



香港科技大學  
THE HONG KONG  
UNIVERSITY OF SCIENCE  
AND TECHNOLOGY

COMP 5212  
Machine Learning  
Lecture 15

# Hidden Markov Models

Junxian He  
April 9, 2026

# Announcements

1. Mid-term exam grades are out, we will hold a paper-check session next week



# i.i.d to sequential data

# i.i.d to sequential data

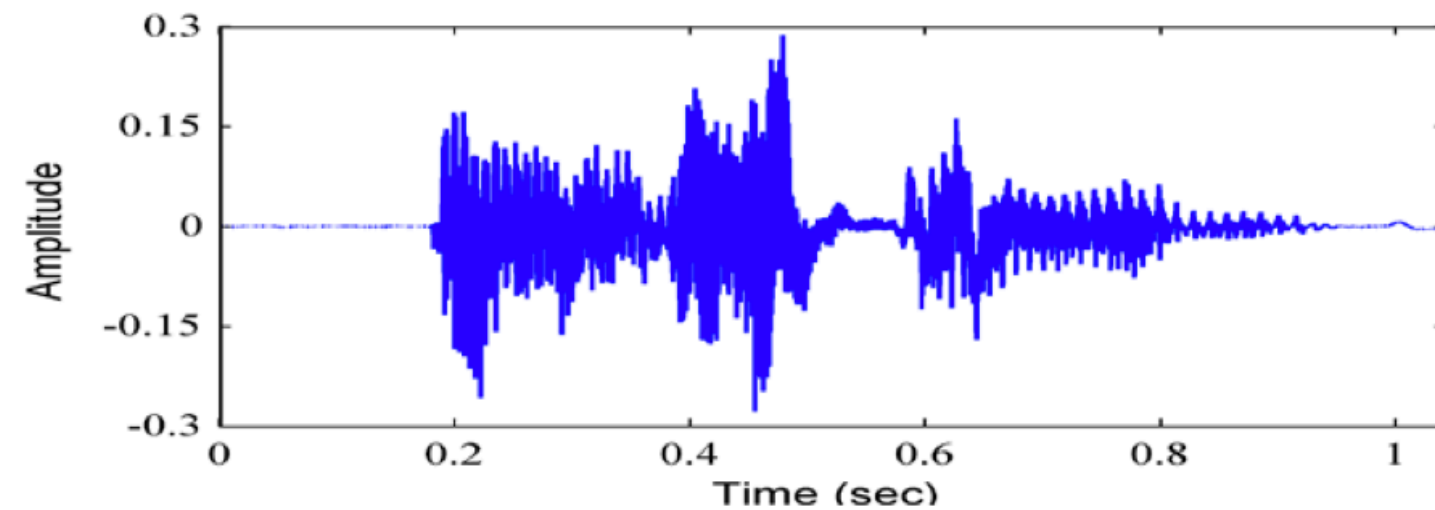
- So far we assumed independent, identically distributed data,  $\{X_i\}_{i=1}^n \stackrel{iid}{\sim} p(\mathbf{X})$

# i.i.d to sequential data

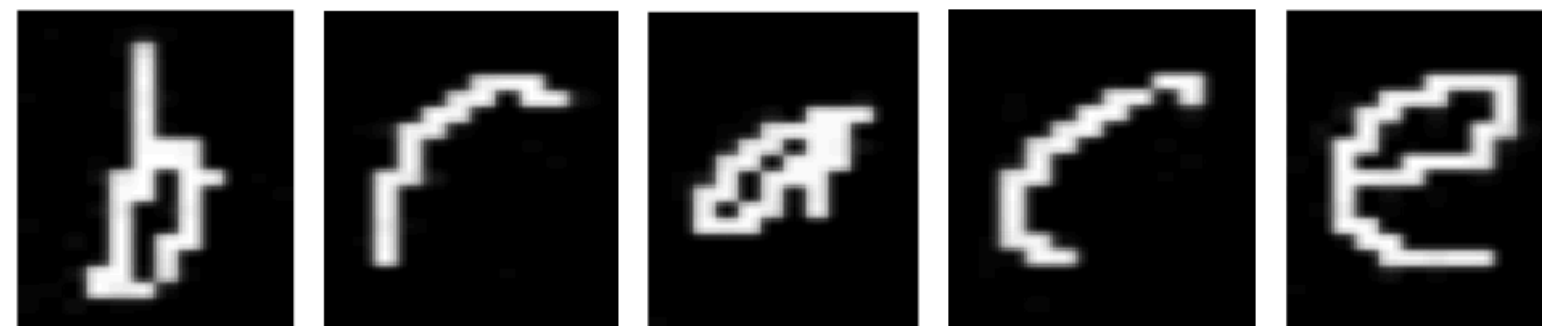
❑ So far we assumed independent, identically distributed data  $\{X_i\}_{i=1}^n \stackrel{iid}{\sim} p(\mathbf{X})$

❑ Sequential (non i.i.d.) data

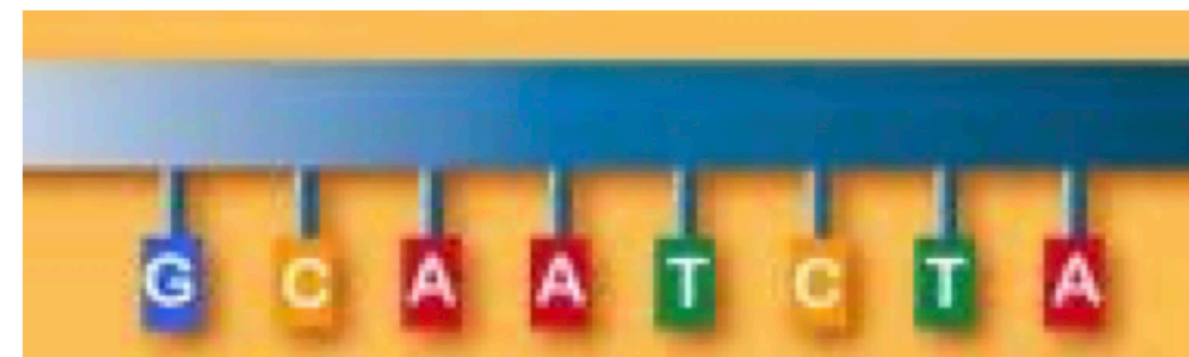
– Time-series data  
E.g. Speech



– Characters in a sentence



– Base pairs along a DNA strand



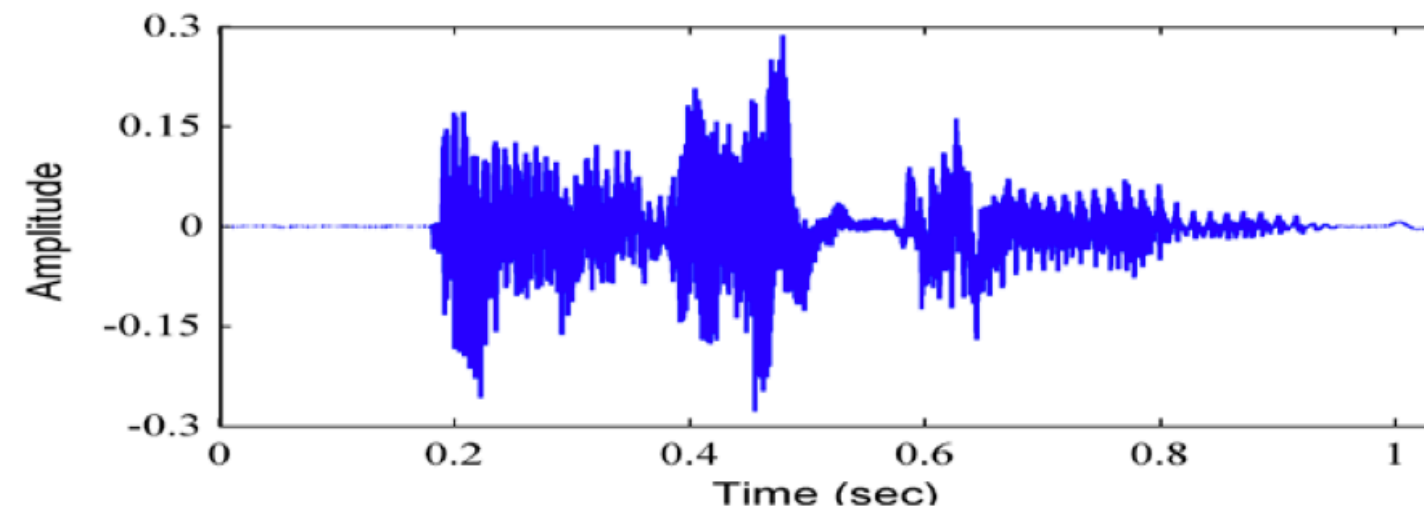
# i.i.d to sequential data

❑ So far we assumed independent, identically distributed data  $\{X_i\}_{i=1}^n \stackrel{iid}{\sim} p(\mathbf{X})$

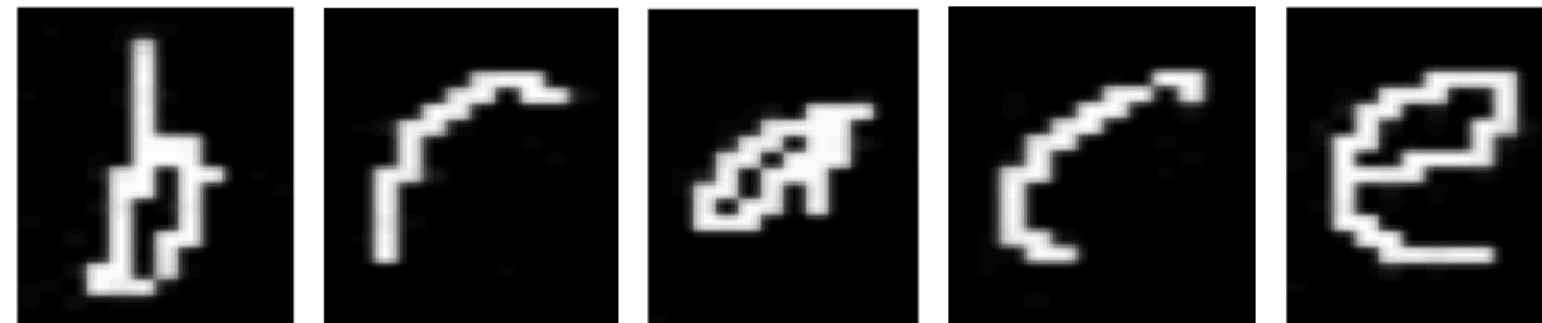
❑ Sequential (non i.i.d.) data

– Time-series data

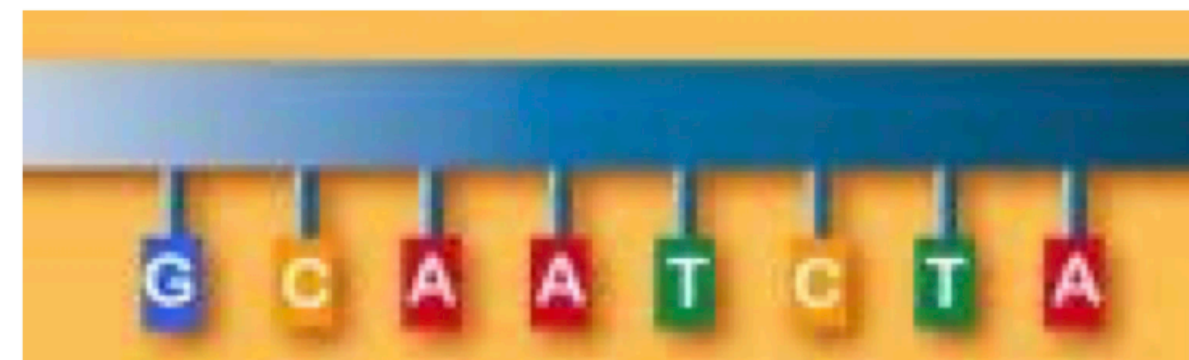
E.g. Speech



– Characters in a sentence



– Base pairs along a DNA strand



(Sequential data is still i.i.d on the sequence level)

# Review: Elimination Algorithm / Marginalization

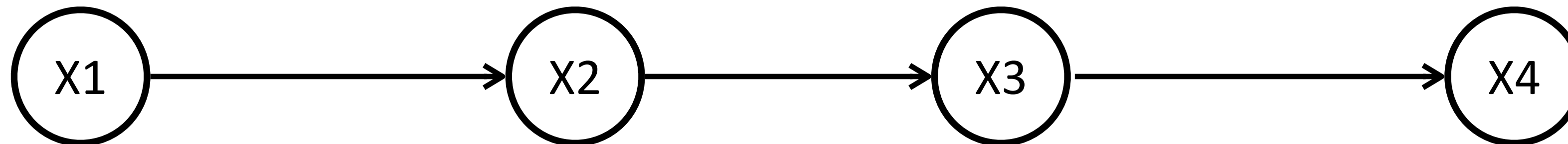
$$P(h) = \sum_g \sum_f \sum_e \sum_d \sum_c \sum_b \sum_a P(a, b, c, d, e, f, g, h)$$

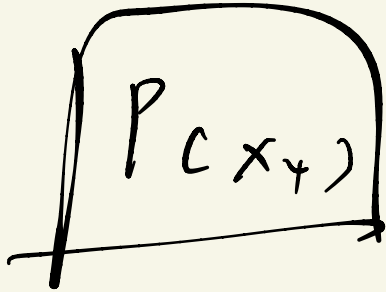
a naive summation needs to enumerate over an exponential number of terms

# Review: Elimination Algorithm / Marginalization

$$P(h) = \sum_g \sum_f \sum_e \sum_d \sum_c \sum_b \sum_a P(a, b, c, d, e, f, g, h)$$

a naïve summation needs to enumerate over an exponential number of terms





$$\sum_{x_1} \sum_{x_2} \sum_{x_3} P(x_2 | x_1) P(x_1) P(x_3 | x_2) P(x_4 | x_3)$$

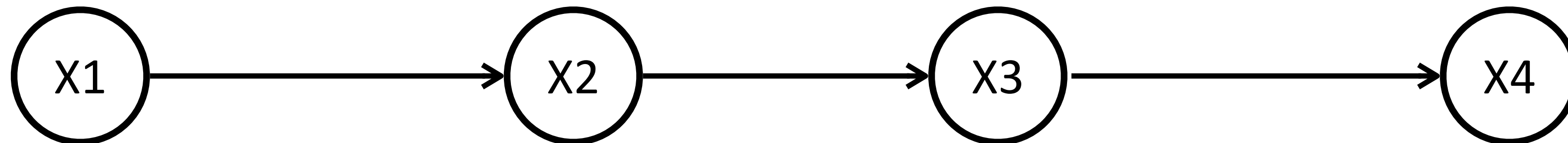
$$\sum_{x_2} \sum_{x_3} P(x_3 | x_2) P(x_4 | x_3) \left( \sum_{x_1} P(x_2 | x_1) P(x_1) \right) = P(x_2) P(x_3)$$

# Review: Elimination Algorithm / Marginalization

$$P(h) = \sum_g \sum_f \sum_e \sum_d \sum_c \sum_b \sum_a P(a, b, c, d, e, f, g, h)$$



a naïve summation needs to enumerate over an exponential number of terms



What if the random variables follow this chain structure?

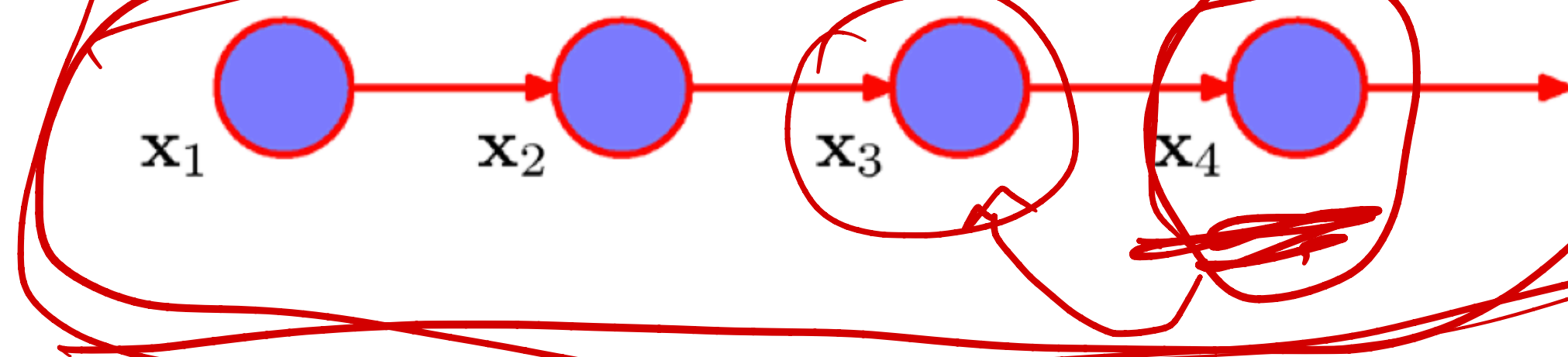
# Review: Markov Models

# Review: Markov Models

## □ Markov Assumption

1<sup>st</sup> order

$$p(\mathbf{X}) = \prod_{i=1}^n p(X_i | X_{i-1})$$



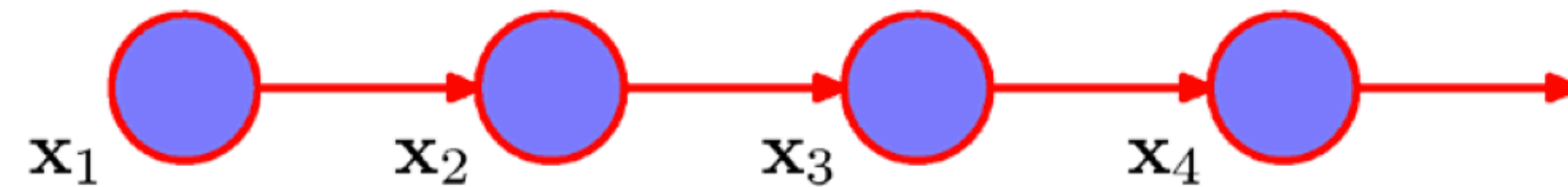
vanilla  
 $p(\mathbf{x}) = \prod_{i=1}^n p(x_i | x_{i-1})$

# Review: Markov Models

## □ Markov Assumption

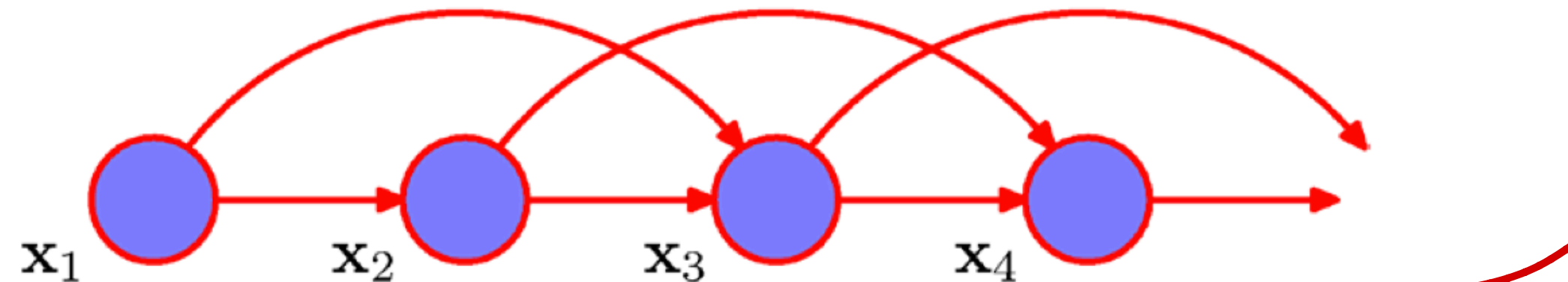
1<sup>st</sup> order

$$p(\mathbf{X}) = \prod_{i=1}^n p(X_i | X_{i-1})$$



2<sup>nd</sup> order

$$p(\mathbf{X}) = \prod_{i=1}^n p(X_i | X_{i-1}, X_{i-2})$$



# Review: Markov Models

# Review: Markov Models

Homogeneous/stationary Markov model (probabilities don't depend on  $n$ )

# Review: Markov Models

Homogeneous/stationary Markov model (probabilities don't depend on  $n$ )

## □ Markov Assumption

1<sup>st</sup> order

$$p(\mathbf{X}) = \prod_{i=1}^n p(X_i | X_{i-1})$$

# parameters in stationary model  
 $K$ -ary variables

$O(K^2)$

$$P(X_n | X_{n-1})$$

$\frac{K^2}{1}$

$$P(X_2 = 3 | X_1 = 0)$$

$$P(X_{10} = 3 | X_9 = 0)$$

# Review: Markov Models

Homogeneous/stationary Markov model (probabilities don't depend on  $n$ )

## □ Markov Assumption

# parameters in  
stationary model  
K-ary variables

1<sup>st</sup> order  $p(\mathbf{X}) = \prod_{i=1}^n p(X_i | X_{i-1})$   $O(K^2)$

m<sup>th</sup> order  $p(\mathbf{X}) = \prod_{i=1}^n p(X_i | X_{i-1}, \dots, X_{i-m})$   $O(K^{m+1})$

$O(K^{m+1})$

# Review: Markov Models

$$P(\mathbf{X}) = \prod_{i=1}^n P(X_i)$$

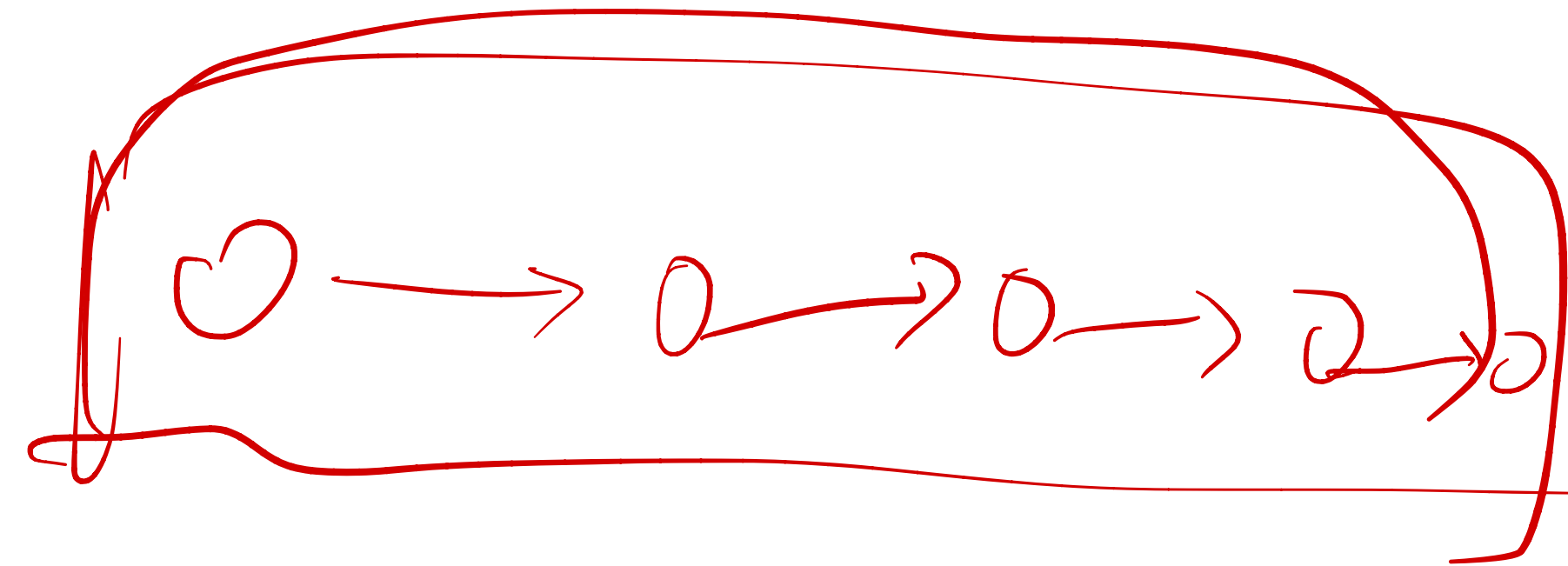
Homogeneous/stationary Markov model (probabilities don't depend on  $n$ )

## □ Markov Assumption

# parameters in stationary model  
K-ary variables

1<sup>st</sup> order  $p(\mathbf{X}) = \prod_{i=1}^n p(X_i | X_{i-1})$

$O(K^2)$



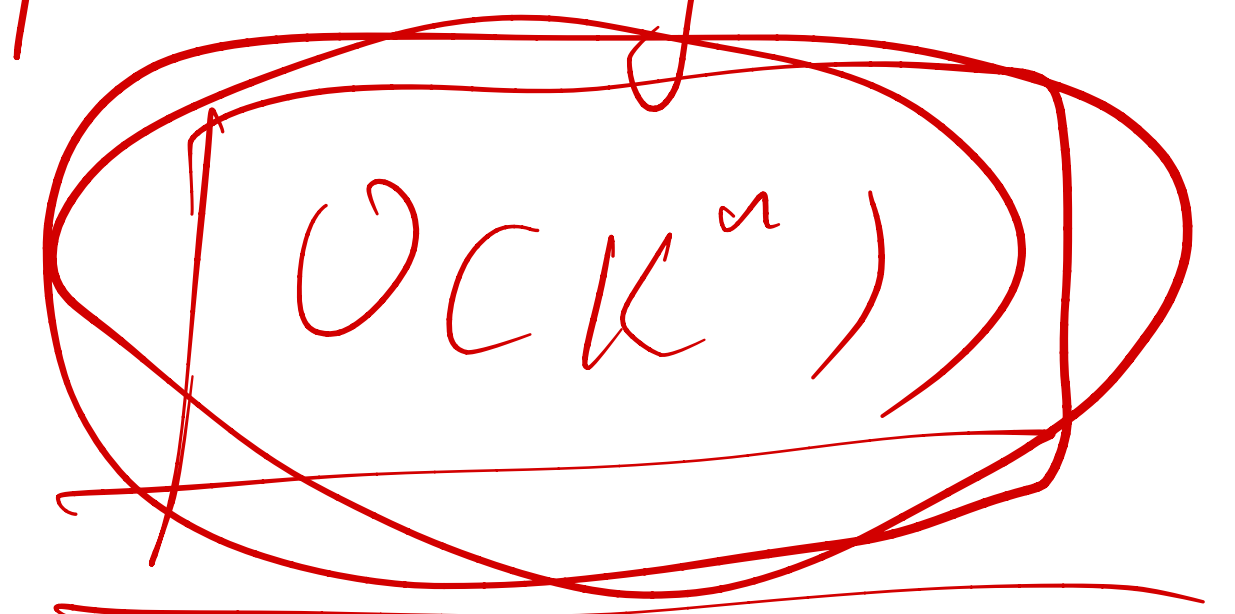
m<sup>th</sup> order  $p(\mathbf{X}) = \prod_{i=1}^n p(X_i | X_{i-1}, \dots, X_{i-m})$

$O(K^{m+1})$

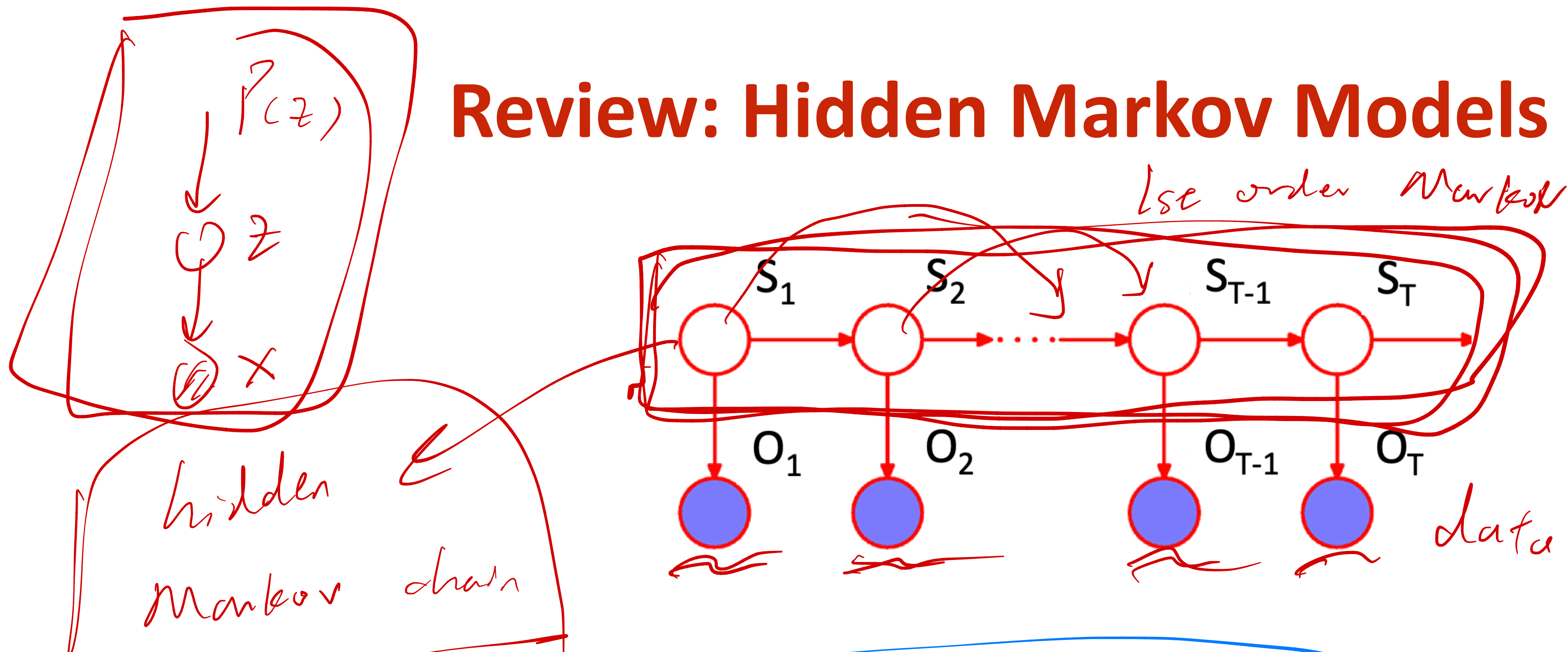
n-1<sup>th</sup> order  $p(\mathbf{X}) = \prod_{i=1}^n p(X_i | X_{i-1}, \dots, X_1)$   $O(K^n)$

≡ no assumptions – complete (but directed) graph

exponential growth



# Review: Hidden Markov Models



Observation space

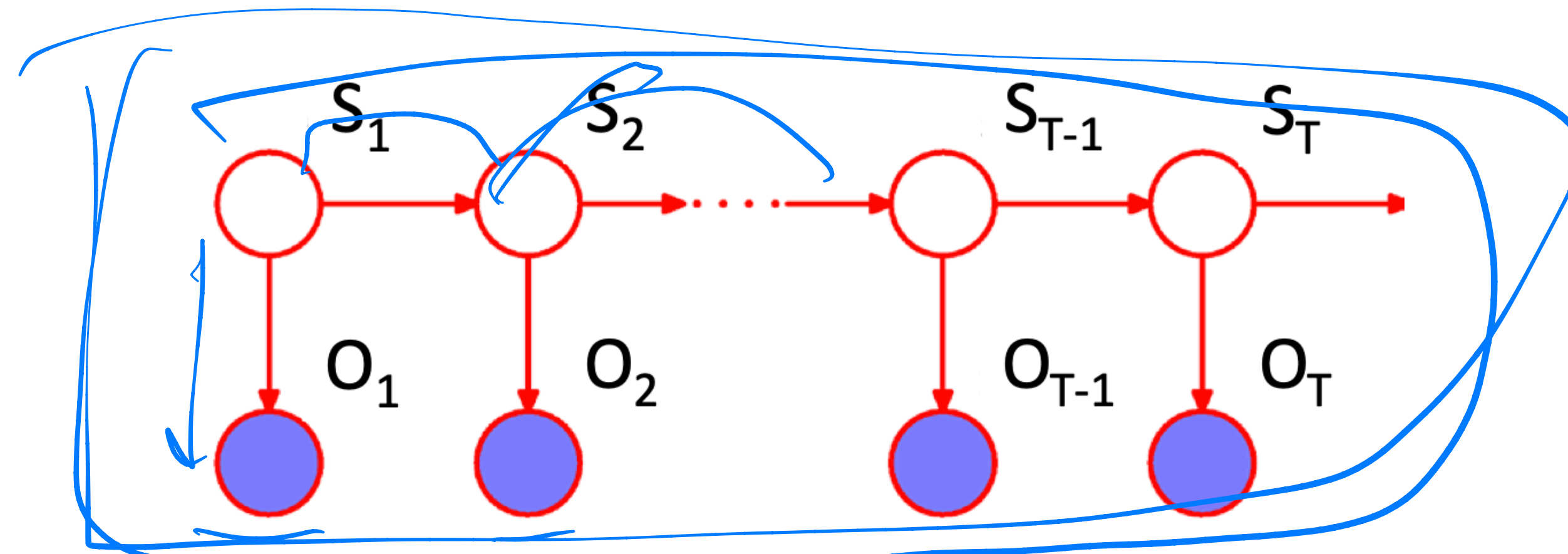
$$O_t \in \{y_1, y_2, \dots, y_K\}$$

Hidden states

$$S_t \in \{1, \dots, I\}$$

# Hidden Markov Models

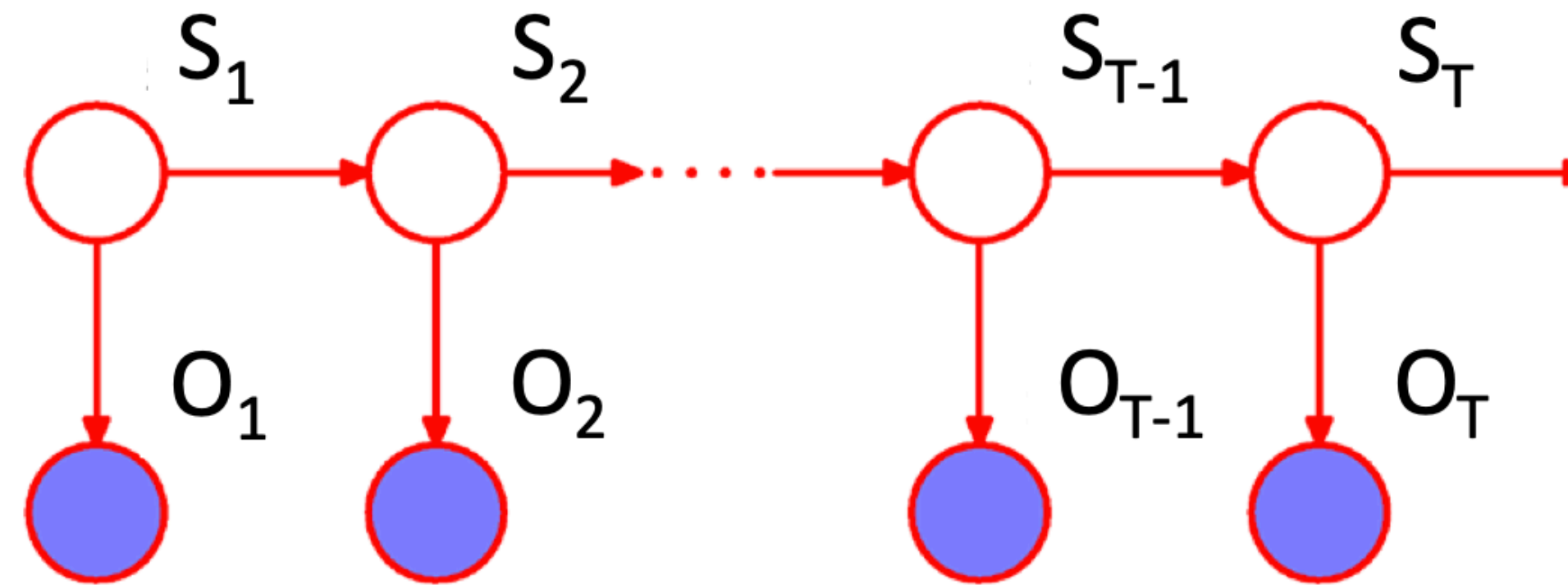
# Hidden Markov Models



$p(o_1, o_2, \dots, o_T)$

$$p(S_1, \dots, S_T, O_1, \dots, O_T) = \prod_{t=1}^T p(O_t | S_t) \prod_{t=1}^T p(S_t | S_{t-1})$$

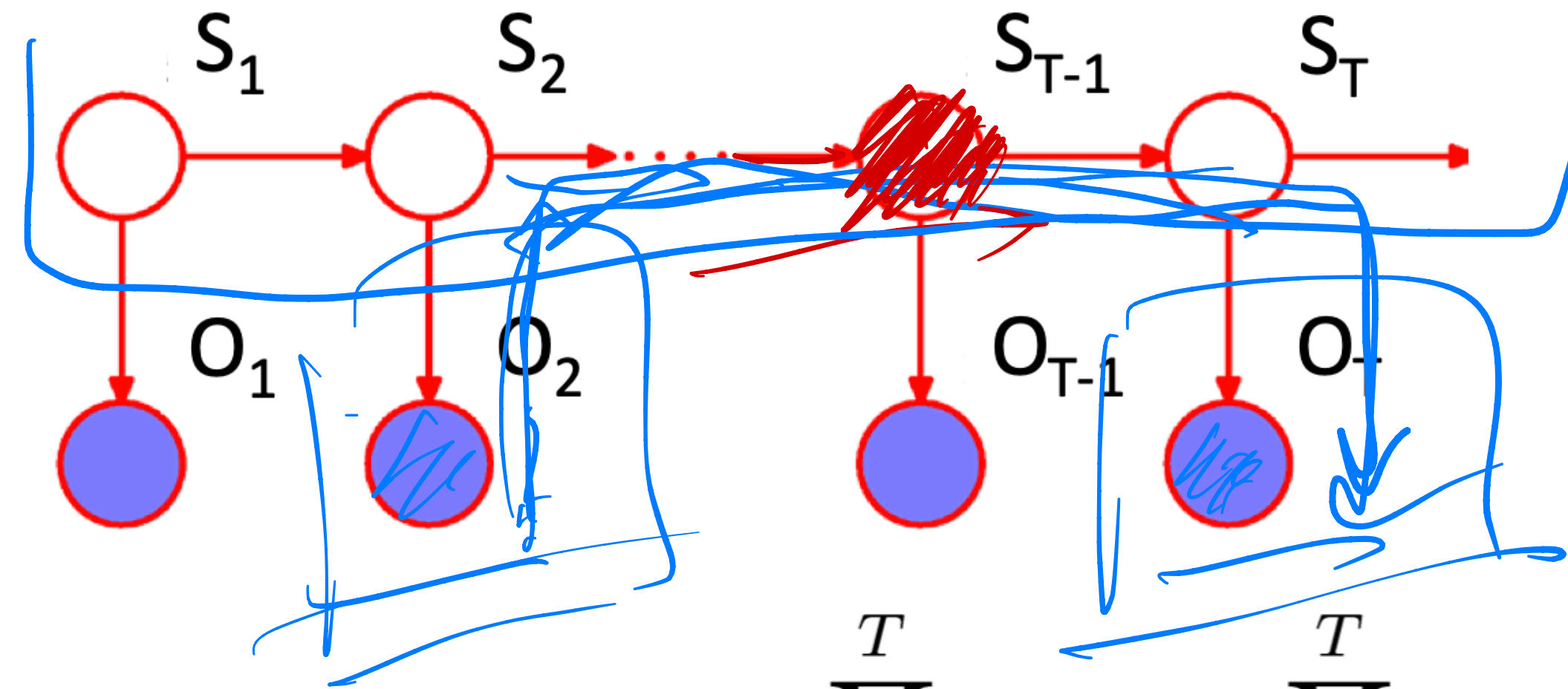
# Hidden Markov Models



$$p(S_1, \dots, S_T, O_1, \dots, O_T) = \prod_{t=1}^T p(O_t | S_t) \prod_{t=1}^T p(S_t | S_{t-1})$$

1<sup>st</sup> order Markov assumption on hidden states  $\{S_t\}$   $t = 1, \dots, T$   
(can be extended to higher order).

# Hidden Markov Models



$$p(S_1, \dots, S_T, O_1, \dots, O_T) = \prod_{t=1}^T p(O_t | S_t) \prod_{t=1}^T p(S_t | S_{t-1})$$

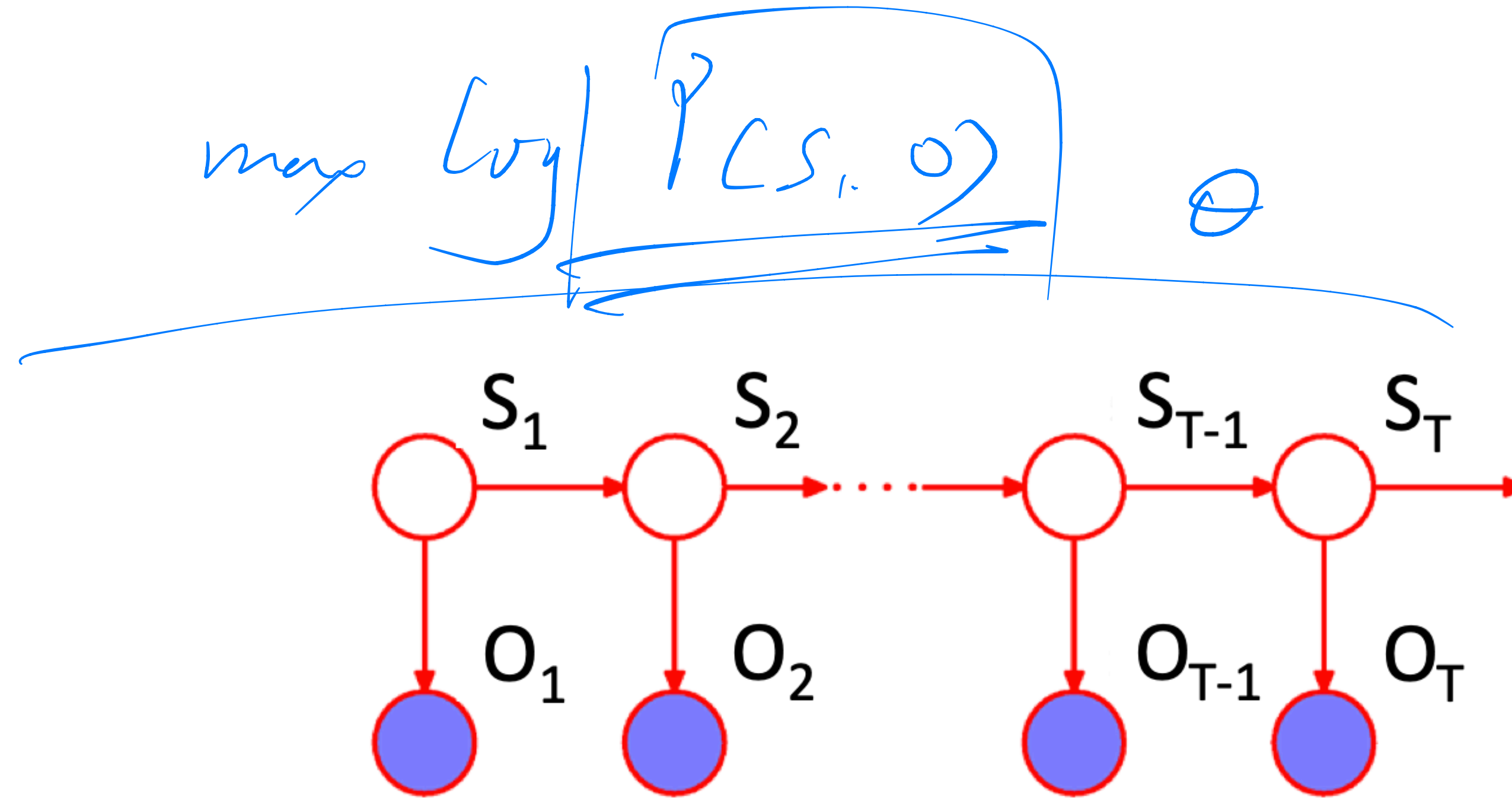
1<sup>st</sup> order Markov assumption on hidden states  $\{S_t\}$   $t = 1, \dots, T$   
 (can be extended to higher order).

Is  $O_T$  and  $O_2$  independent?

buys bull

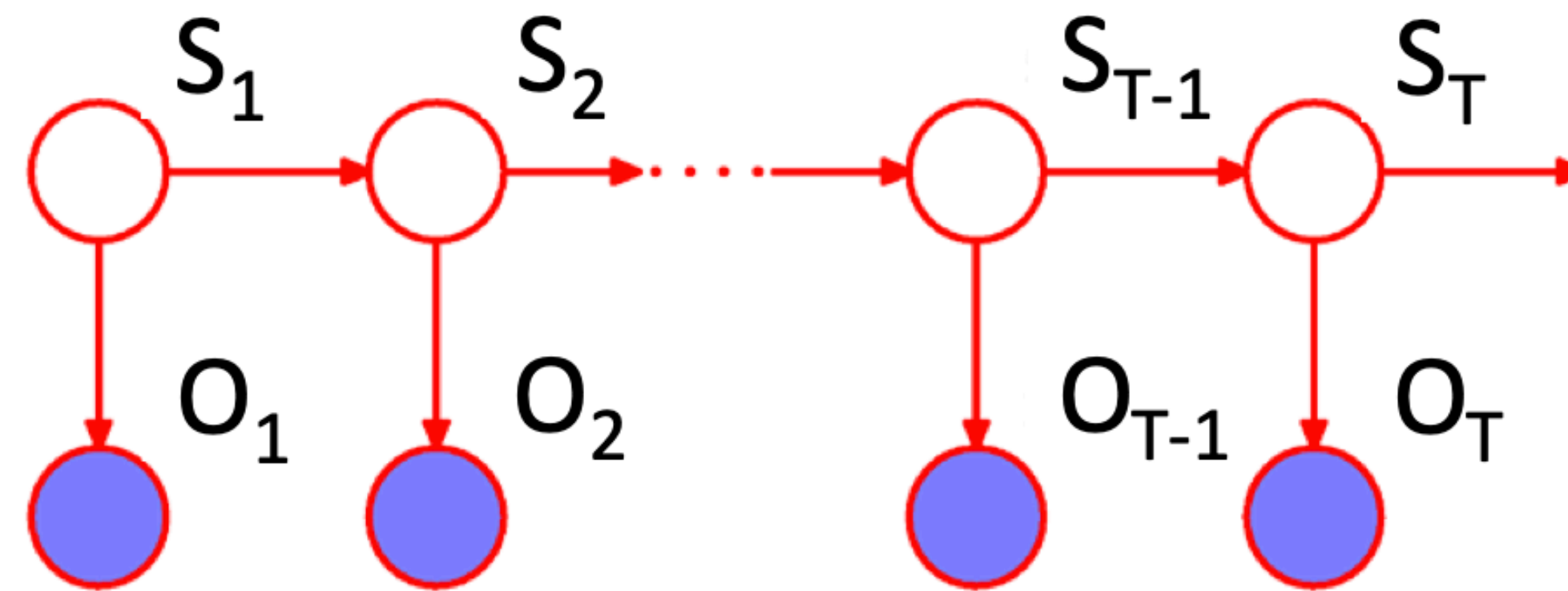
$O_T$  and  $O_2$   
 independent  
 conditioned on  $S_{T-1}$

# Hidden Markov Models



# Hidden Markov Models

- Parameters – stationary/homogeneous markov model  
(independent of time  $t$ )



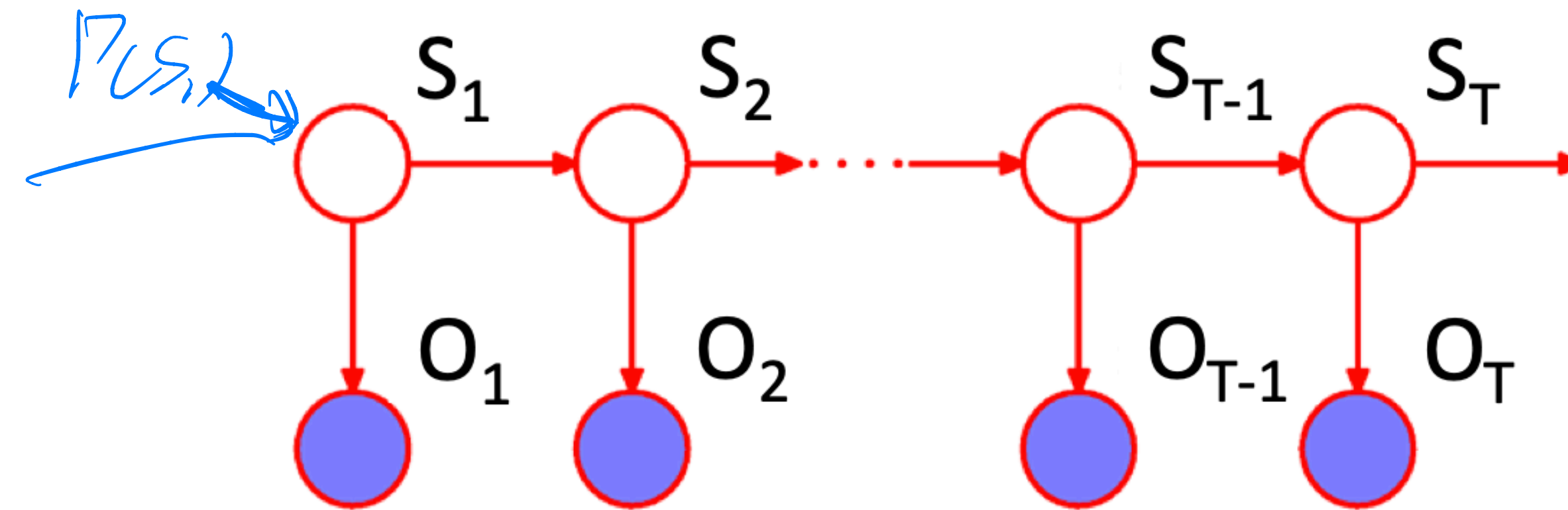
# Hidden Markov Models

- Parameters – stationary/homogeneous markov model  
(independent of time  $t$ )

Initial probabilities

$$p(S_1 = i) = \pi_i$$

*prob*



# Hidden Markov Models

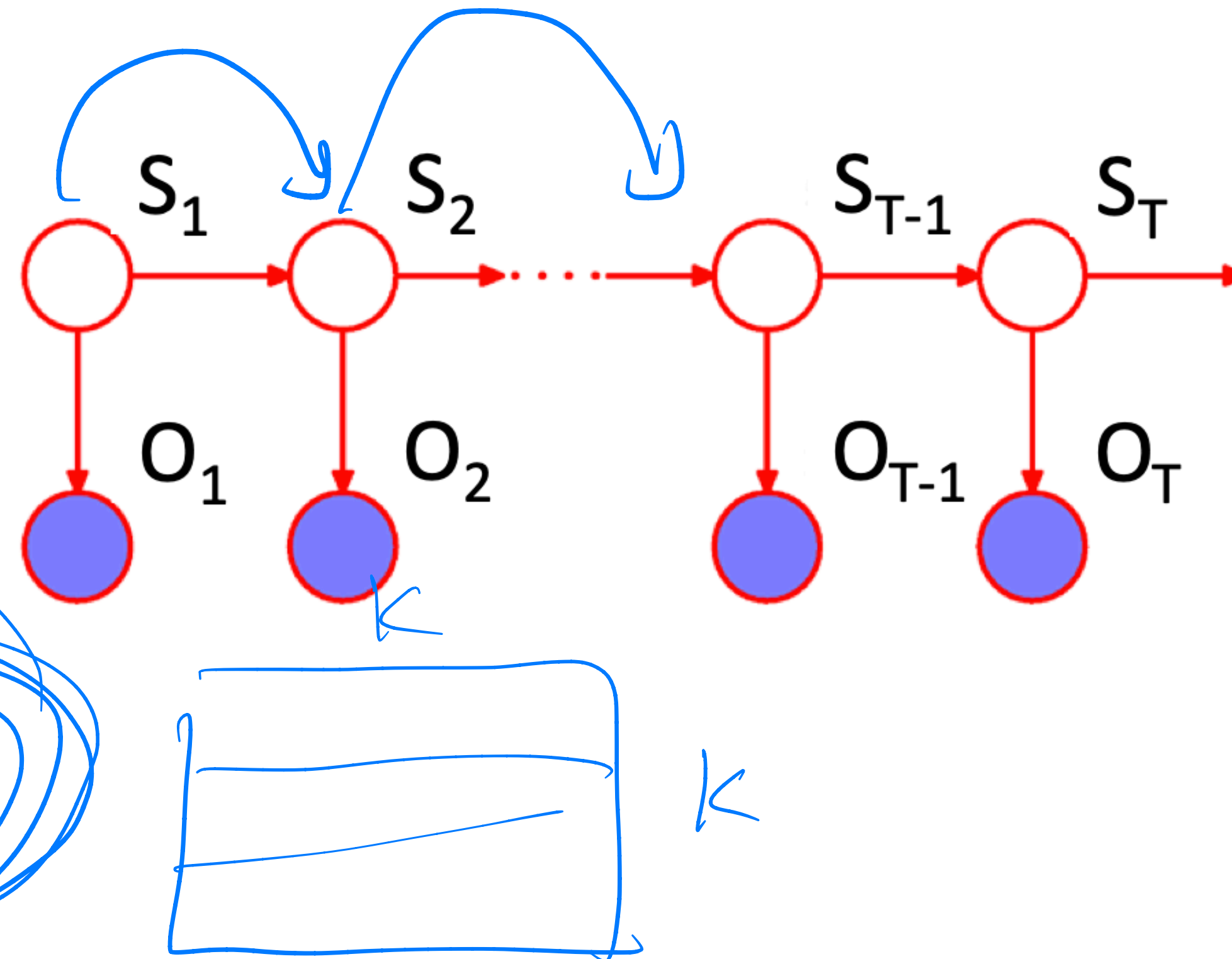
- Parameters – stationary/homogeneous markov model  
(independent of time  $t$ )

Initial probabilities

$$p(S_1 = i) = \pi_i$$

Transition probabilities

$$p(S_t = j | S_{t-1} = i) = p_{ij}$$



# Hidden Markov Models

- Parameters – stationary/homogeneous markov model (independent of time  $t$ )

Initial probabilities

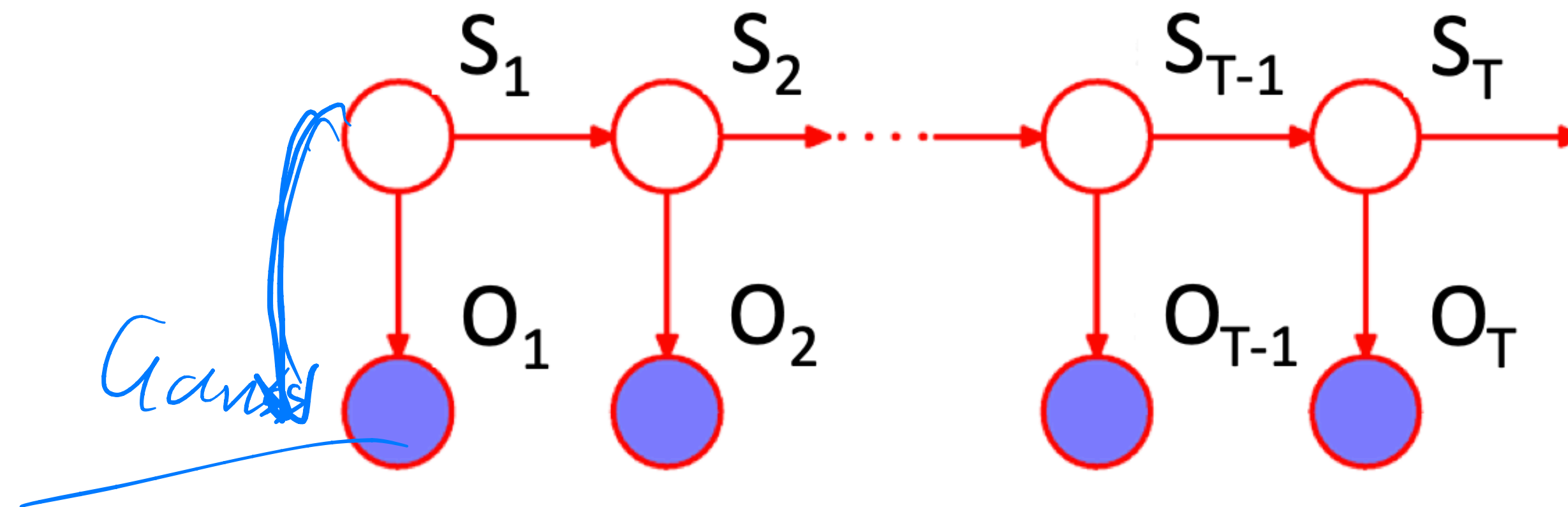
$$p(S_1 = i) = \pi_i$$

Transition probabilities

$$p(S_t = j | S_{t-1} = i) = p_{ij}$$

Emission probabilities

$$p(O_t = y | S_t = i) = q_i^y$$



# Hidden Markov Models

- Parameters – stationary/homogeneous markov model (independent of time  $t$ )

Initial probabilities

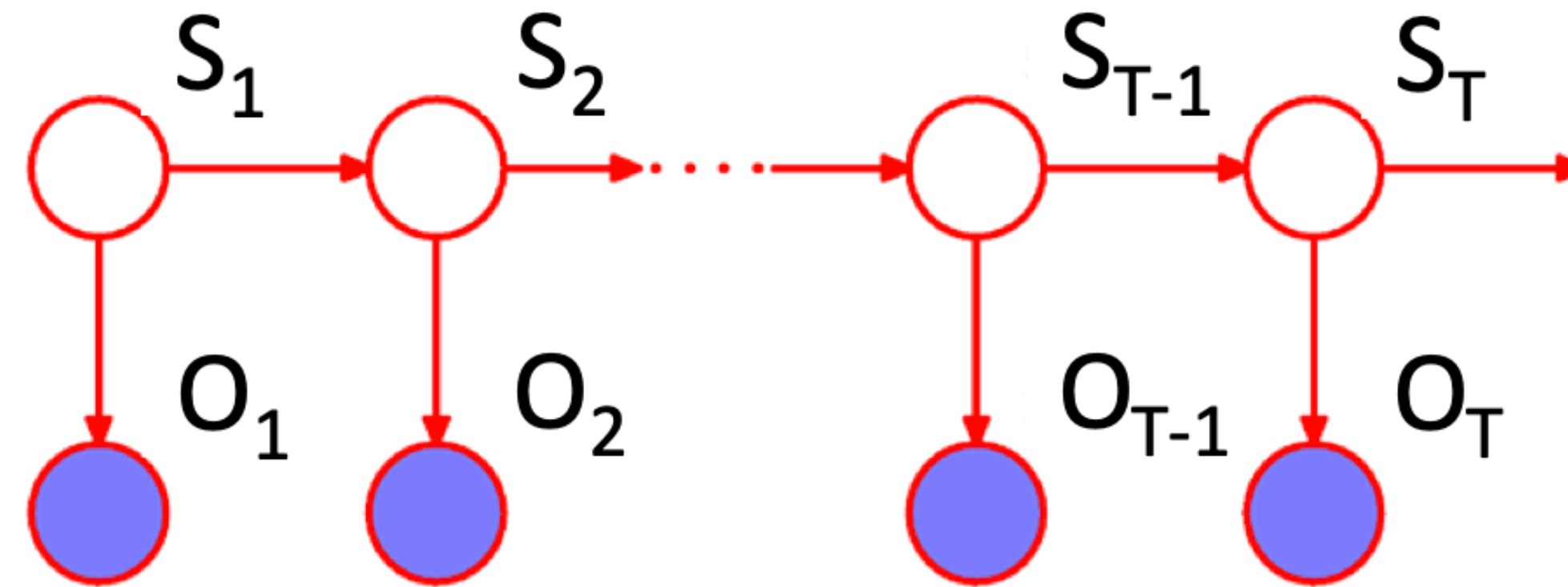
$$p(S_1 = i) = \pi_i$$

Transition probabilities

$$p(S_t = j | S_{t-1} = i) = p_{ij}$$

Emission probabilities

$$p(O_t = y | S_t = i) = q_i^y$$



$$p(\{S_t\}_{t=1}^T, \{O_t\}_{t=1}^T) = p(S_1) \prod_{t=2}^T p(S_t | S_{t-1}) \prod_{t=1}^T p(O_t | S_t)$$

2 loaded dice

# HMM Example

## The Dishonest Casino

A casino has two dices:

Fair dice

$$P(1) = P(2) = P(3) = P(5) = P(6) = 1/6$$

Loaded dice

$$P(1) = P(2) = P(3) = P(5) = 1/10$$

$$P(6) = 1/2$$

Casino player switches back-&-forth between fair and loaded die with 5% probability

$P(c|s)$



evens!

$$P(c_{t+1} = \text{fair} | s_t = \text{fair})$$

$$= 95\%$$

$$P(c_{t+1} = \text{fair} | s_t = \text{loaded})$$

$$= 5\%$$

dice:

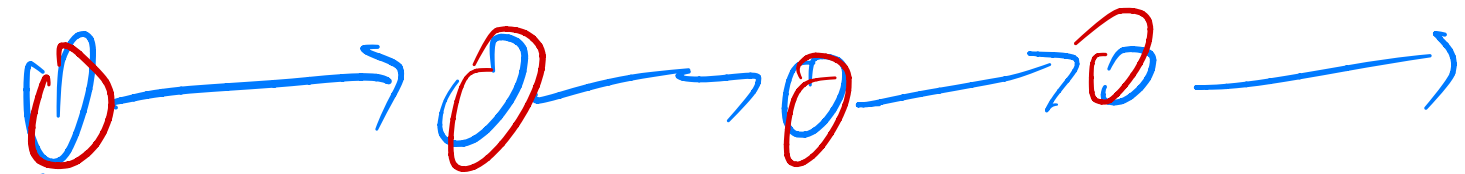
$S$ : binary fair or loaded

$O$ : 6 values

# HMM Example

# HMM Example

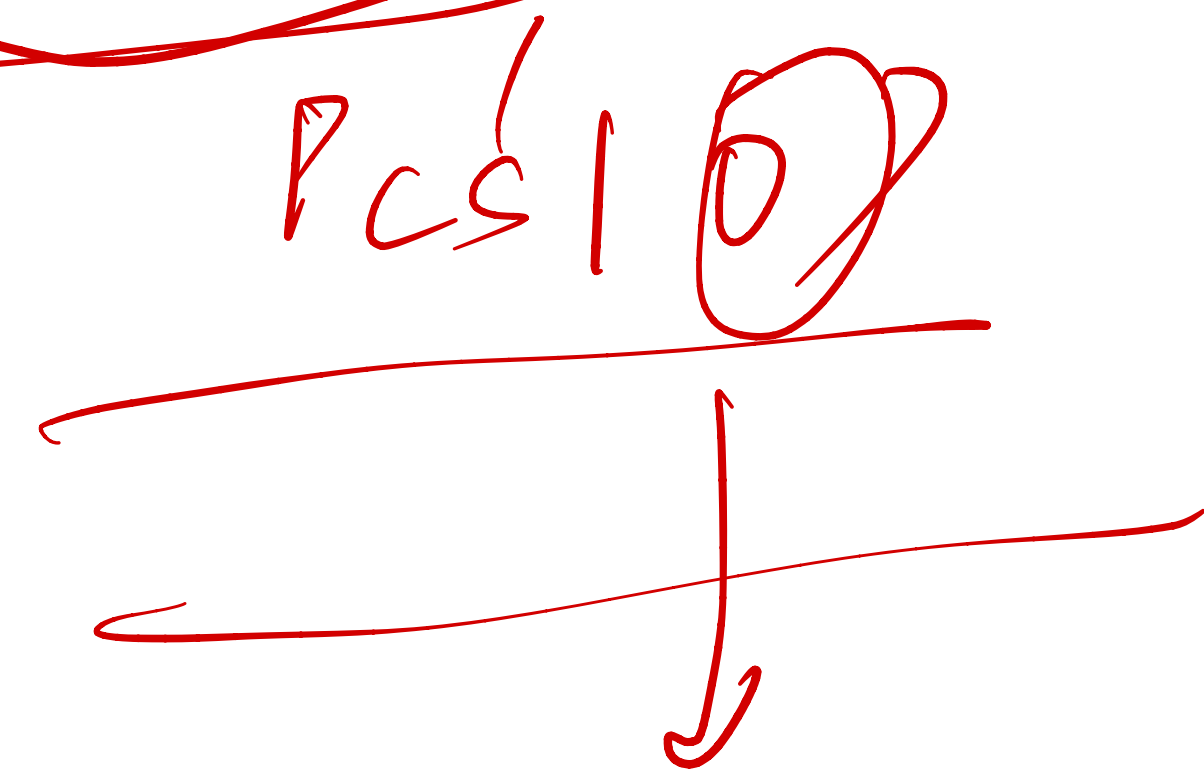
dice



**GIVEN:** A sequence of rolls by the casino player

1245526462146146136136661664661636616366163616515615115146123562344

rolls



# HMM Example

**GIVEN:** A sequence of rolls by the casino player

1245526462146146136136661664661636616366163616515615115146123562344

Question

1. How likely is the sequence given our model?

This is the **evaluation** problem in HMMs

# HMM Example

**GIVEN:** A sequence of rolls by the casino player

1245526462146146136136661664661636616366163616515615115146123562344

Question

1. How likely is the sequence given our model?

This is the evaluation problem in HMMs

2. What portion of the sequence was generated with the fair die, and what portion with the loaded die

This is the decoding question in HMMs

# HMM Example

**GIVEN:** A sequence of rolls by the casino player

1245526462146146136136661664661636616366163616515615115146123562344

Question

1. How likely is the sequence given our model?

This is the evaluation problem in HMMs

2. What portion of the sequence was generated with the fair die, and what portion with the loaded die

This is the decoding question in HMMs

3. How “loaded” is the loaded die? How “fair” is the fair die? How often does the casino player change from fair to loaded, and back?

This is the learning question in HMMs

# Three Main Problems in HMMs

# Three Main Problems in HMMs

- **Evaluation** – Given HMM parameters & observation seqn  $\{O_t\}_{t=1}^T$   
find  $p(\{O_t\}_{t=1}^T | \theta)$  prob of observed sequence

# Three Main Problems in HMMs

- **Evaluation** – Given HMM parameters & observation seqn  $\{O_t\}_{t=1}^T$   
find  $p(\{O_t\}_{t=1}^T | \theta)$  prob of observed sequence
- **Decoding** – Given HMM parameters & observation seqn  $\{O_t\}_{t=1}^T$   
find  $\arg \max_{s_1, \dots, s_T} p(\{S_t\}_{t=1}^T | \{O_t\}_{t=1}^T, \theta)$  most probable  
sequence of hidden states

# Three Main Problems in HMMs

- **Evaluation** – Given HMM parameters & observation seqn  $\{O_t\}_{t=1}^T$   
find  $p(\{O_t\}_{t=1}^T | \theta)$  prob of observed sequence
- **Decoding** – Given HMM parameters & observation seqn  $\{O_t\}_{t=1}^T$   
find  $\arg \max_{s_1, \dots, s_T} p(\{S_t\}_{t=1}^T | \{O_t\}_{t=1}^T, \theta)$  most probable  
sequence of hidden states
- **Learning** – Given HMM with unknown parameters and  $\{O_t\}_{t=1}^T$   
observation sequence  
find  $\arg \max_{\theta} p(\{O_t\}_{t=1}^T | \theta)$  parameters that maximize  
likelihood of observed data

# HMM Algorithms

# HMM Algorithms

- **Evaluation** – What is the probability of the observed sequence? **Forward Algorithm**

# HMM Algorithms

- **Evaluation** – What is the probability of the observed sequence? **Forward Algorithm**
- **Decoding** – What is the probability that the third roll was loaded given the observed sequence? **Forward-Backward Algorithm**

# HMM Algorithms

- **Evaluation** – What is the probability of the observed sequence? **Forward Algorithm**
- **Decoding** – What is the probability that the third roll was loaded given the observed sequence? **Forward-Backward Algorithm**
  - What is the most likely die sequence given the observed sequence? **Viterbi Algorithm**

one state

$S_3$

$S_1, S_2, \dots, S_T$

# HMM Algorithms

- **Evaluation** – What is the probability of the observed sequence? **Forward Algorithm**
- **Decoding** – What is the probability that the third roll was loaded given the observed sequence? **Forward-Backward Algorithm**
  - What is the most likely die sequence given the observed sequence? **Viterbi Algorithm**
- **Learning** – Under what parameterization is the observed sequence most probable? **Baum-Welch Algorithm (EM)**

# Evaluation Problem

# Evaluation Problem

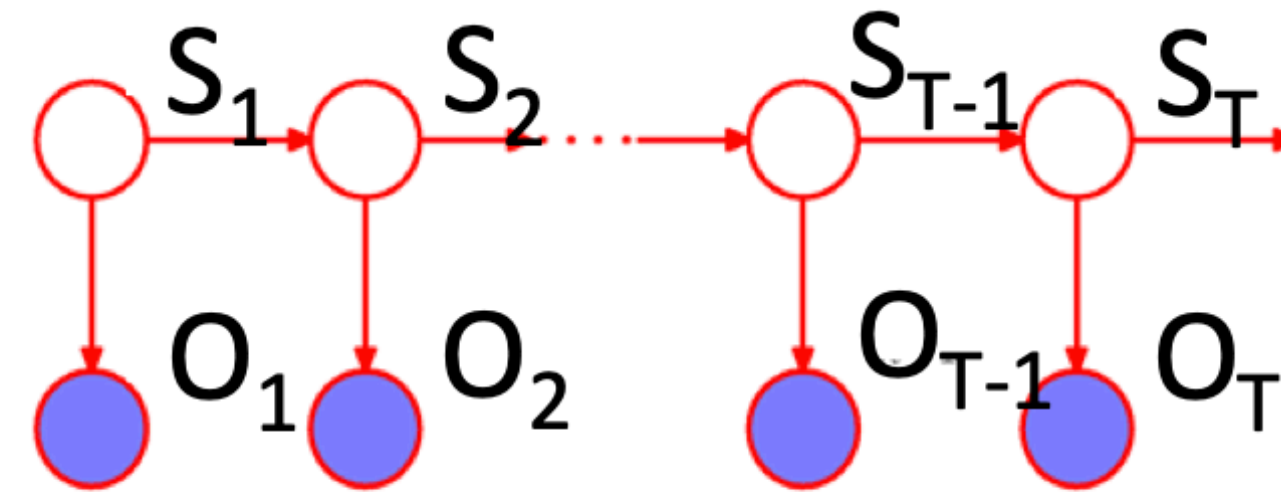
- Given HMM parameters  $p(S_1), p(S_t|S_{t-1}), p(O_t|S_t)$  & observation sequence  $\{O_t\}_{t=1}^T$

# Evaluation Problem

- Given HMM parameters  $p(S_1), p(S_t|S_{t-1}), p(O_t|S_t)$  & observation sequence  $\{O_t\}_{t=1}^T$

find probability of observed sequence

$$p(\{O_t\}_{t=1}^T) = \sum_{S_1, \dots, S_T} p(\{O_t\}_{t=1}^T, \{S_t\}_{t=1}^T)$$
$$= \sum_{S_1, \dots, S_T} p(S_1) \prod_{t=2}^T p(S_t|S_{t-1}) \prod_{t=1}^T p(O_t|S_t)$$

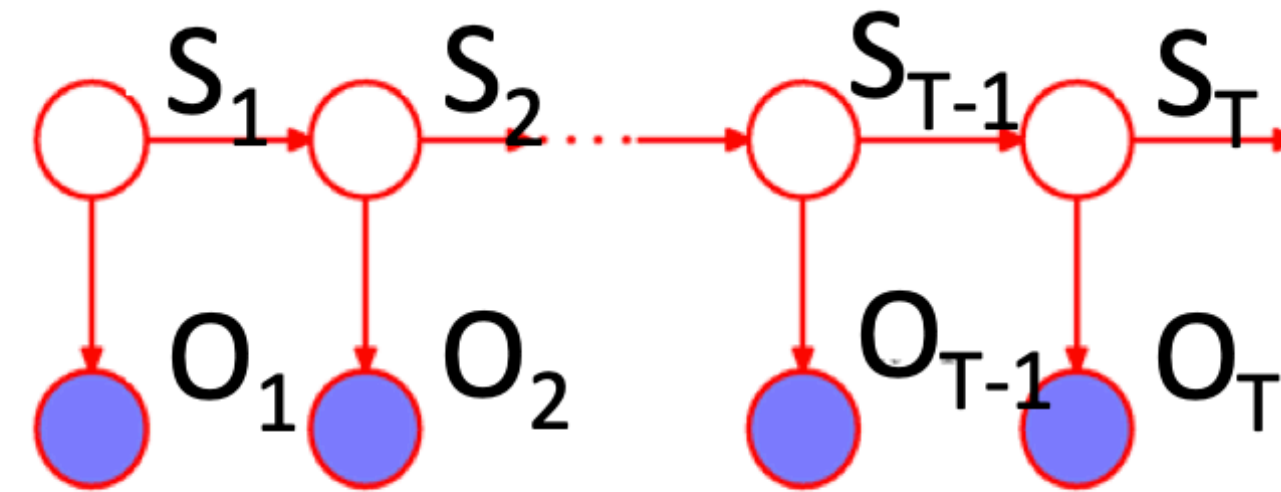


# Evaluation Problem

$P(S_2|S_1)$   $P(S_3|S_2)$   $P(O_1|S_1)$

- Given HMM parameters  $p(S_1), p(S_t|S_{t-1}), p(O_t|S_t)$  & observation sequence  $\{O_t\}_{t=1}^T$

find probability of observed sequence



$$p(\{O_t\}_{t=1}^T) = \sum_{S_1, \dots, S_T} p(\{O_t\}_{t=1}^T, \{S_t\}_{t=1}^T)$$


$$= \sum_{S_1, \dots, S_T} p(S_1) \prod_{t=2}^T p(S_t|S_{t-1}) \prod_{t=1}^T p(O_t|S_t)$$

$S_1$   
 $S_2$

requires summing over all possible hidden state values at all times –  $K^T$  exponential # terms!

# Forward Probability

# Forward Probability

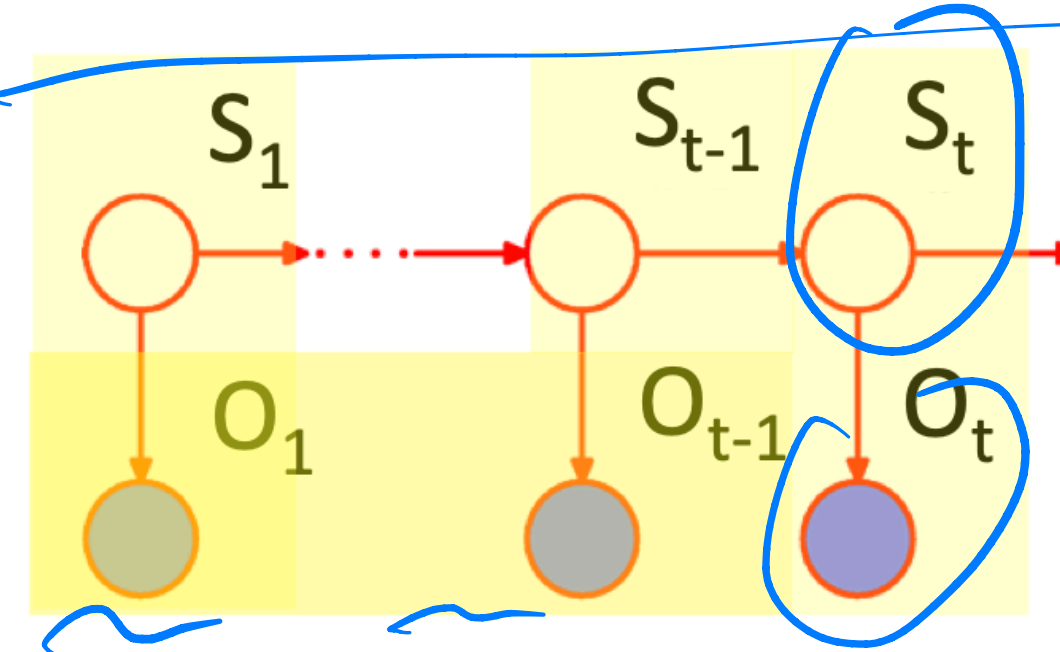
$$p(\{O_t\}_{t=1}^T) = \sum_k p(\{O_t\}_{t=1}^T, S_T = k) = \sum_k \alpha_T^k$$


# Forward Probability

$$p(\{O_t\}_{t=1}^T) = \sum_k p(\{O_t\}_{t=1}^T, S_T = k) = \sum_k \alpha_T^k$$

Compute forward probability  $\alpha_t^k$  recursively over t

$$\alpha_t^k := p(O_1, \dots, O_t, S_t = k)$$



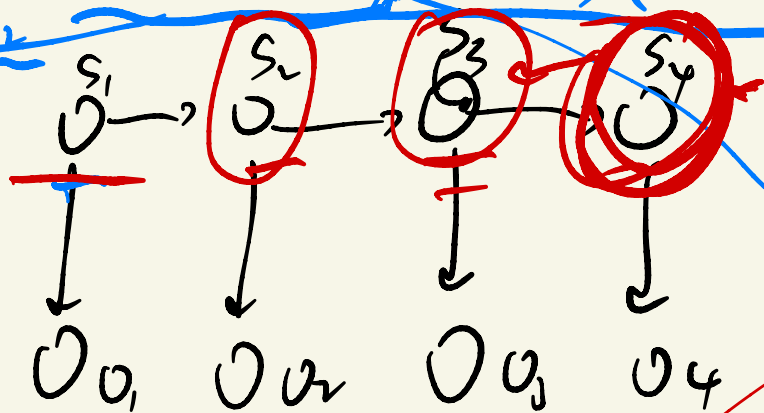
Introduce  $S_{t-1}$

Chain rule

Markov assumption

$$= p(O_t | S_t = k) \sum_i \alpha_{t-1}^i p(S_t = k | S_{t-1} = i)$$

$$\sum_{S_2} P(S_2, O_1, O_2) P(S_3|S_2) P(O_3|S_3) \quad P(S_3|S_2) P(O_3|S_3)$$



$$\sum_{S_3} P(S_1, S_2, O_1, O_2)$$

$$= P(S_2, O_1, O_2)$$

$$P(O_1, \dots, O_4) = \sum_{S_1} \sum_{S_2} \sum_{S_3} \sum_{S_4} P(S_1) P(O_1|S_1) P(S_2|S_1) P(S_3|S_2) P(O_2|S_2) \dots$$

$$= \sum_{S_2} \sum_{S_3} \sum_{S_4} \left( \sum_{S_1} P(S_1) P(O_1|S_1) P(S_2|S_1) P(S_3|S_2) P(O_2|S_2) \right)$$

$$\sum_{S_2} P(S_2, S_3, O_1, O_2, O_3)$$

$$= P(S_3, O_1, O_2, O_3)$$

$$P(O_2|S_2)$$

$$\sum_{s_1} P(s_1) P(u_1 | s_1) P(s_2 | s_1)$$

$$= \sum_{s_1} P(s_1, s_2, u_1)$$

$$P(u_1, \dots, u_n, S_{\text{in}} = s_1)$$

$$= P(s_2, u_1)$$

$$P(s_1, u_1, u_2)$$

$$\cancel{P(s_2, u_1, u_2)}$$

$$\propto P(u_2 | s_2)$$

$$= P(s_2, u_1, u_2)$$

# Forward Algorithm

Can compute  $\alpha_t^k$  for all  $k, t$  using dynamic programming:

# Forward Algorithm

Can compute  $\alpha_t^k$  for all  $k, t$  using dynamic programming:


- Initialize:  $\alpha_1^k = p(O_1 | S_1 = k) p(S_1 = k)$  for all  $k$

# Forward Algorithm

Can compute  $\alpha_t^k$  for all  $k, t$  using dynamic programming:

- Initialize:  $\alpha_1^k = p(O_1 | S_1 = k) p(S_1 = k)$  for all  $k$

- Iterate: for  $t = 2, \dots, T$

$$\alpha_t^k = p(O_t | S_t = k) \sum_i \alpha_{t-1}^i p(S_t = k | S_{t-1} = i) \quad \text{for all } k$$


# Forward Algorithm

Can compute  $\alpha_t^k$  for all  $k, t$  using dynamic programming:

- Initialize:  $\alpha_1^k = p(O_1 | S_1 = k) p(S_1 = k)$  for all  $k$

- Iterate: for  $t = 2, \dots, T$

$$\alpha_t^k = p(O_t | S_t = k) \sum_i \alpha_{t-1}^i p(S_t = k | S_{t-1} = i) \quad \text{for all } k$$

- Termination:

$$p(\{O_t\}_{t=1}^T) = \sum_k \alpha_T^k$$

# Forward Algorithm

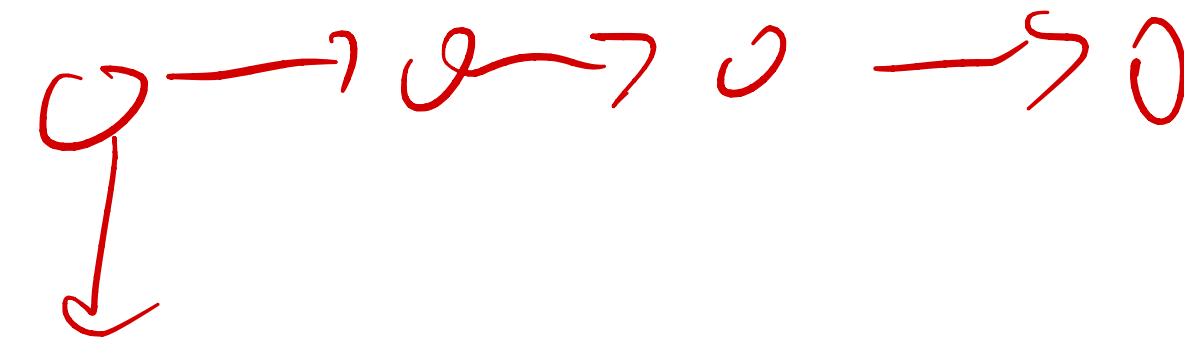
Can compute  $\alpha_t^k$  for all  $k, t$  using dynamic programming:

- Initialize:  $\alpha_1^k = p(O_1 | S_1 = k) p(S_1 = k)$  for all  $k$

- Iterate: for  $t = 2, \dots, T$

$$\alpha_t^k = p(O_t | S_t = k) \sum_i \alpha_{t-1}^i p(S_t = k | S_{t-1} = i) \quad \text{for all } k$$

- Termination:  $p(\{O_t\}_{t=1}^T) = \sum_k \alpha_T^k$



Can we do in the backward direction?

# Forward Algorithm

Can compute  $\alpha_t^k$  for all  $k, t$  using dynamic programming:

- Initialize:  $\alpha_1^k = p(O_1 | S_1 = k) p(S_1 = k)$  for all  $k$

- Iterate: for  $t = 2, \dots, T$

$$\alpha_t^k = p(O_t | S_t = k) \sum_i \alpha_{t-1}^i p(S_t = k | S_{t-1} = i) \quad \text{for all } k$$

- Termination:  $p(\{O_t\}_{t=1}^T) = \sum_k \alpha_T^k$

You will try this in your HW

Can we do in the backward direction?

# Decoding Problem 1

# Decoding Problem 1

- Given HMM parameters  $p(S_1), p(S_t|S_{t-1}), p(O_t|S_t)$  & observation sequence  $\{O_t\}_{t=1}^T$

find probability that hidden state at time t was k  $p(S_t = k | \{O_t\}_{t=1}^T)$

Posterior

# Decoding Problem 1

- Given HMM parameters  $p(S_1), p(S_t|S_{t-1}), p(O_t|S_t)$  & observation sequence  $\{O_t\}_{t=1}^T$

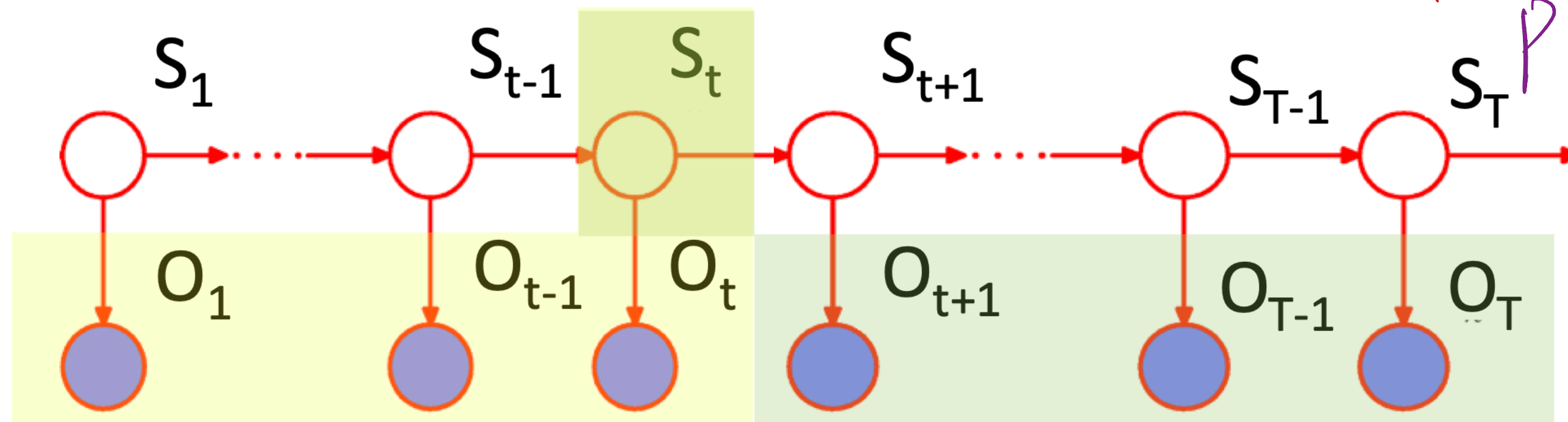
find probability that hidden state at time t was k  $p(S_t = k | \{O_t\}_{t=1}^T)$

$$p(S_t = k, \{O_t\}_{t=1}^T) = p(O_1, \dots, O_t, S_t = k, O_{t+1}, \dots, O_T)$$

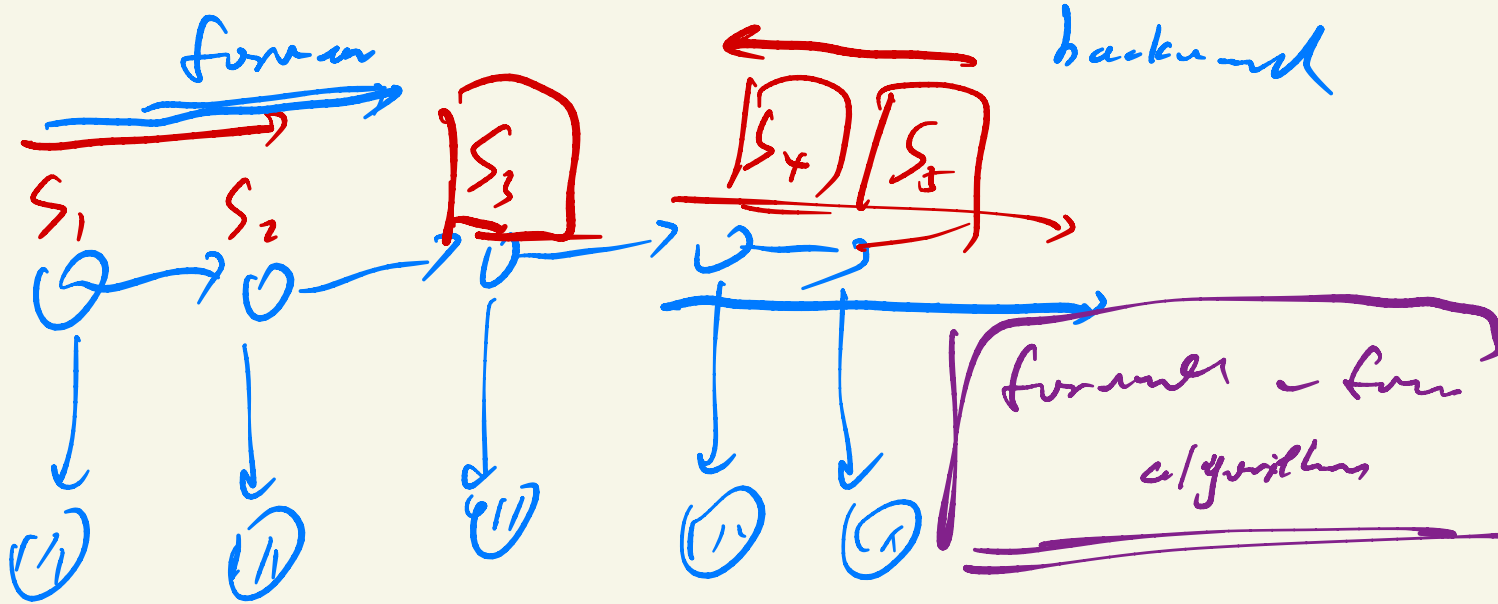
$$= p(O_1, \dots, O_t, S_t = k) p(O_{t+1}, \dots, O_T | S_t = k)$$

$\alpha_t^k$ 
 $\beta_t^k$

Compute recursively



*Handwritten notes:*  
 $p(O_1, \dots, O_t, S_t = k)$   
 $p(O_{t+1}, \dots, O_T | S_t = k)$   
 $S_t = k$



$$P(S_1 | 0_1, \dots, 0_5) \neq P(S_3 | 0_1, \dots, 0_5)$$

$\bar{S}_1$   $\bar{S}_2$   $\bar{S}_4$   $\bar{S}_5$

$$P(S_1, \dots, S_5 | 0_1, \dots, 0_5)$$

# Backward Algorithm

# Backward Algorithm

Can compute  $\beta_t^k$  for all  $k, t$  using dynamic programming:

- Initialize:  $\beta_T^k = 1$  for all  $k$

# Backward Algorithm

Can compute  $\beta_t^k$  for all  $k, t$  using dynamic programming:

- Initialize:  $\beta_T^k = 1$   $P(O_{T+1} | S_T = k)$   
for all  $k$  Why this initialization?

# Backward Algorithm

Can compute  $\beta_t^k$  for all  $k, t$  using dynamic programming:

- Initialize:  $\beta_T^k = 1$  for all  $k$

Why this initialization?

$$\beta_T^k = P(O_{T+1} = O_T | S_T = k)$$

- Iterate: for  $t = T-1, \dots, 1$

$$\beta_t^k = \sum_i p(S_{t+1} = i | S_t = k) p(O_{t+1} | S_{t+1} = i) \beta_{t+1}^i \quad \text{for all } k$$

# Backward Algorithm

Can compute  $\beta_t^k$  for all  $k, t$  using dynamic programming:

- Initialize:  $\beta_T^k = 1$  for all  $k$

Why this initialization?

$$P(S_T = 1, O_1 \dots O_T) = 0.1$$

$$P(S_T = 2, O_1 \dots O_T) = 0.2$$

- Iterate: for  $t = T-1, \dots, 1$

$$\beta_t^k = \sum_i p(S_{t+1} = i | S_t = k) p(O_{t+1} | S_{t+1} = i) \beta_{t+1}^i \quad \text{for all } k$$

- Termination:  $p(S_t = k, \{O_t\}_{t=1}^T) = \alpha_t^k \beta_t^k$

$$p(S_t = k | \{O_t\}_{t=1}^T) = \frac{p(S_t = k, \{O_t\}_{t=1}^T)}{p(\{O_t\}_{t=1}^T)} = \frac{\alpha_t^k \beta_t^k}{\sum_i \alpha_t^i \beta_t^i}$$

$$P(S_T = 1, O_1 \dots O_T)$$

$$\frac{0.1}{0.1 + 0.2}$$

# Backward Algorithm

Can compute  $\beta_t^k$  for all  $k, t$  using dynamic programming:

- Initialize:  $\beta_T^k = 1$  for all  $k$  **Why this initialization?**

- Iterate: for  $t = T-1, \dots, 1$

$$\beta_t^k = \sum_i p(S_{t+1} = i | S_t = k) p(O_{t+1} | S_{t+1} = i) \beta_{t+1}^i \quad \text{for all } k$$

- Termination:  $p(S_t = k, \{O_t\}_{t=1}^T) = \alpha_t^k \beta_t^k$

$$p(S_t = k | \{O_t\}_{t=1}^T) = \frac{p(S_t = k, \{O_t\}_{t=1}^T)}{p(\{O_t\}_{t=1}^T)} = \frac{\alpha_t^k \beta_t^k}{\sum_i \alpha_t^i \beta_t^i}$$

**Can we compute  $\beta$  in a forward manner?**

# Most Likely State vs. Most Likely Sequence

# Most Likely State vs. Most Likely Sequence

- Most likely state assignment at time  $t$

$$\arg \max_k p(S_t = k | \{O_t\}_{t=1}^T) = \arg \max_k \alpha_t^k \beta_t^k$$

E.g. Which die was most likely used by the casino in the third roll given the observed sequence?

# Most Likely State vs. Most Likely Sequence

- Most likely state assignment at time t

$$\arg \max_k p(S_t = k | \{O_t\}_{t=1}^T) = \arg \max_k \alpha_t^k \beta_t^k$$

E.g. Which die was most likely used by the casino in the third roll given the observed sequence?

- Most likely assignment of state sequence

$$\arg \max_{\{S_t\}_{t=1}^T} p(\{S_t\}_{t=1}^T | \{O_t\}_{t=1}^T)$$

# Most Likely State vs. Most Likely Sequence

- Most likely state assignment at time t

$$\arg \max_k p(S_t = k | \{O_t\}_{t=1}^T) = \arg \max_k \alpha_t^k \beta_t^k$$

E.g. Which die was most likely used by the casino in the third roll given the observed sequence?

- Most likely assignment of state sequence

$$\arg \max_{\{S_t\}_{t=1}^T} p(\{S_t\}_{t=1}^T | \{O_t\}_{t=1}^T)$$

Are the solutions the same?

agrees

same only when

$P(S_1, S_2, S_3 | O_1, O_2, O_3) = P(S_1 | O_1, O_2, O_3) \cdot P(S_2 | O_1, O_2, O_3) \cdot P(S_3 | O_1, O_2, O_3)$

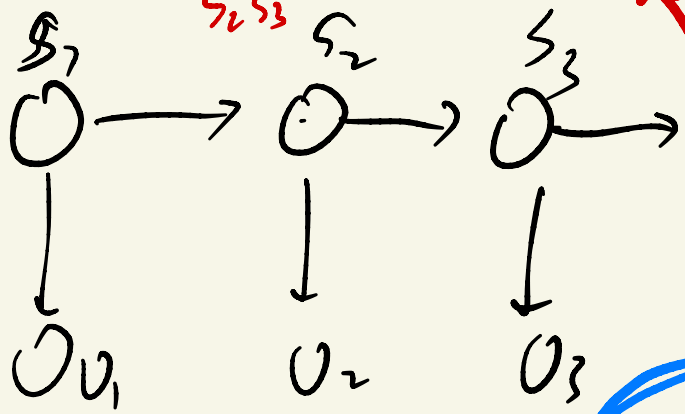
# Decoding Problem 2

# Decoding Problem 2

- Given HMM parameters  $p(S_1), p(S_t|S_{t-1}), p(O_t|S_t)$  & observation sequence  $\{O_t\}_{t=1}^T$

find most likely assignment of state sequence

$$\arg \max_{\{S_t\}_{t=1}^T} p(\{S_t\}_{t=1}^T | \{O_t\}_{t=1}^T) = \arg \max_{\{S_t\}_{t=1}^T} p(\{S_t\}_{t=1}^T, \{O_t\}_{t=1}^T)$$



max  
S<sub>2</sub>S<sub>3</sub>

P(S<sub>1</sub>|S<sub>2</sub>)P(S<sub>2</sub>|S<sub>3</sub>)P(S<sub>3</sub>|S<sub>1</sub>)

max  
S<sub>1</sub>

P(S<sub>1</sub>)P(S<sub>2</sub>|S<sub>1</sub>)P(S<sub>3</sub>|S<sub>1</sub>)

S<sub>2</sub> k

no S<sub>1</sub>

f(S<sub>2</sub>)

f(S<sub>2</sub>=k)

k<sup>2</sup>

argmax

max  
S<sub>1</sub>

max  
S<sub>2</sub>

max  
S<sub>3</sub>

P(S<sub>1</sub>)P(S<sub>2</sub>|S<sub>1</sub>)P(S<sub>3</sub>|S<sub>1</sub>)

S<sub>2</sub> P(S<sub>2</sub>|S<sub>2</sub>)P(S<sub>3</sub>|S<sub>2</sub>)

1	2	3	4
---	---	---	---

argmax

# Decoding Problem 2

- Given HMM parameters  $p(S_1), p(S_t|S_{t-1}), p(O_t|S_t)$  & observation sequence  $\{O_t\}_{t=1}^T$

find most likely assignment of state sequence

$$\begin{aligned} \arg \max_{\{S_t\}_{t=1}^T} p(\{S_t\}_{t=1}^T | \{O_t\}_{t=1}^T) &= \arg \max_{\{S_t\}_{t=1}^T} p(\{S_t\}_{t=1}^T, \{O_t\}_{t=1}^T) \\ &= \arg \max_k \max_{\{S_t\}_{t=1}^{T-1}} p(S_T = k, \underbrace{\{S_t\}_{t=1}^{T-1}, \{O_t\}_{t=1}^T}_{V_T^k}) \end{aligned}$$

Compute recursively

# Decoding Problem 2

- Given HMM parameters  $p(S_1), p(S_t|S_{t-1}), p(O_t|S_t)$  & observation sequence  $\{O_t\}_{t=1}^T$

find most likely assignment of state sequence

$$\begin{aligned} \arg \max_{\{S_t\}_{t=1}^T} p(\{S_t\}_{t=1}^T | \{O_t\}_{t=1}^T) &= \arg \max_{\{S_t\}_{t=1}^T} p(\{S_t\}_{t=1}^T, \{O_t\}_{t=1}^T) \\ &= \arg \max_k \max_{\{S_t\}_{t=1}^{T-1}} p(S_T = k, \underbrace{\{S_t\}_{t=1}^{T-1}, \{O_t\}_{t=1}^T}_{V_T^k}) \end{aligned}$$

Compute recursively

$V_T^k$  - probability of most likely sequence of states ending at state  $S_T = k$

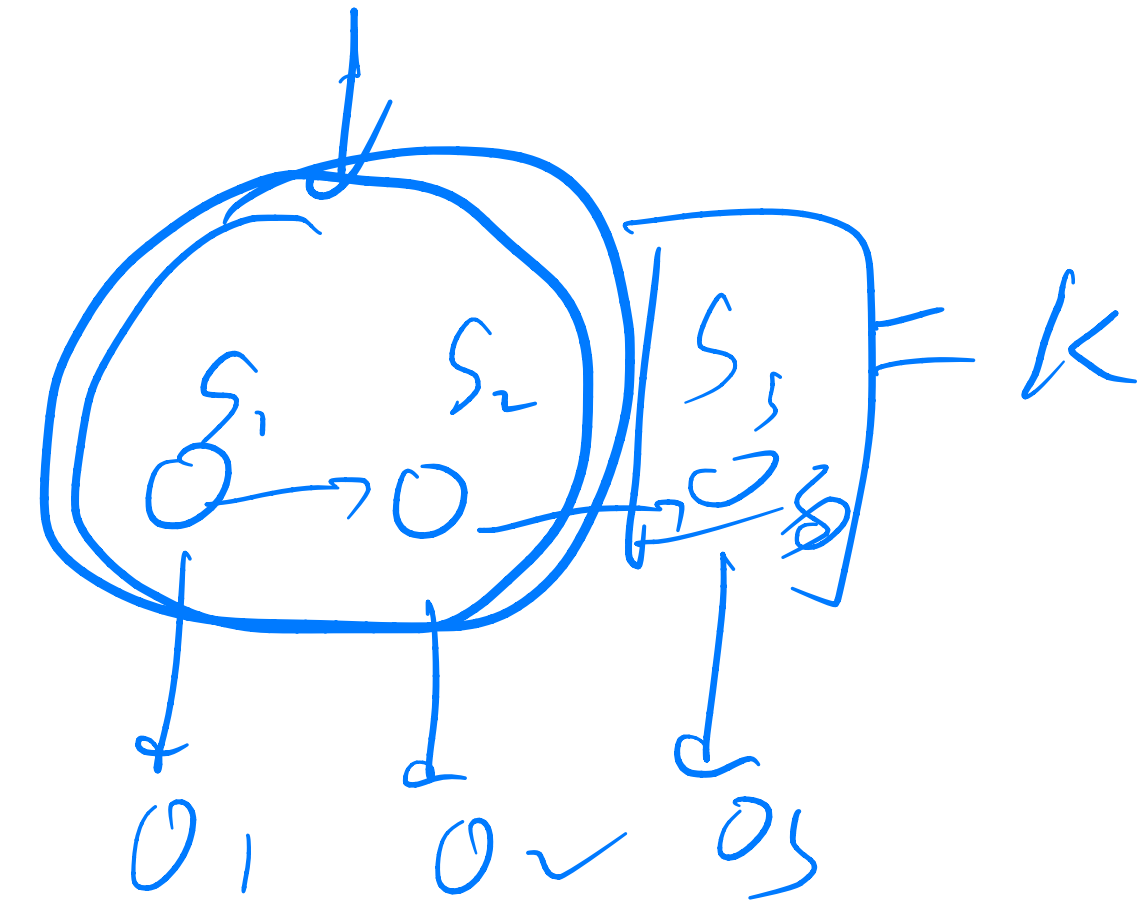
# Viterbi Decoding

# Viterbi Decoding

$$\max_{\{S_t\}_{t=1}^T} p(\{S_t\}_{t=1}^T, \{O_t\}_{t=1}^T) = \max_k V_T^k$$

# Viterbi Decoding

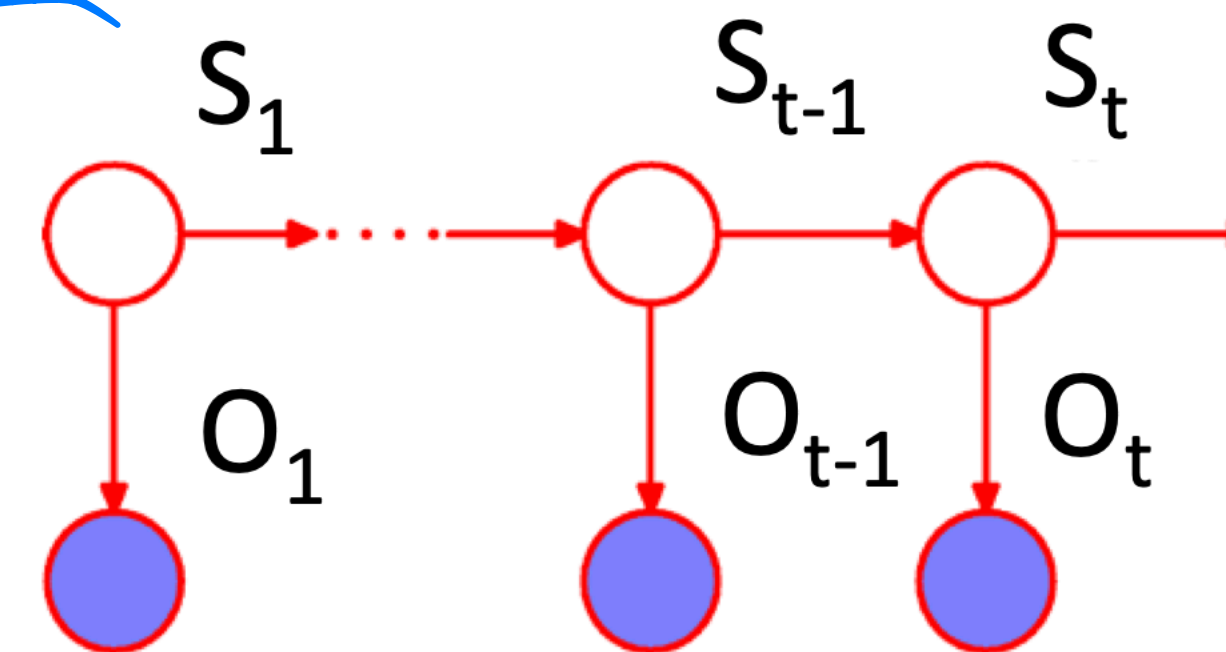
$$\max_{\{S_t\}_{t=1}^T} p(\{S_t\}_{t=1}^T, \{O_t\}_{t=1}^T) = \max_k V_T^k$$



Compute probability  $V_t^k$  recursively over  $t$

$$V_t^k := \max_{S_1, \dots, S_{t-1}} p(S_t = k, S_1, \dots, S_{t-1}, O_1, \dots, O_t)$$

- Bayes rule
- Markov assumption



$$= p(O_t | S_t = k) \max_i p(S_t = k | S_{t-1} = i) V_{t-1}^i$$

# Viterbi Algorithm

# Viterbi Algorithm

Can compute  $V_t^k$  for all  $k, t$  using dynamic programming:

- Initialize:  $V_1^k = p(O_1|S_1=k)p(S_1 = k)$  for all  $k$

- Iterate: for  $t = 2, \dots, T$

$$V_t^k = p(O_t|S_t = k) \max_i p(S_t = k|S_{t-1} = i) V_{t-1}^i \quad \text{for all } k$$

- Termination:  $\max_{\{S_t\}_{t=1}^T} p(\{S_t\}_{t=1}^T, \{O_t\}_{t=1}^T) = \max_k V_T^k$

forward

# Viterbi Algorithm

Can compute  $V_t^k$  for all  $k, t$  using dynamic programming:

- Initialize:  $V_1^k = p(O_1|S_1=k)p(S_1 = k)$  for all  $k$

- Iterate: for  $t = 2, \dots, T$

$$V_t^k = p(O_t|S_t = k) \max_i p(S_t = k|S_{t-1} = i) V_{t-1}^i \quad \text{for all } k$$

- Termination:  $\max_{\{S_t\}_{t=1}^T} p(\{S_t\}_{t=1}^T, \{O_t\}_{t=1}^T) = \max_k V_T^k$

Traceback:

$$S_T^* = \arg \max_k V_T^k$$

$$S_{t-1}^* = \arg \max_i p(S_t^*|S_{t-1} = i) V_{t-1}^i$$

# Viterbi Algorithm

Can compute  $V_t^k$  for all  $k, t$  using dynamic programming:

- Initialize:  $V_1^k = p(O_1|S_1=k)p(S_1 = k)$  for all  $k$

- Iterate: for  $t = 2, \dots, T$

$$V_t^k = p(O_t|S_t = k) \max_i p(S_t = k|S_{t-1} = i) V_{t-1}^i \quad \text{for all } k$$

- Termination:  $\max_{\{S_t\}_{t=1}^T} p(\{S_t\}_{t=1}^T, \{O_t\}_{t=1}^T) = \max_k V_T^k$

Traceback:

$$S_T^* = \arg \max_k V_T^k$$

$$S_{t-1}^* = \arg \max_i p(S_t^*|S_{t-1} = i) V_{t-1}^i$$

Can we do in the  
backward direction?

# Computational Complexity

# Computational Complexity

- What is the running time for Forward, Backward, Viterbi?

$$\alpha_t^k = q_k^{O_t} \sum_i \alpha_{t-1}^i p_{i,k}$$

$$\beta_t^k = \sum_i p_{k,i} q_i^{O_{t+1}} \beta_{t+1}^i$$

$$V_t^k = q_k^{O_t} \max_i p_{i,k} V_{t-1}^i$$

$O(K^2T)$  linear in  $T$  instead of  $O(K^T)$  exponential in  $T$ !

# Learning with EM

- Start with random initialization of parameters
- **E-step** – Fix parameters, find expected state assignments

$$\gamma_i(t) = p(S_t = i | \mathbf{O}, \theta) = \frac{\alpha_t^i \beta_t^i}{\sum_j \alpha_t^j \beta_t^j} \quad \mathbf{O} = \{O_t\}_{t=1}^T$$

Forward-Backward algorithm

# Learning with EM

- Start with random initialization of parameters
- **E-step** – Fix parameters, find expected state assignments

$$\gamma_i(t) = p(S_t = i | \mathbf{O}, \theta) = \frac{\alpha_t^i \beta_t^i}{\sum_j \alpha_t^j \beta_t^j} \quad \mathbf{O} = \{O_t\}_{t=1}^T$$

## Forward-Backward algorithm

$$\begin{aligned} \xi_{ij}(t) &= p(S_{t-1} = i, S_t = j | \mathbf{O}, \theta) \\ &= \frac{p(S_{t-1} = i | \mathbf{O}, \theta) p(S_t = j, O_t, \dots, O_T | S_{t-1} = i, \theta)}{p(O_t, \dots, O_T | S_{t-1} = i, \theta)} \\ &= \frac{\gamma_i(t-1) p_{ij} q_j^{O_t} \beta_t^j}{\beta_{t-1}^i} \end{aligned}$$

# Learning with EM

- Start with random initialization of parameters
- **E-step** – Fix parameters, find expected state assignments

$$\gamma_i(t) = p(S_t = i | \mathbf{O}, \theta) = \frac{\alpha_t^i \beta_t^i}{\sum_j \alpha_t^j \beta_t^j} \quad \mathbf{O} = \{O_t\}_{t=1}^T$$

## Forward-Backward algorithm

$$\begin{aligned} \xi_{ij}(t) &= p(S_{t-1} = i, S_t = j | \mathbf{O}, \theta) \\ &= \frac{p(S_{t-1} = i | \mathbf{O}, \theta) p(S_t = j, O_t, \dots, O_T | S_{t-1} = i, \theta)}{p(O_t, \dots, O_T | S_{t-1} = i, \theta)} \\ &= \frac{\gamma_i(t-1) p_{ij} q_j^{O_t} \beta_t^j}{\beta_{t-1}^i} \end{aligned}$$

You will derive the EM  
in your HW

# If you still remember why we do EM in the first place...

$$\begin{aligned}\ell(\phi, \mu, \Sigma) &= \sum_{i=1}^n \log p(x^{(i)}; \phi, \mu, \Sigma) \\ &= \sum_{i=1}^n \log \sum_{z^{(i)}=1}^k p(x^{(i)} | z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi).\end{aligned}$$

# If you still remember why we do EM in the first place...

$$\begin{aligned} \ell(\phi, \mu, \Sigma) &= \sum_{i=1}^n \log p(x^{(i)}; \phi, \mu, \Sigma) \\ &\neq \sum_{i=1}^n \log \sum_{z^{(i)}=1}^k p(x^{(i)} | z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi). \end{aligned}$$

1. Intractable (no closed-form for the solution)

$$\sum_z p(z, x)$$

# If you still remember why we do EM in the first place...

$$\begin{aligned}\ell(\phi, \mu, \Sigma) &= \sum_{i=1}^n \log p(x^{(i)}; \phi, \mu, \Sigma) \\ &= \sum_{i=1}^n \log \sum_{z^{(i)}=1}^k p(x^{(i)} | z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi).\end{aligned}$$

1. Intractable (no closed-form for the solution)
2. Large variance in gradient descent

# If you still remember why we do EM in the first place...

$$\begin{aligned}\ell(\phi, \mu, \Sigma) &= \sum_{i=1}^n \log p(x^{(i)}; \phi, \mu, \Sigma) \\ &= \sum_{i=1}^n \log \sum_{z^{(i)}=1}^k p(x^{(i)} | z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi).\end{aligned}$$

1. Intractable (no closed-form for the solution)
2. Large variance in gradient descent

Expectation Maximization is to address the MLE optimization problem

# If you still remember why we do EM in the first place...

$$\begin{aligned}\ell(\phi, \mu, \Sigma) &= \sum_{i=1}^n \log p(x^{(i)}; \phi, \mu, \Sigma) \\ &= \sum_{i=1}^n \log \sum_{z^{(i)}=1}^k p(x^{(i)} | z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi).\end{aligned}$$

*log PCO*



Wait, HMM has closed-form likelihood?

1. Intractable (no closed-form for the solution)
2. Large variance in gradient descent

Expectation Maximization is to address the MLE optimization problem

# If you still remember why we do EM in the first place...

$$\begin{aligned} \ell(\phi, \mu, \Sigma) &= \sum_{i=1}^n \log p(x^{(i)}; \phi, \mu, \Sigma) && \text{forward} \\ &= \sum_{i=1}^n \log \sum_{z^{(i)}=1}^k p(x^{(i)} | z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi). \end{aligned}$$

Wait, HMM has closed-form likelihood?

1. Intractable (no closed-form for the solution)
2. Large variance in gradient descent

non-convex

Expectation Maximization is to address the MLE optimization problem

Can we do MLE directly for HMM using gradient descent, without EM?

**Thank You!**