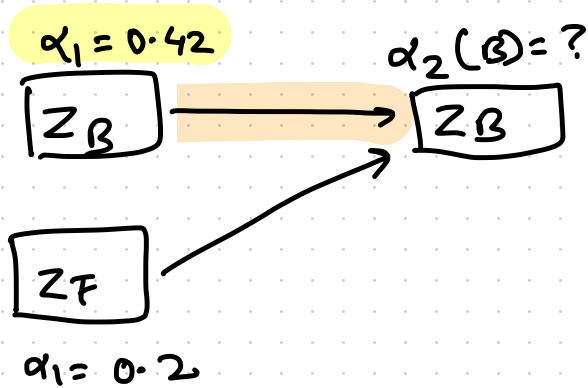
Define $\alpha_t(i) = P(x_{1:t}, z_t = i)$



2 paths (k)
to Bias state
at time
Step 2

FORWARD PROCEDURE (INDUCTION STEP)

Define $\alpha_t(i) = P(x_{1:t}, z_t = i)$

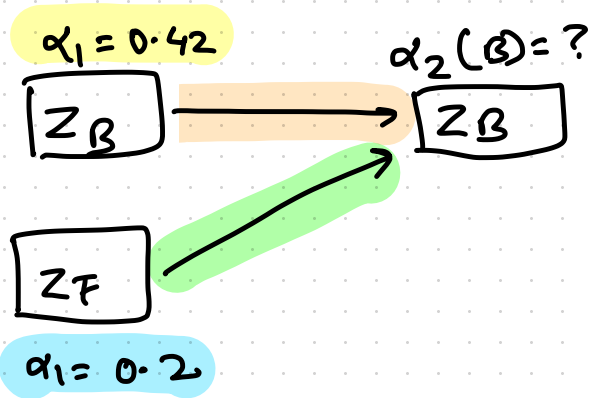


2 paths (k)
to Bias state
at time
Step 2

$$P(x_1 = H, x_2 = H, z_2 = B) = [P(x_1 = H, z_1 = B) \cdot A_{BB}] P(H|Bias)$$

FORWARD PROCEDURE (INDUCTION STEP)

Define $\alpha_t(i) = P(x_{1:t}, z_t = i)$



2 paths (k)
to Bias state
at time
Step 2

$$P(x_1 = H, x_2 = H, z_2 = B) = [P(x_1 = H, z_1 = B) \cdot A_{BB}] P(H|Bias) + [P(x_1 = H, z_1 = F) \cdot A_{FB}] P(H|B)$$

FORWARD PROCEDURE (INDUCTION STEP)

Define $\alpha_t(i) = P(x_{1:t}, z_t = i)$

$$P(x_1 = H, x_2 = H, z_2 = B) = [P(x_1 = H, z_1 = B) \cdot A_{BB}] P(H|B_{\text{bias}}) + [P(x_1 = H, z_1 = F) \cdot A_{FB}] P(H|F)$$

$$\Rightarrow \alpha_2(B) = (\alpha_1(B) A_{BB}) \phi_B(x=H) + (\alpha_1(F) A_{FB}) \phi_B(x=H)$$

$$= (0.42) * (0.9) * (0.7) + (0.2) * (0.1) * (0.7)$$

$$= 0.2646 + 0.014$$

$$= 0.2786$$

FORWARD PROCEDURE (INDUCTION STEP)

Define $\alpha_t(i) = P(x_{1:t}, z_t = i)$

$$\alpha_{t+1}(j) = \left[\sum_i \alpha_t(i) * A_{ij} \right] \phi_j(x_{t+1})$$

FORWARD PROCEDURE (INDUCTION STEP)

Define $\alpha_t(i) = P(x_{1:t}, z_t = i)$

$$\alpha_1 = 0.42$$

Z_B

Z_F

$$\alpha_1 = 0.2$$

$$\alpha_2(B) = 0.2786$$

Z_B

Z_F

$$\alpha_2(F) = 0.1110$$

FORWARD PROCEDURE (Termination)

$$\alpha_T(B) = P(x_{1:T}, z_T = B)$$

$$\alpha_T(F) = P(x_{1:T}, z_T = F)$$

$$\therefore P(x_{1:T}) = \alpha_T(B) + \alpha_T(F)$$

FORWARD PROCEDURE (Termination)

$$\alpha_T(B) = P(x_{1:T}, z_T = B)$$

$$\alpha_T(F) = P(x_{1:T}, z_T = F)$$

$$\therefore P(x_{1:T}) = \alpha_T(B) + \alpha_T(F)$$

In general,

$$P(x_{1:T} | \theta) = \sum_i \alpha_T(i)$$

or

$$L(x_{1:T} | \theta) = \sum_i \alpha_T(i)$$

P III HMM Filtering

Compute belief state

$p(z_t | x_{1:t})$ online (or recursively)

P III HMM Filtering

Compute belief state

$p(z_t | x_{1:t})$ online (or recursively)

$$= \frac{p(z_t, x_{1:t})}{p(x_{1:t})} = \frac{\alpha_t(i)}{\sum_i \alpha_t(i)}$$

P III HMM Filtering.

$$P(z_1 = B | x_1 = H) = \frac{\alpha_1(B)}{\alpha_1(B) + \alpha_1(F)} = \frac{0.42}{0.62} \approx \frac{2}{3}$$

$$P(z_1 = F | x_1 = H) = \frac{\alpha_1(F)}{\alpha_1(B) + \alpha_1(F)} = \frac{0.2}{0.62} \approx \frac{1}{3}$$

P III KMM Filtering.

$$P(z_2 = B | x = \{nn\}) = \frac{\alpha_2(B)}{\alpha_2(B) + \alpha_2(F)} = 0.715$$

$$P(z_2 = F | x = \{nn\}) = \frac{\alpha_2(F)}{\alpha_2(B) + \alpha_2(F)} = 0.289$$

P III HMM Filtering.

$$P(z_2 = B | x = \{hh\}) = \frac{\alpha_2(B)}{\alpha_2(B) + \alpha_2(F)} = 0.715$$

$$P(z_2 = F | x = \{hh\}) = \frac{\alpha_2(F)}{\alpha_2(B) + \alpha_2(F)} = 0.289$$

Note

$$P(z_2 = B | \{hh\}) > P(z_1 = B | \{h\})$$

\therefore Biased coin \Rightarrow more heads!

$$\underline{PIV} : P(X_{t+1:T} | Z_t = i, \theta)?$$

Backwards algorithm

$$\beta_t(i) = P(X_{t+1:T} | Z_t = i, \theta)$$

= Probability of sequence $X_{t+1:T}$
given state is 'i' at t^{th}
time instance

$$\underline{P^i} : P(X_{t+1:T} | Z_t = i, \theta)?$$

$$\beta_t(i) = P(X_{t+1:T} | Z_t = i, \theta)$$

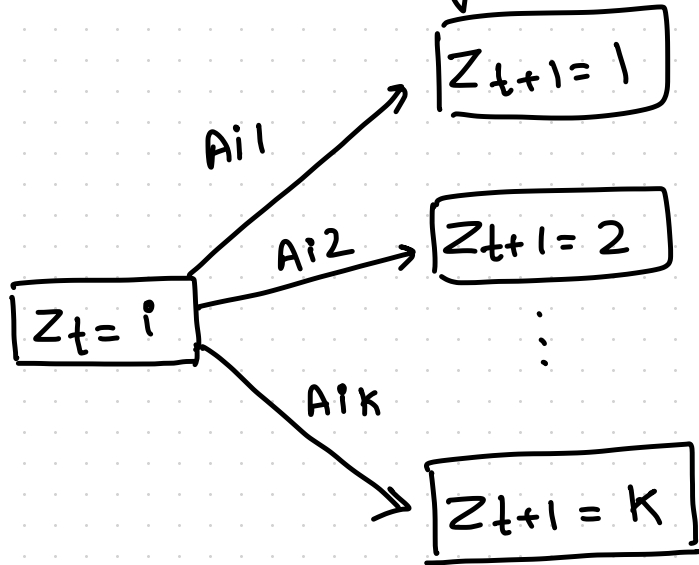
e.g. $X = \{H, H, H\}$

$$T = 3$$

$$t = 1$$

$$Z_1 = \text{Biased}$$

Backwards Algorithm



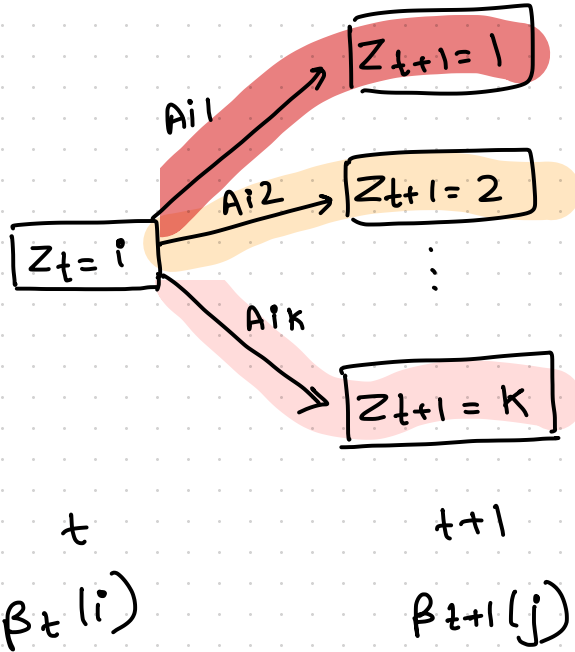
t

$\beta_t(i)$

$t+1$

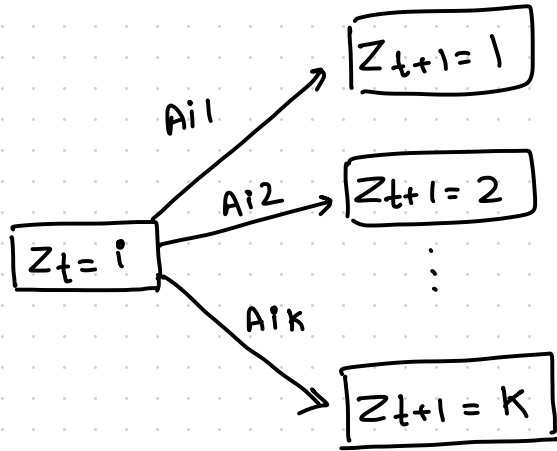
$\beta_{t+1}(j)$

Backwards Algorithm



$$\begin{aligned} \beta_t(i) &= P(x_{t+1:T} | Z_t = i) \\ &= P(x_{t+2:T} | Z_{t+1} = 1) \cdot A_{i1} \cdot \phi_1(x_{t+1}) \\ &+ P(x_{t+2:T} | Z_{t+1} = 2) \cdot A_{i2} \cdot \phi_2(x_{t+1}) \\ &\vdots \\ &+ P(x_{t+2:T} | Z_{t+1} = k) \cdot A_{ik} \cdot \phi_k(x_{t+1}) \end{aligned}$$

Backwards Algorithm



t

$\beta_t(i)$

$t+1$

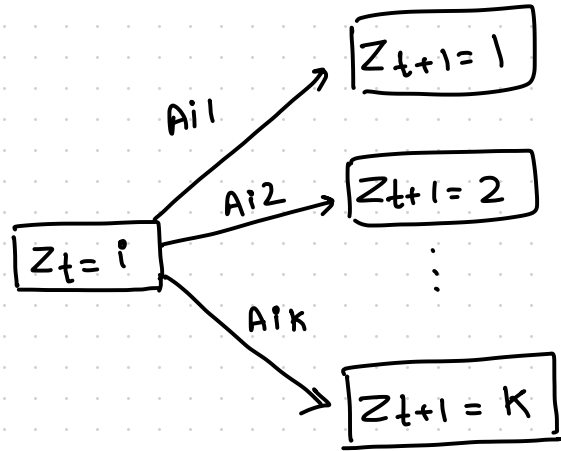
$\beta_{t+1}(j)$

INDUCTION STEP

$$\beta_t(i) = P(x_{t+1:T} | z_t = i)$$

$$= \sum_{j=1}^k \beta_{t+1}(j) \cdot A_{ij} \cdot \phi_j(x_{t+1})$$

Backwards Algorithm



t
 $\beta_t(i)$

$t+1$
 $\beta_{t+1}(j)$

INDUCTION STEP

$$\beta_t(i) = P(x_{t+1:T} | z_t = i)$$

$$= \sum_{j=1}^k \beta_{t+1}(j) \cdot A_{ij} \cdot \phi_j(x_{t+1})$$

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

INIT. STEP

$$\beta_T(i) = 1 \quad \forall i \in \{1, \dots, k\}$$

(Arbitrarily defined)

Backwards Algorithm (example)

$$X = \{H, H, H\}$$

Find $\beta_1, \beta_2, \beta_3$

INIT

$$\beta_3(B) = 1 = \beta_3(F) \quad (\text{Arbitrary})$$

Backwards Algorithm (example)

$$X = \{H, H, H\} \quad \text{Find } \beta_1, \beta_2, \beta_3$$

INIT

$$\beta_3(B) = 1 = \beta_3(F) \quad (\text{Arbitrary})$$

STEP

$$\beta_2(B) = \sum_{j \in \{B, F\}} \beta_3(j) \cdot A_{Bj} \cdot \phi_j(H)$$

$$\begin{aligned} &= \beta_3(B) \cdot A_{BB} \cdot \phi_B(H) + \beta_3(F) \cdot A_{BF} \cdot \phi_F(H) \\ &= 1 \cdot (0.9) \cdot (0.7) + (1) \cdot (0.1) \cdot (0.5) = 0.68 \end{aligned}$$

Backwards Algorithm (example)

$$X = \{H, H, H\} \quad \text{Find } \beta_1, \beta_2, \beta_3$$

INIT

$$\beta_3(B) = 1 = \beta_3(F) \quad (\text{Arbitrary})$$

STEP

$$\beta_2(B) = \sum_{j \in \{B, F\}} \beta_3(j) \cdot A_{Bj} \cdot \phi_j(H)$$

$$\begin{aligned} &= \beta_3(B) \cdot A_{BB} \cdot \phi_B(H) + \beta_3(F) \cdot A_{BF} \cdot \phi_F(H) \\ &= 1 \cdot (0.9) \cdot (0.7) + (1) (0.1) (0.5) = 0.68 \end{aligned}$$

$$\begin{aligned} \beta_2(F) &= \beta_3(B) \cdot A_{FB} \cdot \phi_B(H) + \beta_3(F) \cdot A_{FF} \cdot \phi_F(H) \\ &= (1) \cdot (0.1) (0.7) + (1) (0.9) (0.5) = 0.52 \end{aligned}$$

Backwards Algorithm (example)

$$X = \{H, H, H\} \quad \text{Find } \beta_1, \beta_2, \beta_3$$

INIT

$$\beta_3(B) = 1 = \beta_3(F) \quad (\text{Arbitrary})$$

STEP

$$\begin{aligned} \beta_1(B) &= \beta_2(B) \cdot A_{BB} \cdot \phi_B(H) + \beta_2(F) \cdot A_{BF} \cdot \phi_F(H) \\ &= (0.62)(0.9)(0.7) + (0.52)(0.1)(0.5) = 0.4544 \end{aligned}$$

$$\begin{aligned} \beta_1(F) &= \beta_2(B) \cdot A_{FB} \cdot \phi_B(H) + \beta_2(F) \cdot A_{FF} \cdot \phi_F(H) \\ &= (0.62)(0.1)(0.7) + (0.52)(0.9)(0.5) = 0.2816 \end{aligned}$$

PV 'OPTIMAL' sequence of states given model
and observations.

Given $x = \{x_1, \dots, x_T\}$

$\theta = \{\pi, A, \phi\}$

How to choose $z = \{z_1, \dots, z_T\}$
optimally?

PV 'OPTIMAL' sequence of states : MPM

OPTIMALITY DEFINITION #1

Choose state Z_t which is individually most likely

PV 'OPTIMAL' sequence of states: MPM

OPTIMALITY DEFINITION #1

Choose state z_t which is individually most likely

$$\hat{z} = \left(\underset{z_1}{\operatorname{arg\,max}} P(z_1 | x_{1:T}), \dots, \underset{z_T}{\operatorname{arg\,max}} P(z_T | x_{1:T}) \right)$$

PV 'OPTIMAL' sequence of states : MPM

OPTIMALITY DEFINITION #1

Choose state z_t which is individually most likely

$$\hat{z} = \left(\underset{z_1}{\operatorname{arg\,max}} P(z_1 | X_{1:T}), \dots, \underset{z_T}{\operatorname{arg\,max}} P(z_T | X_{1:T}) \right)$$

Also called Sequence of marginally most
Probable states (MPM)

PV 'OPTIMAL' sequence of states: mpm

OPTIMALITY DEFINITION #1

Choose state z_t which is individually most likely

$$\text{Let } \gamma_t(i) = p(z_t = i | x_{1:T})$$

PV 'OPTIMAL' sequence of states: mpm

OPTIMALITY DEFINITION #1

Choose state z_t which is individually most likely

$$\text{Let } \gamma_t(i) = p(z_t = i | x_{1:T})$$

$$\propto p(z_t = i | x_{1:t}) \cdot p(x_{t+1:T} | z_t = i, x_{1:t})$$

PV 'OPTIMAL' sequence of states: MPM

OPTIMALITY DEFINITION #1

Choose state z_t which is individually most likely

$$\text{Let } \gamma_t(i) = p(z_t = i | x_{1:T})$$

$$\propto p(z_t = i | x_{1:t}) \cdot p(x_{t+1:T} | z_t = i, x_{1:t})$$

Independent

$$\propto p(z_t = i | x_{1:t}) \cdot p(x_{t+1:T} | z_t = i)$$

$$\gamma_t(i) \propto \alpha_t(i) \cdot \beta_t(i)$$

PV 'OPTIMAL' sequence of states : MPM

FORWARD BACKWARD ALGORITHM

$$\gamma_t(i) = p(z_t = i | x_{1:T})$$

$$\gamma_t(i) \propto \alpha_t(i) \cdot \beta_t(i)$$

$$\gamma_t(i) = \frac{\alpha_t(i) \cdot \beta_t(i)}{\sum_i \alpha_t(i) \cdot \beta_t(i)}$$

PV 'OPTIMAL' sequence of states : MPM

FORWARD BACKWARD ALGORITHM

Example: For some θ , $x = \{H, H, H\}$ Find 'smooth'
set of states

$$\gamma_t(i) = \frac{\alpha_t(i) \cdot \beta_t(i)}{\sum_i \alpha_t(i) \cdot \beta_t(i)}$$

PV 'OPTIMAL' sequence of states: MPM

FORWARD BACKWARD ALGORITHM

Example: For some θ , $x = \{H, H, H\}$ Find 'smooth' set of states

$$\gamma_t(i) = \frac{\alpha_t(i) \cdot \beta_t(i)}{\sum_i \alpha_t(i) \cdot \beta_t(i)}$$

$$\begin{aligned} \gamma_1(B) &= \frac{\alpha_1(B) \cdot \beta_1(B)}{\alpha_1(B) \cdot \beta_1(B) + \alpha_1(F) \cdot \beta_1(F)} = \frac{0.42 \times 0.4544}{0.42 \times 0.4544 + 0.2 \times 0.2816} \\ &= \frac{0.190848}{0.247168} = 0.77 \end{aligned}$$

PV 'OPTIMAL' sequence of states: MPM

FORWARD BACKWARD ALGORITHM

Example: For some θ , $x = \{H, H, H\}$ Find 'smooth' set of states

$$\gamma_t(i) = \frac{\alpha_t(i) \cdot \beta_t(i)}{\sum_i \alpha_t(i) \cdot \beta_t(i)}$$

$$\gamma_1(B) = 0.77$$

$$\Rightarrow \gamma_1(F) = 0.23$$

Much more likely to have $z_1 = B$ if $\{H, H, H\}$ and θ given

PV 'OPTIMAL' sequence of states: MPM

FORWARD BACKWARD ALGORITHM

Example: For some Θ , $x = \{H, H, H\}$ Find 'smooth' set of states

$$\gamma_t(i) = \frac{\alpha_t(i) \cdot \beta_t(i)}{\sum_i \alpha_t(i) \cdot \beta_t(i)}$$

$$\hat{z} = (\arg\max_i \gamma_1(i), \dots) = (\text{Biased}, \text{Biased}, \text{Biased})$$

PV 'OPTIMAL' sequence of states: MAP

OPTIMALITY DEFINITION #2

CHOOSE MOST PROBABLE SEQUENCE OF STATES

$$z^* = \arg \max_{z_{1:T}} p(z_{1:T} | x_{1:T})$$

PV 'OPTIMAL' sequence of states: MAP vs MAPM

CONSIDER FOLLOWING 2 timesteps joint prob

	$Z_1 = 1$	$Z_1 = 2$
$Z_2 = 1$	0.04	0.3
$Z_2 = 2$	0.36	0.3

PV 'OPTIMAL' sequence of states: MAP vs MPM

CONSIDER FOLLOWING 2 timesteps joint prob

	$Z_1 = 1$	$Z_1 = 2$	
$Z_2 = 1$	0.04	0.3	0.34
$Z_2 = 2$	0.36	0.3	0.66
	.4	.6	

MPM = Sequence of Marginally MOST PROB. STATES

$$= \left(\begin{array}{l} Z_1 = 2 \\ \vdots \\ 0.66 > 0.4 \end{array} , \begin{array}{l} Z_2 = 2 \\ \vdots \\ 0.66 > 0.34 \end{array} \right)$$

PV 'OPTIMAL' sequence of states: MAP vs MAPM

CONSIDER FOLLOWING 2 timesteps joint prob

	$Z_1 = 1$	$Z_1 = 2$	
$Z_2 = 1$	0.04	0.3	0.34
$Z_2 = 2$	0.36	0.3	0.66
	.4	.6	

MAP = $(Z_1 = 1, Z_2 = 2)$

$\therefore .36$ is highest number

PV 'OPTIMAL' sequence of states: MAP

OPTIMALITY DEFINITION # 2

CHOOSE MOST PROBABLE SEQUENCE OF STATES

$$z^* = \underset{z_{1:T}}{\operatorname{arg\,max}} P(z_{1:T} | x_{1:T})$$

VITERBI ALGORITHM (Dynamic Programming)

VITERBI ALGORITHM

Define $\delta_t(i) = \max_{z_1, z_2, \dots, z_{t-1}} P[z_1, z_2, z_3, \dots, z_{t-1}, z_t = i, x_1, x_2, \dots, x_t | \theta]$

VITERBI ALGORITHM

Define $\delta_t(i) = \max_{z_1, z_2, \dots, z_{t-1}} P[z_1, z_2, z_3, \dots, z_{t-1}, z_t = i, x_1, x_2, \dots, x_t | \theta]$

Best score (highest prob.) along a single path at time t , which accounts for first t observations and ends in $z_t = i$

Difference Blw $\alpha_t(i)$ and $\delta_t(i)$

$$\delta_t(i) = \max_{z_1, z_2, \dots, z_{t-1}} P[z_1, z_2, z_3, \dots, z_{t-1}, z_t = i, x_1, x_2, \dots, x_t | \theta]$$

$$\alpha_t(i) = P(x_{1:t}, z_t = i | \theta)$$

Difference Blw $\alpha_t(i)$ and $\delta_t(i)$

$$\delta_t(i) = \max_{z_1, z_2, \dots, z_{t-1}} P[z_1, z_2, z_3, \dots, z_{t-1}, z_t = i, x_1, x_2, \dots, x_t | \theta]$$

$$\alpha_t(i) = P[x_{1:t}, z_t = i | \theta]$$

FOCUS ON MOST LIKELY SEQUENCE

FOCUS ON MOST LIKELY STATE AT 't'

VITERBI ALGORITHM

Define $\delta_t(i) = \max_{z_1, z_2, \dots, z_{t-1}} P[z_1, z_2, z_3, \dots, z_{t-1}, z_t = i, x_1, x_2, \dots, x_t | \theta]$

Relation b/w $\delta_t(i)$ AND $\delta_{t+1}(j)$?

VITERBI ALGORITHM

Define $\delta_t(i) = \max_{z_1, z_2, \dots, z_{t-1}} P[z_1, z_2, z_3, \dots, z_{t-1}, z_t = i, x_1, x_2, \dots, x_t | \theta]$

Relation b/w $\delta_t(i)$ AND $\delta_{t+1}(j)$?

$$z_1 = 1$$

$$z_2 = 1$$

...

$$z_t = 1$$

$$z_{t+1} = 1$$

$$z_1 = 2$$

$$z_2 = 2$$

$$z_t = 2$$

$$z_{t+1} = 2$$

VITERBI ALGORITHM

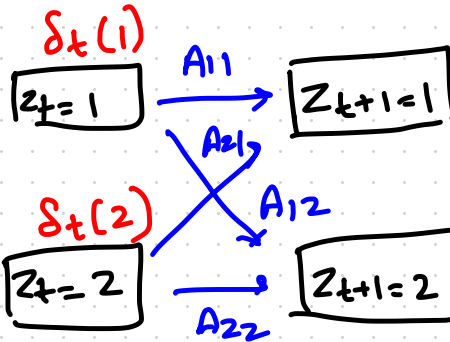
Define $\delta_t(i) = \max_{z_1, z_2, \dots, z_{t-1}} P[z_1, z_2, z_3, \dots, z_{t-1}, z_t = i, x_1, x_2, \dots, x_t | \theta]$

Relation b/w $\delta_t(i)$ AND $\delta_{t+1}(j)$?

$z_1 = 1$

$z_2 = 1$

...

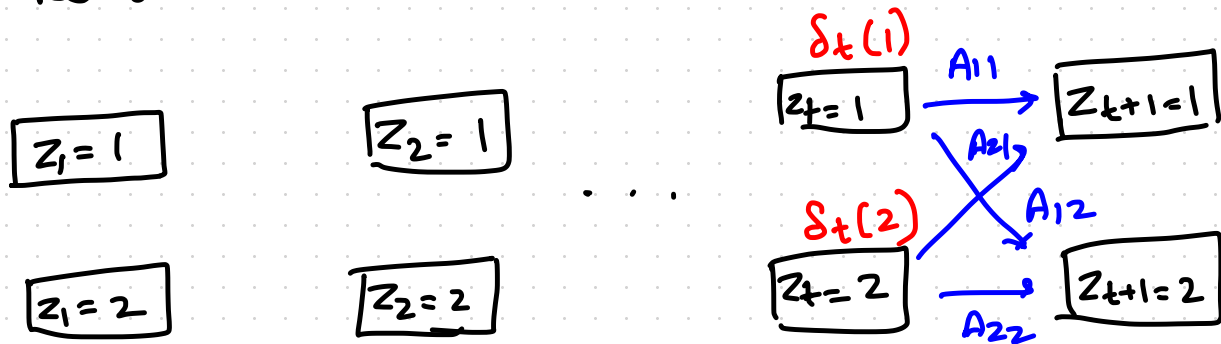


$z_1 = 2$

$z_2 = 2$

VITERBI ALGORITHM

Relation b/w $\delta_t(i)$ AND $\delta_{t+1}(j)$?

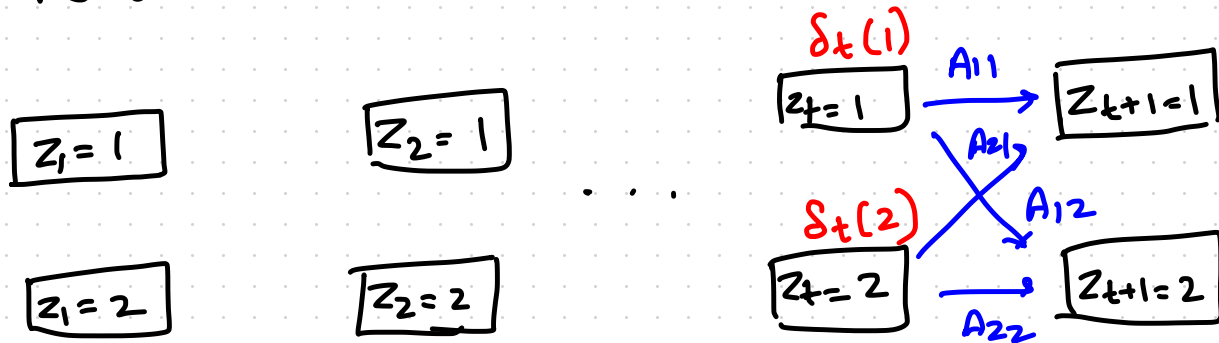


* We could reach $z_{t+1}=j$ from any $i \in \{1, \dots, K\}$ via a transition with prob. A_{ij}

* Once we reach $z_{t+1}=j$, prob. of observing x_{t+1} is $\phi_j(x_{t+1})$

VITERBI ALGORITHM

Relation b/w $\delta_t(i)$ AND $\delta_{t+1}(j)$?



* We could reach $z_{t+1}=j$ from any $i \in \{1, \dots, k\}$ via a transition with prob. A_{ij} . We take 'best path'

* Once we reach $z_{t+1}=j$, prob. of observing x_{t+1} is

$$\phi_j^o(x_{t+1})$$

$$\delta_{t+1}(j) = \left[\max_i \{ \delta_t(i) A_{ij} \} * \phi_j^o(x_{t+1}) \right]$$

VITERBI ALGORITHM

$$\delta_{t+1}(j) = \left[\max_i \{ \delta_t(i) \cdot A_{ij} \} * \phi_j(x_{t+1}) \right]$$

FOR EACH t and j , we need to store argument ' i ' which maximized above equation.

(called $\psi_t(j)$)

VITERBI ALGORITHM

① INITIALISATION

$$\delta_1(i) =$$

VITERBI ALGORITHM

① INITIALISATION

$$\delta_1(i) = \pi_i * \phi_i(x_1)$$

$$\psi_1(i) = 0 \quad (\text{ARBITRARY})$$

VITERBI ALGORITHM

① INITIALISATION

$$\delta_1(i) = \pi_i * \phi_i(x_1)$$

$$\psi_1(i) = 0 \quad (\text{ARBITRARY})$$

② RECURSION

FOR t in 2 to T

FOR j in 1 to K

$$\delta_t(j) = \left(\max_{i \in \{1, \dots, K\}} \left(\delta_{t-1}(i) \cdot A_{ij} \right) \right) \phi_j(x_t)$$

$$\psi_t(j) = \underset{i \in \{1, \dots, K\}}{\operatorname{argmax}} \delta_{t-1}(i) \cdot A_{ij}$$

VITERBI ALGORITHM

③ Termination

$$p^* = \max_{i \in \{1, \dots, K\}} \delta_T(i)$$

$$z_T^* = \operatorname{argmax}_{i \in \{1, \dots, K\}} \delta_T(i)$$

VITERBI ALGORITHM

③ Termination

$$p^* = \max_{i \in \{1, \dots, K\}} \delta_T(i)$$

$$z_T^* = \operatorname{argmax}_{i \in \{1, \dots, K\}} \delta_T(i)$$

④ Backtracking

For t in $T-1, \dots, 1$:

$$z_t^* = \psi_{t+1}(z_{t+1}^*)$$

VITERBI ALGORITHM

Example

Given $\Theta = \{\pi, A, \phi\}$ and $x_{1:T} = \{H, H, H\}$

determine MAP z^*

$$\pi = \begin{bmatrix} 0.6 & 0.4 \\ \text{Bias} & \text{Fair} \end{bmatrix}$$

$$A = \begin{array}{c} \text{Bias} \\ \text{Fair} \end{array} \begin{array}{cc} \text{B} & \text{F} \\ \left[\begin{array}{cc} 0.9 & 0.1 \\ 0.1 & 0.9 \end{array} \right] \end{array}$$

$$\phi = \begin{array}{cc} \text{Bias} & \text{Fair} \\ p(H) = 0.7 & p(H) = 0.5 \\ p(T) = 0.3 & p(T) = 0.5 \end{array}$$

VITERBI ALGORITHM

Example

Given $\Theta = \{\pi, A, \phi\}$ and $x_{1:T} = \{H, H, H\}$

determine MAP z^*

$$\delta_1(\text{Bias}) = \pi_{\text{Bias}} * \phi_{\text{Bias}}(H) = 0.6 * 0.7 = 0.42$$

$$\delta_1(\text{Fair}) = \pi_{\text{Fair}} * \phi_{\text{Fair}}(H) = 0.4 * 0.5 = 0.20$$

VITERBI ALGORITHM

Example

Given $\Theta = \{\pi, A, \phi\}$ and $x_{1:T} = \{H, H, H\}$

determine MAP z^*

$$\delta_1(\text{Bias}) = \pi_{\text{Bias}} * \phi_{\text{Bias}}(H) = 0.6 * 0.7 = 0.42$$

$$\delta_1(\text{Fair}) = \pi_{\text{Fair}} * \phi_{\text{Fair}}(H) = 0.4 * 0.5 = 0.20$$

$$\psi_1(\text{Bias}) = \psi_1(\text{Fair}) = 0 \text{ (Arbitrary)}$$

VITERBI ALGORITHM

Example

$$\delta_1(\text{Bias}) = \pi_{\text{Bias}} * \phi_{\text{Bias}}(H) = 0.6 * 0.7 = 0.42$$

$$\delta_1(\text{Fair}) = \pi_{\text{Fair}} * \phi_{\text{Fair}}(H) = 0.4 * 0.5 = 0.20$$

$$\psi_1(\text{Bias}) = \psi_1(\text{Fair}) = 0 \text{ (Arbitrary)}$$

$$\delta_2(\text{Bias}) = \max \left\{ \begin{array}{l} \delta_1(\text{Fair}) \cdot A_{\text{Fair}, \text{Bias}} \\ \delta_1(\text{Bias}) \cdot A_{\text{Bias}, \text{Bias}} \end{array} \right\} * \phi_{\text{Bias}}(H)$$

$$= \max \left\{ \begin{array}{l} 0.2 * 0.1 \\ 0.42 * 0.9 \end{array} \right\} * 0.7 = 0.42 * 0.9 * 0.7$$

VITERBI ALGORITHM

$$\delta_1(\text{Bias}) = \pi_{\text{Bias}} * \phi_{\text{Bias}}(H) = 0.6 * 0.7 = 0.42$$

$$\delta_1(\text{Fair}) = \pi_{\text{Fair}} * \phi_{\text{Fair}}(H) = 0.4 * 0.5 = 0.20$$

$$\psi_1(\text{Bias}) = \psi_1(\text{Fair}) = 0 \text{ (Arbitrary)}$$

$$\begin{aligned} \delta_2(\text{Bias}) &= \max \left\{ \begin{array}{l} \delta_1(\text{Fair}) \cdot A_{\text{Fair}, \text{Bias}} \\ \delta_1(\text{Bias}) \cdot A_{\text{Bias}, \text{Bias}} \end{array} \right\} * \phi_{\text{Bias}}(H) \\ &= \max \left\{ \begin{array}{l} 0.2 * 0.1 \\ 0.42 * 0.9 \end{array} \right\} * 0.7 = 0.42 * 0.9 * 0.7 \\ &= 0.26 \end{aligned}$$

$$\psi_2(\text{Bias}) = \text{Bias}$$

VITERBI ALGORITHM

$$\delta_1(\text{Bias}) = \pi_{\text{Bias}} * \phi_{\text{Bias}}(H) = 0.6 * 0.7 = 0.42$$

$$\delta_1(\text{Fair}) = \pi_{\text{Fair}} * \phi_{\text{Fair}}(H) = 0.4 * 0.5 = 0.20$$

$$\psi_1(\text{Bias}) = \psi_1(\text{Fair}) = 0 \text{ (Arbitrary)}$$

$$\delta_2(\text{Fair}) = \max \left\{ \begin{array}{l} \delta_1(\text{Fair}) \cdot A_{\text{Fair}, \text{Fair}} \\ \delta_1(\text{Bias}) \cdot A_{\text{Bias}, \text{Fair}} \end{array} \right\} * \phi_{\text{Fair}}(H)$$

$$= \max \left\{ \begin{array}{l} 0.2 * 0.9 \\ 0.42 * 0.1 \end{array} \right\} * 0.5 = 0.2 * 0.9 * 0.5 = 0.09$$

$$\psi_2(\text{Fair}) = \text{Fair}$$

VITERBI ALGORITHM

$$\delta_2(\text{Bias}) = 0.2646$$

$$\delta_2(\text{Fair}) = 0.09$$

$$\psi_2(\text{Bias}) = \text{Bias}$$

$$\psi_2(\text{Fair}) = \text{Fair}$$

$$\delta_3(\text{Bias}) = \max \left\{ \begin{array}{l} \delta_2(\text{Bias}) * A_{\text{Bias}, \text{Bias}} \\ \delta_2(\text{Fair}) * A_{\text{Fair}, \text{Bias}} \end{array} \right\} * \phi_{\text{Bias}}(H)$$

$$\delta_3(\text{Bias}) = \max \left\{ \begin{array}{l} 0.2646 * 0.9 \\ 0.09 * 0.1 \end{array} \right\} * 0.7 = 0.166698$$

$$\psi_3(\text{Bias}) = \text{Bias}$$

VITERBI ALGORITHM

$$\delta_2(\text{Bias}) = 0.2646$$

$$\delta_2(\text{Fair}) = 0.09$$

$$\psi_2(\text{Bias}) = \text{Bias}$$

$$\psi_2(\text{Fair}) = \text{Fair}$$

$$\delta_3(\text{Fair}) = \max \left\{ \begin{array}{l} \delta_2(\text{Fair}) \cdot A_{\text{Fair}, \text{Fair}} \\ \delta_2(\text{Bias}) \cdot A_{\text{Bias}, \text{Fair}} \end{array} \right\} * \phi_{\text{Fair}}(H)$$

$$= \max \left\{ \begin{array}{l} 0.09 * 0.9 \\ 0.2646 * 0.1 \end{array} \right\} * 0.5 = \max \left\{ \begin{array}{l} 0.081 \\ 0.02646 \end{array} \right\} * 0.5 \\ = 0.0405$$

$$\psi_3(\text{Fair}) = \text{Fair}$$

VITERBI ALGORITHM

$$\delta_0(\text{Bias}) = 0.42$$

$$\delta_0(\text{Fair}) = 0.2$$

$$\delta_1(\text{Bias}) = 0.2646$$

$$\delta_1(\text{Fair}) = 0.09$$

$$\delta_2(\text{Bias}) = 0.167$$

$$\delta_2(\text{Fair}) = 0.04$$

$$\psi_1(\text{Bias}) = \text{Bias}$$

$$\psi_1(\text{Fair}) = \text{Fair}$$

$$\psi_2(\text{Bias}) = \text{Bias}$$

$$\psi_2(\text{Fair}) = \text{Fair}$$

VITERBI ALGORITHM

$$\delta_0(\text{Bias}) = 0.42$$

$$\delta_0(\text{Fair}) = 0.2$$

$$\delta_1(\text{Bias}) = 0.2646$$

$$\delta_1(\text{Fair}) = 0.09$$

$$\delta_2(\text{Bias}) = 0.167$$

$$\delta_2(\text{Fair}) = 0.04$$

$$\psi_1(\text{Bias}) = \text{Bias}$$

$$\psi_1(\text{Fair}) = \text{Fair}$$

$$\psi_2(\text{Bias}) = \text{Bias}$$

$$\psi_2(\text{Fair}) = \text{Fair}$$

$$z_3^* = \operatorname{argmax} [\delta_3(\text{Bias}), \delta_3(\text{Fair})] = \text{Bias}$$

VITERBI ALGORITHM

$$\delta_0(\text{Bias}) = 0.42$$

$$\delta_0(\text{Fair}) = 0.2$$

$$\delta_1(\text{Bias}) = 0.2646$$

$$\delta_1(\text{Fair}) = 0.09$$

$$\delta_2(\text{Bias}) = 0.167$$

$$\delta_2(\text{Fair}) = 0.04$$

$$\psi_1(\text{Bias}) = \text{Bias}$$

$$\psi_1(\text{Fair}) = \text{Fair}$$

$$\psi_2(\text{Bias}) = \text{Bias}$$

$$\psi_2(\text{Fair}) = \text{Fair}$$

$$z_3^* = \operatorname{argmax} [\delta_3(\text{Bias}), \delta_3(\text{Fair})] = \text{Bias}$$

$$z_2^* = \psi_3(\text{Bias}) = \text{Bias}$$

$$z_1^* = \psi_2(z_2^*) = \psi_2(\text{Bias}) = \text{Bias} \Rightarrow z^* = \left. \begin{array}{l} \text{Bias, Bias} \\ \text{Bias} \end{array} \right\}$$

PVI : HMM PARAMETER LEARNING GIVEN
FULLY OBSERVED DATA

* Fully observed: observe all hidden state sequences

* $Q = \{\pi, A, \phi\}$

PVI: HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

* Fully observed: observe all hidden state sequences

* $\Theta = \{\pi, A, \Phi\}$

* Given $D = \left\{ \begin{aligned} &((z_{11}, z_{12}, \dots, z_{1T_1}), (x_{11}, x_{12}, \dots, x_{1T_1})), \\ &\vdots \\ &((z_{N1}, \dots, z_{NT_N}), (x_{N1}, \dots, x_{NT_N})) \end{aligned} \right\}$

PVI: HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

* Fully observed: observe all hidden state sequences

* $\theta = \{\pi, A, \phi\}$

* Given $D = \left\{ \left((z_{11}, z_{12}, \dots, z_{1T_1}), (x_{11}, x_{12}, \dots, x_{1T_1}) \right), \right.$

N sequences \leftarrow

$\left. \left((z_{N1}, \dots, z_{NT_N}), (x_{N1}, \dots, x_{NT_N}) \right) \right\}$

PVI: HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

* Fully observed: observe all hidden state sequences

* $\theta = \{\pi, A, \phi\}$

* Given $D = \left\{ \left((z_{11}, z_{12} \dots z_{1T_1}), (x_{11}, x_{12} \dots x_{1T_1}) \right), \right.$
 \vdots

N sequences ←
each of
 T_i length
(may be different)

$\left((z_{N1}, \dots z_{NT_N}), (x_{N1}, \dots x_{NT_N}) \right)$

PVI : HMM PARAMETER LEARNING GIVEN
FULLY OBSERVED DATA

SOME MATHS FOR A SINGLE SEQUENCE FIRST

Likelihood of a sequence $Z_{1:T}$

PVI : HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

Likelihood of a sequence $z_{1:T}$

$$= \pi(z_1) A(z_1, z_2) \dots A(z_{T-1}, z_T)$$

PVI : HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

Likelihood of a sequence $z_{1:T}$

$$= \pi(z_1) A(z_1, z_2) \dots A(z_{T-1}, z_T)$$

Let us write $\pi(z_1)$ in terms of $j \in \{1, \dots, K\}$

PVI : HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

Likelihood of a sequence $z_{1:T}$

$$= \pi(z_1) A(z_1, z_2) \dots A(z_{T-1}, z_T)$$

Let us write $\pi(z_1)$ in terms of $j \in \{1, \dots, K\}$

$$\pi = \left[\begin{array}{cccc} \dots & \dots & \dots & \dots \\ z_1=1 & z_1=2 & z_1=j & z_1=K \end{array} \right]$$

$$\pi(z_1) = \pi_j : \text{if } z_1=j$$

PVI : HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

Likelihood of a sequence $z_{1:T}$

$$= \pi(z_1) A(z_1, z_2) \dots A(z_{T-1}, z_T)$$

Let us write $\pi(z_1)$ in terms of $j \in \{1, \dots, k\}$

$$\pi = \left[\begin{array}{cccc} \dots & \dots & \dots & \dots \\ z_1=1 & z_1=2 & z_1=j & z_1=k \end{array} \right]$$

$$\pi(z_1) = \pi_j : \text{if } z_1=j$$

$$\text{or } \pi(z_1) = \prod_{j=1}^k (\pi_j)^{\mathbb{I}(z_1=j)}$$

PVI : HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

Likelihood of a sequence $z_{1:T}$

$$= \pi(z_1) A(z_1, z_2) \dots A(z_{T-1}, z_T)$$

Let us write $A(z_1, z_2)$ in terms of $j, k \in \{1, \dots, K\}$

PVI : HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

Likelihood of a sequence $z_{1:T}$

$$= \pi(z_1) A(z_1, z_2) \dots A(z_{T-1}, z_T)$$

Let us write $A(z_1, z_2)$ in terms of $j, k \in \{1, \dots, K\}$

$$A(z_1, z_2) = \prod_{j=1}^K \prod_{k=1}^K (A_{jk})^{I(z_1=j, z_2=k)}$$

PVI: HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

Likelihood of a sequence $z_{1:T}$

$$= \pi(z_1) A(z_1, z_2) \dots A(z_{T-1}, z_T)$$

Now, $A(z_1, z_2) = \prod_{j=1}^K \prod_{k=1}^K (A_{jk})^{I(z_1=j, z_2=k)}$

$$\Rightarrow A(z_1, z_2) \dots A(z_{T-1}, z_T) = \prod_{t=2}^T \prod_{j=1}^K \prod_{k=1}^K (A_{jk})^{I(z_t=k, z_{t-1}=j)}$$

PVI : HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

Likelihood of a sequence $Z_{1:T}$

$$= \pi(z_1) A(z_1, z_2) \dots A(z_{T-1}, z_T)$$

$$LL(O|D = \{A, \pi\}) = \log \left(\prod_{j=1}^K (\pi_j)^{I(z_1=j)} \prod_{t=2}^T \prod_{j=1}^K \prod_{k=1}^K (A_{jk})^{I(z_{t-1}=j, z_t=k)} \right)$$

$$= \log(F \cdot G)$$

$$= \log F + \log G$$

PVI : HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

$$LL(O|\theta = \{A, \pi\}) = \log \left(\prod_{j=1}^K (\pi_j)^{I(z_1=j)} \prod_{t=2}^T \prod_{j=1}^K \prod_{k=1}^K (A_{jk})^{I(z_{t-1}=j, z_t=k)} \right)$$

$$= \log(F \cdot G)$$

$$= \log F + \log G$$

$$\text{Now } \log F = \log \left(\pi_1^{I(z_1=1)} \pi_2^{I(z_1=2)} \dots \pi_k^{I(z_1=k)} \right)$$

PVI : HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

$$LL(O|\theta = \{A, \pi\}) = \log \left(\prod_{j=1}^K (\pi_j)^{I(z_1=j)} \prod_{t=2}^T \prod_{j=1}^K \prod_{k=1}^K (A_{jk})^{I(x_{t-1}=j, x_t=k)} \right)$$

$$= \log(F \cdot G)$$

$$= \log F + \log G$$

$$\text{Now } \log F = \log \left(\pi_1^{I(z_1=1)} \pi_2^{I(z_1=2)} \dots \pi_k^{I(z_1=k)} \right)$$

$$= I(z_1=1) \log \pi_1 + \dots$$

$$\log F = \sum_{j=1}^K I(z_1=j) \log \pi_j$$

PVI : HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

$$LL(O|\theta = \{A, \pi\}) = \log \left(\prod_{j=1}^K (\pi_j)^{I(z_1=j)} \prod_{t=2}^T \prod_{j=1}^K \prod_{k=1}^K (A_{jk})^{I(z_{t-k}, z_{t-1}=j)} \right)$$

$$= \log(F \cdot G)$$

$$\log G = \log \left(\prod_{j=1}^K \prod_{k=1}^K A_{jk}^{I(z_1=j, z_2=k)} \times \prod_{j=1}^K \prod_{k=1}^K A_{jk}^{I(z_2=j, z_3=k)} \dots \right)$$

PVI: HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

$$LL(O|\theta = \{A, \pi\}) = \log \left(\prod_{j=1}^K (\pi_j)^{I(z_1=j)} \prod_{t=2}^T \prod_{j=1}^K \prod_{k=1}^K (A_{jk})^{I(z_t=k, z_{t-1}=j)} \right)$$

$$= \log(F \cdot G)$$

$$\log G = \log \left(\prod_{j=1}^K \prod_{k=1}^K A_{jk}^{I(z_1=j, z_2=k)} \times \prod_{j=1}^K \prod_{k=1}^K A_{jk}^{I(z_2=j, z_3=k)} \dots \right)$$

$$= \sum_{t=2}^T \log \left(\prod_{j=1}^K \prod_{k=1}^K A_{jk}^{I(z_t=k, z_{t-1}=j)} \right)$$

PVI: HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

$$\log G = \log \left(\prod_{j=1}^K \prod_{k=1}^K A_{jk}^{I(x_1=j, x_2=k)} \times \prod_{j=1}^K \prod_{k=1}^K A_{jk}^{I(x_2=j, x_3=k)} \dots \right)$$

$$= \sum_{t=2}^T \log \left(\prod_{j=1}^K \prod_{k=1}^K A_{jk}^{I(x_t=k, x_{t+1}=j)} \right)$$

$$\log G = \sum_{t=2}^T \sum_{j=1}^K \sum_{k=1}^K I(x_t=k, x_{t+1}=j) \log A_{jk}$$

PVI: HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

BACK TO 'N' sequences now.

* Given $D = \left\{ (z_{11}, z_{12} \dots z_{1T_1}), (x_{11}, x_{12} \dots x_{1T_1}) \right\},$

\vdots

N sequences
each of
 T_i length
(may be different)

$((z_{N1}, \dots z_{NT_N}), (x_{N1}, \dots x_{NT_N}))$

PVI : HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

STEP I

Estimate π and A given $D' = ((z_{i1} \dots z_{iT_i}) \forall i \in (1, N))$

$$\log P(D' | \theta' = \{A, \pi\}) = \sum_{i=1}^N \log P(z_i | \theta')$$

PVI: HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

STEP I

Estimate π and A given $D' = ((z_{i1} \dots z_{iT_i}) \forall i \in (1, N))$

$$\log P(D' | \theta' = \{A, \pi\}) = \sum_{i=1}^N \log P(z_i | \theta')$$

$$N_j \stackrel{\Delta}{=} \sum_{i=1}^N I(z_{i1} = j)$$

$$N_{j,k} \stackrel{\Delta}{=} \sum_{i=1}^N \sum_{t=1}^{T_i-1} I(x_{i,t} = j, x_{i,t+1} = k)$$

PVI: HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

STEP I

Estimate π and A given $D' = ((z_{i1} \dots z_{iT_i}) \forall i \in (1, N))$

$$\log P(D' | \theta' = \{A, \pi\}) = \sum_{i=1}^N \log P(z_i | \theta')$$

$$N_j \stackrel{\Delta}{=} \sum_{i=1}^N I(z_i = j) \leftarrow \text{COUNT OF SEQUENCES WHERE } z_i = j$$

$$N_{j,k} \stackrel{\Delta}{=} \sum_{i=1}^N \sum_{t=1}^{T_i-1} I(x_{i,t} = j, x_{i,t+1} = k)$$

PVI: HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

STEP I

Estimate π and A given $D' = ((z_{i1}, \dots, z_{iT_i}) \forall i \in (1, N))$

$$\log P(D' | \theta' = \{A, \pi\}) = \sum_{i=1}^N \log P(z_i | \theta')$$

$$N_j \stackrel{\Delta}{=} \sum_{i=1}^D I(z_i = j) \leftarrow \text{COUNT OF SEQUENCES WHERE } z_i = j$$

COUNT OF TRANSITIONS FROM 'j' to 'k'

$$N_{j,k} \stackrel{\Delta}{=} \sum_{i=1}^N \sum_{t=1}^{T_i-1} I(x_{i,t} = j, x_{i,t+1} = k)$$

PVI : HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

FROM MLE,

$$\hat{\pi}_j = \frac{N_j^1}{\sum_{i=1}^K N_j^i}$$

$$\hat{A}_{jk} = \frac{N_{jk}}{\sum_k N_{jk}}$$

$N_j^1 \triangleq I(Z_1 = j)$ ← COUNT OF SEQUENCES WHERE $Z_1 = j$

COUNT OF TRANSITIONS FROM 'j' to 'k'

→ $N_{jk} \triangleq \sum_{i=1}^N \sum_{t=1}^{T_i-1} I(x_{i,t} = j, x_{i,t+1} = k)$

PVI: HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

FROM MLE,

$$\hat{\pi}_j = \frac{N_j}{\sum_{j=1}^K N_j}$$

$$\hat{A}_{jk} = \frac{N_{jk}}{\sum_k N_{jk}}$$

HOW TO ESTIMATE ϕ :

① ASSUME MULTINOULLI EMISSION

$$\Rightarrow \phi_{jl} = P(x_t = l | z_t = j) \quad l \in \{1, \dots, L\}$$

PVI: HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

FROM MLE,

$$\hat{\pi}_j = \frac{N_j^i}{\sum_{j=1}^K N_j^i}$$

$$\hat{A}_{jk} = \frac{N_{jk}}{\sum_k N_{jk}}$$

HOW TO ESTIMATE ϕ :

① ASSUME MULTINOULLI EMISSION

$$\Rightarrow \phi_{jl} = p(x_t = l | z_t = j) \quad l \in \{1, \dots, L\}$$

FROM MLE,

$$\hat{\phi}_{jl} = \frac{N_{jl}^x}{N_j}$$

$$N_{jl}^x \triangleq \sum_{i=1}^N \sum_{t=1}^{T_i} I(z_{i,t} = j, x_{i,t} = l)$$

PVI: HMM PARAMETER LEARNING GIVEN FULLY OBSERVED DATA

FROM MLE,

$$\hat{\pi}_j = \frac{N_j}{\sum_{j=1}^K N_j}$$

$$\hat{A}_{jk} = \frac{N_{jk}}{\sum_k N_{jk}}$$

HOW TO ESTIMATE ϕ FOR NORMALLY DISTRIBUTED ϕ

$$\hat{\mu}_k = \frac{\overline{x}_k}{N_k} \quad \hat{\Sigma}_k = \frac{\left(\overline{xx}\right)_k - N_k \hat{\mu}_k \hat{\mu}_k^T}{N_k}$$

$$\overline{x}_k \triangleq \sum_{i=1}^N \sum_{t=1}^{T_i} I(z_{i,t}=k) x_{i,t}$$

$$\left(\overline{xx}\right)_k \triangleq \sum_{i=1}^N \sum_{t=1}^{T_i} I(z_{i,t}=k) x_{i,t} x_{i,t}^T$$

HMM

LEARNING

EXAMPLE

(ESTIMATE π, A, Φ)

s_1

Bias

Fair

Fair

Bias

H

T

H

H

s_2 :

F

F

F

F

H

T

T

T

s_3 :

B

B

B

B

T

H

H

H

HMM LEARNING EXAMPLE (ESTIMATE π, A, Φ)

	Bias	Fair	Fair	Bias
s_1	H	T	H	H
s_2 :	F	F	F	F
	H	T	T	T
s_3 :	B	B	B	B
	T	H	H	H

$$N_j^t = \# \text{ Times } z_j \text{ occurs at } t=1$$

$$= \sum_{i=1}^N I(z_i = j)$$

HMM LEARNING EXAMPLE (ESTIMATE π, A, Φ)

	Bias	Fair	Fair	Bias
s_1	H	T	H	H
s_2 :	F	F	F	F
	H	T	T	T
s_3 :	B	B	B	B
	T	H	H	H

$N_j^1 = \# \text{ Times } z_j \text{ occurs at } t=1$

$$= \sum_{i=1}^N I(z_i = j)$$

$$N_B^1 = 2 \text{ (} s_1 \text{ and } s_3 \text{)}$$

$$N_F^1 = 1 \text{ (} s_2 \text{)}$$

HMM LEARNING EXAMPLE (ESTIMATE π, A, Φ)

	Bias	Fair	Fair	Bias
s_1	H	T	H	H
s_2 :	F	F	F	F
	H	T	T	T
s_3 :	B	B	B	B
	T	H	H	H

$$N'_B = 2 \quad , \quad N'_F = 1$$

$$\hat{\pi} = \left[\hat{\pi}_B, \hat{\pi}_F \right] = \left[\frac{N'_B}{N'_B + N'_F}, \frac{N'_F}{N'_F + N'_B} \right] = \left[\frac{2}{3}, \frac{1}{3} \right]$$

HMM LEARNING EXAMPLE (ESTIMATE π, A, ϕ)

s_1	Bias	Fair	fair	Bias
	H	T	H	H
s_2 :	F	F	F	F
	H	T	T	T
s_3 :	B	B	B	B
	T	H	H	H

$N_{jk} =$ COUNT OF TRANSITIONS FROM Z_j TO Z_k .

HMM LEARNING EXAMPLE (ESTIMATE π, A, ϕ)

s_1	Bias	Fair	fair	Bias
	H	T	H	H
s_2 :	F	F	F	F
	H	T	T	T
s_3 :	B	B	B	B
	T	H	H	H

N_{jk} = COUNT OF TRANSITIONS FROM Z_j TO Z_k .

$$N_{BB} = 3$$

$$N_{FF} = 4$$

$$N_{BF} = 1$$

$$N_{FB} = 1$$

HMM LEARNING EXAMPLE (ESTIMATE π, A, Φ)

	Bias	Fair	fair	Bias
s_1	H	T	H	H
s_2	F	F	F	F
	H	T	T	T
s_3	B	B	B	B
	T	H	H	H

$N_{BB} = 3$	$N_{FF} = 4$
$N_{BF} = 1$	$N_{FB} = 1$

$$\hat{A}_{BB} = \frac{N_{BB}}{N_{BB} + N_{BF}} = \frac{3}{4}$$

$$\hat{A}_{BF} = \frac{N_{BF}}{N_{BB} + N_{BF}} = \frac{1}{4}$$

$$\hat{A} = \begin{bmatrix} \hat{A}_{BB} & \hat{A}_{BF} \\ \hat{A}_{FB} & \hat{A}_{FF} \end{bmatrix}$$

$$\hat{A} = \begin{bmatrix} 0.75 & 0.25 \\ 0.2 & 0.8 \end{bmatrix}$$

HMM

LEARNING

EXAMPLE

(ESTIMATE π, A, Φ)

s_1	B	F	F	B
	H	T	H	H
s_2 :	F	F	F	F
	H	T	T	T
s_3 :	B	B	B	B
	T	H	H	H

$$\hat{\Phi} = \begin{bmatrix} \hat{\Phi}_{BH} & \hat{\Phi}_{BT} \\ \hat{\Phi}_{FH} & \hat{\Phi}_{FT} \end{bmatrix}$$

HMM LEARNING EXAMPLE (ESTIMATE π, A, Φ)

s_1	B	F	F	B
	H	T	H	H
s_2 :	F	F	F	F
	H	T	T	T
s_3 :	B	B	B	B
	T	H	H	H

$$\hat{\Phi} = \begin{bmatrix} \hat{\Phi}_{BH} & \hat{\Phi}_{BT} \\ \hat{\Phi}_{FH} & \hat{\Phi}_{FT} \end{bmatrix}$$

$$\hat{\Phi}_{BH} = \frac{\text{COUNT OF (Bias \& H)}}{\text{COUNT OF BIAS}}$$

HMM LEARNING EXAMPLE (ESTIMATE π, A, Φ)

s_1	B H	F T	F H	B H
s_2 :	F H	F T	F T	F T
s_3 :	B T	B H	B H	B H

$$\hat{\Phi} = \begin{bmatrix} \hat{\Phi}_{BH} & \hat{\Phi}_{BT} \\ \hat{\Phi}_{FH} & \hat{\Phi}_{FT} \end{bmatrix}$$

$$\begin{aligned} \hat{\Phi}_{BH} &= \frac{\text{COUNT OF (BIAS \& H)}}{\text{COUNT OF BIAS}} \\ &= \frac{\text{COUNT (B,H)}}{\text{COUNT (B,T) + COUNT (B,H)}} \\ &= \frac{4}{1 + 4} \\ &= 0.8 \end{aligned}$$

$$\hat{\Phi}_{BT} = 0.2$$

$$\hat{\Phi}_{FH} = 0.4$$

$$\hat{\Phi}_{FT} = 0.6$$

PVII: HMM PARAMETER LEARNING WITHOUT FULLY OBSERVED DATA

* WITH MISSING DATA | LATENT VARIABLES (z_t)
COMPUTING MLE | MAP IS HARD

* COULD USE GRADIENT BASED METHODS

PVII: HMM PARAMETER LEARNING WITHOUT FULLY OBSERVED DATA

* WITH MISSING DATA | LATENT VARIABLES (z_t)
COMPUTING MLE | MAP IS HARD

* COULD USE GRADIENT BASED METHODS
BUT TRICKY TO ENFORCE CONSTRAINTS

Like $\sum_j \pi_j = 1$ etc

* OFTEN USE EXPECTATION MAXIMISATION (EM)
→ ITERATIVE, CLOSED FORM UPDATES

EM ALGORITHM

Let x_i be observed variable

z_i be latent (hidden)

{ General case
DON'T CONFUSE
NOTATION WITH HMM }

GOAL: MAXIMIZE LOG. LIKELIHOOD

$$l(\theta) = \sum_{i=1}^N \log p(x_i | \theta)$$

EM ALGORITHM

Let x_i be observed variable

z_i be latent (hidden)

{ General case
DON'T CONFUSE
NOTATION WITH HMM }

* GOAL: MAXIMIZE LOG. LIKELIHOOD

$$l(\theta) = \sum_{i=1}^N \log p(x_i | \theta) = \sum_{i=1}^N \log \left[\sum_{z_i} p(x_i, z_i | \theta) \right]$$

HARD TO SOLVE! (LOG CAN'T BE
PUSHED INSIDE)

EM ALGORITHM

* GOAL: MAXIMIZE LOG. LIKELIHOOD

$$l(\theta) = \sum_{i=1}^N \log p(x_i | \theta) = \sum_{i=1}^N \log \left[\sum_{z_i} p(x_i, z_i | \theta) \right]$$

* DEFINE COMPLETE DATA LOG. LIKELIHOOD

$$l_c(\theta) \triangleq \sum_{i=1}^N \log p(x_i, z_i | \theta)$$

EM ALGORITHM

* GOAL: MAXIMIZE LOG. LIKELIHOOD

$$l(\theta) = \sum_{i=1}^N \log p(x_i | \theta) = \sum_{i=1}^N \log \left[\sum_{z_i} p(x_i, z_i | \theta) \right]$$

* DEFINE COMPLETE DATA LOG. LIKELIHOOD

$$l_c(\theta) \triangleq \sum_{i=1}^N \log p(x_i, z_i | \theta)$$

$l_c(\theta)$ can not be computed $\because z_i$ is unknown

EM ALGORITHM

* GOAL: MAXIMIZE LOG. LIKELIHOOD

$$l(\theta) = \sum_{i=1}^N \log p(x_i | \theta) = \sum_{i=1}^N \log \left[\sum_{z_i} p(x_i, z_i | \theta) \right]$$

* DEFINE COMPLETE DATA LOG. LIKELIHOOD

$$l_c(\theta) \triangleq \sum_{i=1}^N \log p(x_i, z_i | \theta)$$

* DEFINE EXPECTED COMPLETE DATA LIKELIHOOD

$$Q(\theta, \theta^{t-1}) = E[l_c(\theta) | D, \theta^{t-1}]$$

t : CURRENT ITERATION

Q : AUXILIARY FUNCTION

EM ALGORITHM

* GOAL: MAXIMIZE LOG. LIKELIHOOD

$$l(\theta) = \sum_{i=1}^N \log p(x_i | \theta) = \sum_{i=1}^N \log \left[\sum_{z_i} p(x_i, z_i | \theta) \right]$$

* DEFINE COMPLETE DATA LOG. LIKELIHOOD

$$l_c(\theta) \triangleq \sum_{i=1}^N \log p(x_i, z_i | \theta)$$

* DEFINE EXPECTED COMPLETE DATA LIKELIHOOD

$$Q(\theta, \theta^{t-1}) = E[l_c(\theta) | D, \theta^{t-1}]$$

* E-step: COMPUTE $Q \Rightarrow$ COMPUTE EXPECTED SUFFICIENT STATS (ESS)

EM ALGORITHM

* GOAL: MAXIMIZE LOG. LIKELIHOOD

$$l(\theta) = \sum_{i=1}^N \log p(x_i | \theta) = \sum_{i=1}^N \log \left[\sum_{z_i} p(x_i, z_i | \theta) \right]$$

* DEFINE COMPLETE DATA LOG. LIKELIHOOD

$$l_c(\theta) \triangleq \sum_{i=1}^N \log p(x_i, z_i | \theta)$$

* DEFINE EXPECTED COMPLETE DATA LIKELIHOOD

$$Q(\theta, \theta^{t-1}) = E[l_c(\theta) | D, \theta^{t-1}]$$

* E-step: COMPUTE Q

* M-step: OPTIMIZE Q wrt θ , i.e. $\theta^t = \underset{\theta}{\operatorname{argmax}} (\theta, \theta^{t-1})$

EM ALGORITHM (WITHOUT THE DERIVATION)

TWO STEPS (TILL CONVERGENCE)

① E : INFERRING MISSING VALUES ($Z_{i,t}$)
GIVEN MODEL PARAMETERS ($\theta = (\pi, A, \phi)$)

EM ALGORITHM

TWO STEPS (TILL CONVERGENCE)

① E : INFERRING MISSING VALUES ($Z_{i,t}$)
GIVEN MODEL PARAMETERS ($\theta = \{\pi, A, \phi\}$)

② M : OPTIMIZING PARAMETERS ($\theta = \{\pi, A, \phi\}$)
USING MLG AND
FILLED IN DATA ($Z_{i,t}$)

COMPUTING EXPECTED SUFFICIENT STATS (ESS)
(OR TERMS FROM AUX. FUNCTION ON
WHICH MLE DEPENDS ON)

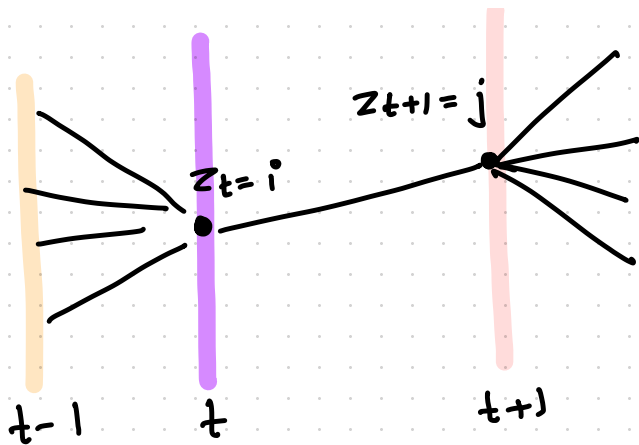
COMPUTING EXPECTED SUFFICIENT STATS (ESS)
(OR TERMS FROM AUX. FUNCTION ON WHICH MLE DEPENDS ON)

$$q_t(i, j) = P(z_t = i, z_{t+1} = j | x_{1:T}, \theta)$$

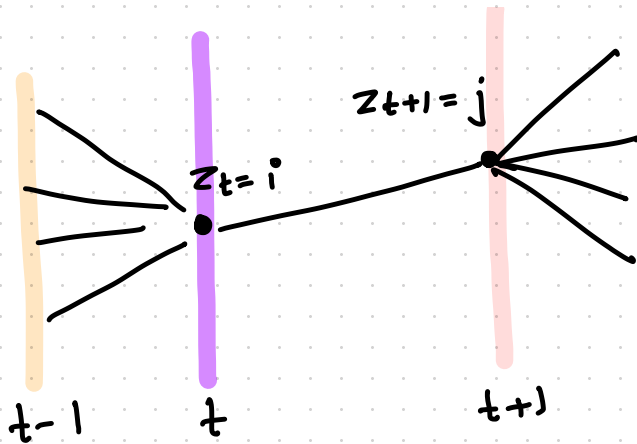
Probability of being in state 'i' at 't' and state 'j' at 't+1' given model and observations

COMPUTING EXPECTED SUFFICIENT STATS (ESS)
(OR TERMS FROM AUX. FUNCTION ON WHICH MLE DEPENDS ON)

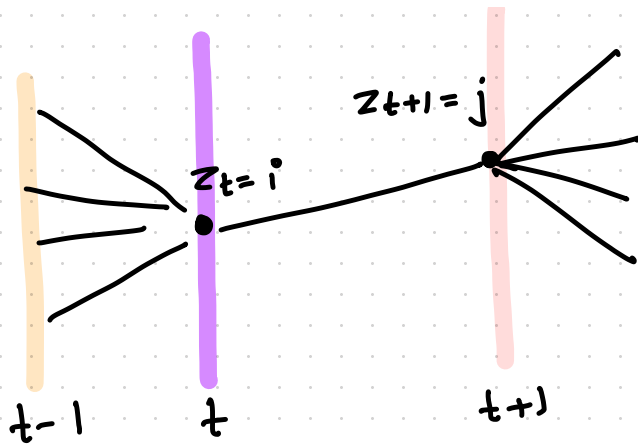
$$q_t(i, j) = P(z_t = i, z_{t+1} = j | x_{1:T}, \theta)$$



$$q_t(i, j) = P(z_t = i, z_{t+1} = j | x_{1:T}, \theta)$$
$$= \frac{P(z_t = i, z_{t+1} = j, x_{1:T} | \theta)}{P(x_{1:T} | \theta)}$$

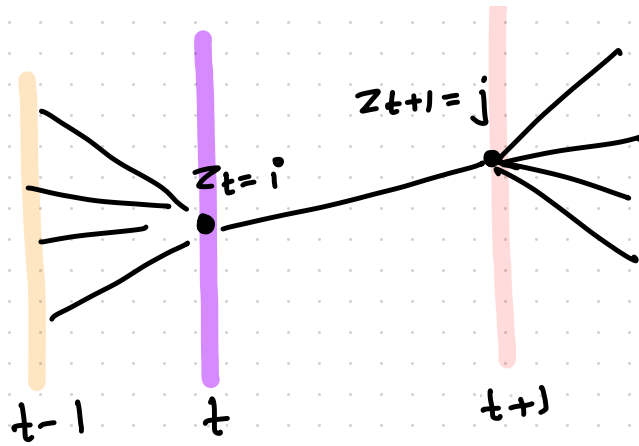


$$\begin{aligned}
 q_t(i, j) &= P(z_t = i, z_{t+1} = j \mid x_{1:T}, \theta) \\
 &= \frac{P(z_t = i, z_{t+1} = j, x_{1:T} \mid \theta)}{P(x_{1:T} \mid \theta)}
 \end{aligned}$$



$$q_t(i, j) = \frac{P(z_t = i, x_{1:t}) A_{ij} P(x_{t+1} \mid z_{t+1} = j) P(x_{t+1:T} \mid z_{t+1} = j)}{P(x_{1:T} \mid \theta)}$$

$$q_t(i, j) = P(z_t = i, z_{t+1} = j | x_{1:T}, \theta)$$



$$\begin{aligned}
 q_t(i, j) &= \frac{P(z_t = i, x_{1:t}) A_{ij} P(x_{t+1} | z_{t+1} = j) P(x_{t+1:T} | z_{t+1} = j)}{P(x_{1:T} | \theta)} \\
 &= \frac{\alpha_t(i) \phi_j(x_{t+1}) \beta_{t+1}(j) A_{ij}}{\sum_i \sum_j \alpha_t(i) \phi_j(x_{t+1}) \beta_{t+1}(j) A_{ij}}
 \end{aligned}$$

$$\ell_{y_t}(i, j) = P(z_t = i, z_{t+1} = j | x_{1:T}, \theta)$$

now; $\gamma_t(i) =$ Probability of being in state 'i' at time 't' given observations and model.

$$\eta_t(i, j) = P(z_t = i, z_{t+1} = j \mid x_{1:T}, \theta)$$

now; $\gamma_t(i) =$ Probability of being in state 'i' at time 't' given observations and model.

$$\Rightarrow \gamma_t(i) = \sum_{j=1}^K \eta_t(i, j)$$

$$q_t(i, j) = P(z_t = i, z_{t+1} = j | x_{1:T}, \theta)$$

Now; $\gamma_t(i) =$ Probability of being in state 'i' at time 't' given observations and model.

$$\Rightarrow \gamma_t(i) = \sum_{j=1}^K q_t(i, j)$$

Q: What does $\sum_t \gamma_t(i)$ mean?

$$q_t(i, j) = P(z_t = i, z_{t+1} = j | x_{1:T}, \theta)$$

Now; $\gamma_t(i) =$ Probability of being in state 'i' at time 't' given observations and model.

$$\Rightarrow \gamma_t(i) = \sum_{j=1}^K q_t(i, j)$$

Q: What does $\sum_t \gamma_t(i)$ mean?

Expected # of visits to state i

$$q_t(i, j) = P(z_t = i, z_{t+1} = j | x_{1:T}, \theta)$$

Now; $\gamma_t(i) =$ Probability of being in state 'i' at time 't' given observations and model.

$$\Rightarrow \gamma_t(i) = \sum_{j=1}^K q_t(i, j)$$

Q: What does $\sum_{t=1}^{T-1} \gamma_t(i)$ mean

Expected # of visits to state i

Or, Expected # of transitions from state i (ignoring $t=T$)

$$l_{y_t}(i, j) = P(z_t = i, z_{t+1} = j | x_{1:T}, \theta)$$

Q: What does $\sum_{t=1}^{T-1} \gamma_t(i)$ mean

Expected # of visits to state i

Or, Expected # of transitions from state i (ignoring $z=T$)

Q: What does $\sum_{t=1}^{T-1} l_{y_t}(i, j)$ mean?

$$l_{y_t}(i, j) = P(z_t = i, z_{t+1} = j | x_{1:T}, \theta)$$

Q: What does $\sum_{t=1}^{T-1} \gamma_t(i)$ mean

Expected # of visits to state i

Or, Expected # of transitions from state i (ignoring $z=T$)

Q: What does $\sum_{t=1}^{T-1} l_{y_t}(i, j)$ mean?

Expected # of transitions from state i to j .

HMM LEARNING EXAMPLE (ESTIMATE π, A, Φ)

s_1	Bias	Fair	fair	Bias	Unknown
	H	T	H	H	
s_2 :	F	F	F	F	
	H	T	T	T	
s_3 :	B	B	B	B	
	T	H	H	H	

Let us assume :

$$\pi^1 = \begin{bmatrix} 0.9 & 0.1 \end{bmatrix}$$

$$\Phi_B^1 = \begin{bmatrix} H & T \\ 0.8 & 0.2 \end{bmatrix}$$

$$\Phi_F^1 = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$$

$$A^1 = \begin{bmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{bmatrix}$$

← 1st iteration (RANDOM INIT.)

HMM LEARNING EXAMPLE (ESTIMATE π, A, Φ)

CONSIDER S1 ONLY FOR NOW

S1

H T H H

$$\pi' = [0.9 \quad 0.1]$$

$$A' = \begin{bmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{bmatrix}$$

$$\Phi'_B = \begin{bmatrix} H & T \\ 0.8 & 0.2 \end{bmatrix}$$

$$\Phi'_F = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$$

$$u_1(B, B) = ?$$

$$u_1(B, F) = ?$$

$$u_1(F, F) = ? \quad u_1(F, B) = ?$$

HMM LEARNING EXAMPLE (ESTIMATE π, A, Φ)

CONSIDER S1 ONLY FOR NOW

S1	H	T	H	H
----	---	---	---	---

$$\pi^1 = [0.9 \quad 0.1] \quad A^1 = \begin{bmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{bmatrix}$$

$$\phi_B^1 = \begin{matrix} H & T \\ [0.8 & 0.2] \end{matrix} \quad \phi_F^1 = \begin{matrix} H & T \\ [0.4 & 0.6] \end{matrix}$$

$$e_{y_1}(B, B) = \frac{\alpha_1(B) \cdot A_{BB} \phi_B(x_2=T) \cdot \beta_2(B)}{\text{NORMALIZATION CONSTANT}}$$

HMM LEARNING EXAMPLE (ESTIMATE π, A, ϕ)

CONSIDER S1 ONLY FOR NOW

S1	H	T	H	H
----	---	---	---	---

$$\pi' = \begin{bmatrix} 0.9 & 0.1 \end{bmatrix} \quad A' = \begin{bmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{bmatrix}$$

$$\phi_B' = \begin{bmatrix} H & T \\ 0.8 & 0.2 \end{bmatrix} \quad \phi_F' = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$$

$$e_{y_1}(B, B) = \frac{\alpha_1(B) \cdot A_{BB} \phi_B(x_2=T) \cdot \beta_2(B)}{\text{NORMALIZATION CONSTANT}}$$

$$\alpha_1(B) = \pi_B \phi_B(x_1=H) = (0.9)(0.8) = 0.72$$

$$\beta_2(B) = ?$$

HMM LEARNING EXAMPLE (ESTIMATE π, A, ϕ)

CONSIDER S1 ONLY FOR NOW

S1	H	T	H	H
----	---	---	---	---

$$\pi' = \begin{bmatrix} 0.9 & 0.1 \end{bmatrix} \quad A' = \begin{bmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{bmatrix}$$

$$\phi'_B = \begin{bmatrix} H & T \\ 0.8 & 0.2 \end{bmatrix} \quad \phi'_F = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$$

$$\beta_4(B) = 1 \quad \& \quad \beta_4(F) = 1$$

$$\beta_3(B) = A_{BB} \phi_B(x_4=H) \beta_4(B) + A_{BF} \phi_F(x_4=H) \beta_4(F)$$

$$= (0.7)(0.8) + (0.3)(0.4) = .56 + .12 = 0.68$$

HMM LEARNING EXAMPLE (ESTIMATE π, A, ϕ)

CONSIDER S1 ONLY FOR NOW

S1	H	T	H	H
----	---	---	---	---

$$\pi' = \begin{bmatrix} 0.9 & 0.1 \end{bmatrix} \quad A' = \begin{bmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{bmatrix}$$

$$\phi'_B = \begin{bmatrix} H & T \\ 0.8 & 0.2 \end{bmatrix} \quad \phi'_F = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$$

$$\beta_4(B) = 1 \quad \& \quad \beta_4(F) = 1$$

$$\begin{aligned} \beta_3(B) &= A_{BB} \phi_B(x_4=H) \beta_4(B) + A_{BF} \phi_F(x_4=H) \beta_4(F) \\ &= (0.7)(0.8) + (0.3)(0.4) = .56 + .12 = 0.68 \end{aligned}$$

$$\begin{aligned} \beta_2(B) &= A_{BB} \phi_B(x_3=H) \beta_3(B) + A_{BF} \phi_F(x_3=H) \beta_3(F) \\ &= (0.7)(0.8)(0.68) + (0.3)(0.4)(1 - .68) = 0.4192 \end{aligned}$$

HMM LEARNING EXAMPLE (ESTIMATE π, A, Φ)

CONSIDER S1 ONLY FOR NOW

S1	H	T	H	H
----	---	---	---	---

$$\pi^1 = \begin{bmatrix} 0.9 & 0.1 \end{bmatrix} \quad A^1 = \begin{bmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{bmatrix}$$

$$\phi_B^1 = \begin{bmatrix} H & T \\ 0.8 & 0.2 \end{bmatrix} \quad \phi_F^1 = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$$

$$e_{y_1}(B, B) = \frac{\alpha_1(B) \cdot A_{BB} \phi_B(x_2=T) \cdot \beta_2(B)}{\text{NORMALIZATION CONSTANT (C)}}$$

$$= \frac{0.72 * 0.7 * 0.2 * 0.4192}{C}$$

HMM LEARNING EXAMPLE (ESTIMATE π, A, Φ)

CONSIDER S1 ONLY FOR NOW

S1	H	T	H	H
----	---	---	---	---

$$\pi^1 = [0.9 \quad 0.1] \quad A^1 = \begin{bmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{bmatrix}$$

$$\phi_B^1 = \begin{bmatrix} H & T \\ 0.8 & 0.2 \end{bmatrix} \quad \phi_F^1 = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$$

we can similarly find $\mu_1(B, F)$, $\mu_1(F, F)$, $\mu_1(F, B)$ for this sequence

HMM LEARNING EXAMPLE (ESTIMATE π, A, Φ)

CONSIDER S1 ONLY FOR NOW

S1	H	T	H	H
----	---	---	---	---

$$\pi^1 = \begin{bmatrix} 0.9 & 0.1 \end{bmatrix} \quad A^1 = \begin{bmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{bmatrix}$$

$$\phi_B^1 = \begin{bmatrix} H & T \\ 0.8 & 0.2 \end{bmatrix} \quad \phi_F^1 = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$$

we can find.

$$\gamma_1(B) = \mu_1(B, F) + \mu_1(B, B)$$

$$\gamma_1(F) = 1 - \gamma_1(B)$$

HMM LEARNING EXAMPLE (ESTIMATE π, A, Φ)

s_1

H

T

H

H

CONSIDER ONLY THE STATES AT $t=1$

$$P(z_1 = B) = \gamma_1(B)$$

$$P(z_1 = F) = \gamma_1(F)$$

HMM LEARNING EXAMPLE (ESTIMATE π, A, Φ)

s_1

H

T

H

H

CONSIDER ONLY THE STATES AT $t=1$

$$P(z_1 = B) = \gamma_1(B)$$

$$P(z_1 = F) = \gamma_1(F)$$

We don't know for sure whether z_1 is B or F,
but we know the probabilities.

KEY INSIGHT: CONSIDER $\gamma_1(j) = \text{EXPECTED \# OF TIMES}$
 $z_1 = j$

HMM LEARNING EXAMPLE (ESTIMATE π, A, ϕ)

KEY INSIGHT: CONSIDER $\gamma_i(j) =$ EXPECTED # OF TIMES
 $z_i = j$

NOW AS WE DID FOR FULLY OBSERVED CASE,
WE CAN USE MLE TO
estimate

$$\pi^2, A^2, \phi^2$$

EM ALGORITHM

TWO STEPS (TILL CONVERGENCE)

① E : INFERRING MISSING VALUES ($Z_{i,t}$)
GIVEN MODEL PARAMETERS ($\theta = \{\pi, A, \Phi\}$)

COMPUTE ESS

$$\gamma_{i,t}(j) \stackrel{\Delta}{=} P(Z_t = j \mid x_{i,1:T_i}, \theta)$$

$$\eta_{i,t}(j,k) \stackrel{\Delta}{=} P(Z_{t-1} = j, Z_t = k \mid x_{i,1:T_i}, \theta)$$

EM ALGORITHM

TWO STEPS (TILL CONVERGENCE)

② M : OPTIMIZING PARAMETERS ($\theta = \{\pi, A, D\}$)
USING MLB AND
FILLED IN DATA ($Z_{i,t}$) OR ESS

$$\hat{\pi}_k = \text{Expected Freq. of } Z_{i,t} = \text{state } k \\ = \frac{\sum_{i=1}^N \gamma_{i,1}(k)}{N}$$

EM ALGORITHM

TWO STEPS (TILL CONVERGENCE)

② M: OPTIMIZING PARAMETERS ($\theta = \{\pi, A, \phi\}$)

USING MLB AND

FILLED IN DATA (Z_i, t) OR ESS

$$\hat{A}_{jR} = \frac{\text{Expected \# transitions } z_t = j \text{ to } z_{t+1} = R}{\text{Expected \# transitions from } z_t = j}$$

EM ALGORITHM

TWO STEPS (TILL CONVERGENCE)

(2) M: OPTIMIZING PARAMETERS ($\theta = \{\pi, A, D\}$)

USING MLG AND

FILLED IN DATA ($Z_{i,t}$) OR ESS

$\hat{A}_{j,k} =$ Expected # transitions $Z_t = j$ to $Z_{t+1} = k$

Expected # transitions from $Z_t = j$

$$= \frac{\sum_{i=1}^N \sum_{t=2}^{T_i} y_{i,t}(j,k)}{\sum_{k=1}^K \sum_{i=1}^N \sum_{t=2}^{T_i} y_{i,t}(j,k)} \leftarrow \text{NORMALIZATION}$$

EM ALGORITHM

TWO STEPS (TILL CONVERGENCE)

(2) M: OPTIMIZING PARAMETERS ($\theta = \{\pi, A, \phi\}$)

USING MLG AND

FILLED IN DATA ($Z_{i,t}$) OR ESS

$\hat{\phi}_{j,l} = \frac{\text{Expected \# times in state } j' \text{ and observing } l'}{\text{Expected \# times in state } j}$

$$= \frac{\sum_{i=1}^N \sum_{t: x_{i,t}=l} \gamma_{i,t}(j)}{\sum_{i=1}^N \sum_{t=1}^T \gamma_{i,t}(j)}$$

HOW TO INITIALIZE θ_s ?

HOW TO INITIALIZE θ_s ?

- ① RANDOMLY INITIALIZE W/ MULTIPLE RESTARTS
- CHOOSE THE "BEST" (MAX LL) PARAMS

HOW TO INITIALIZE θ s?

- ① RANDOMLY INITIALIZE W/ MULTIPLE RESTARTS
- CHOOSE THE "BEST" (MAX LL) PARAMS
- ② INIT. USING "SOMF" TRAINING DATA
- KNOWN z s

HOW TO INITIALIZE θ s?

- ① RANDOMLY INITIALIZE W/ MULTIPLE RESTARTS
- CHOOSE THE "BEST" (MAX LL) PARAMS
- ② INIT. USING "SOMF" TRAINING DATA
- KNOWN Zs
- ③ INIT. USING VITERBI TRAINING