



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

COMP 5212  
Machine Learning  
Lecture 16

# Hidden Markov Models

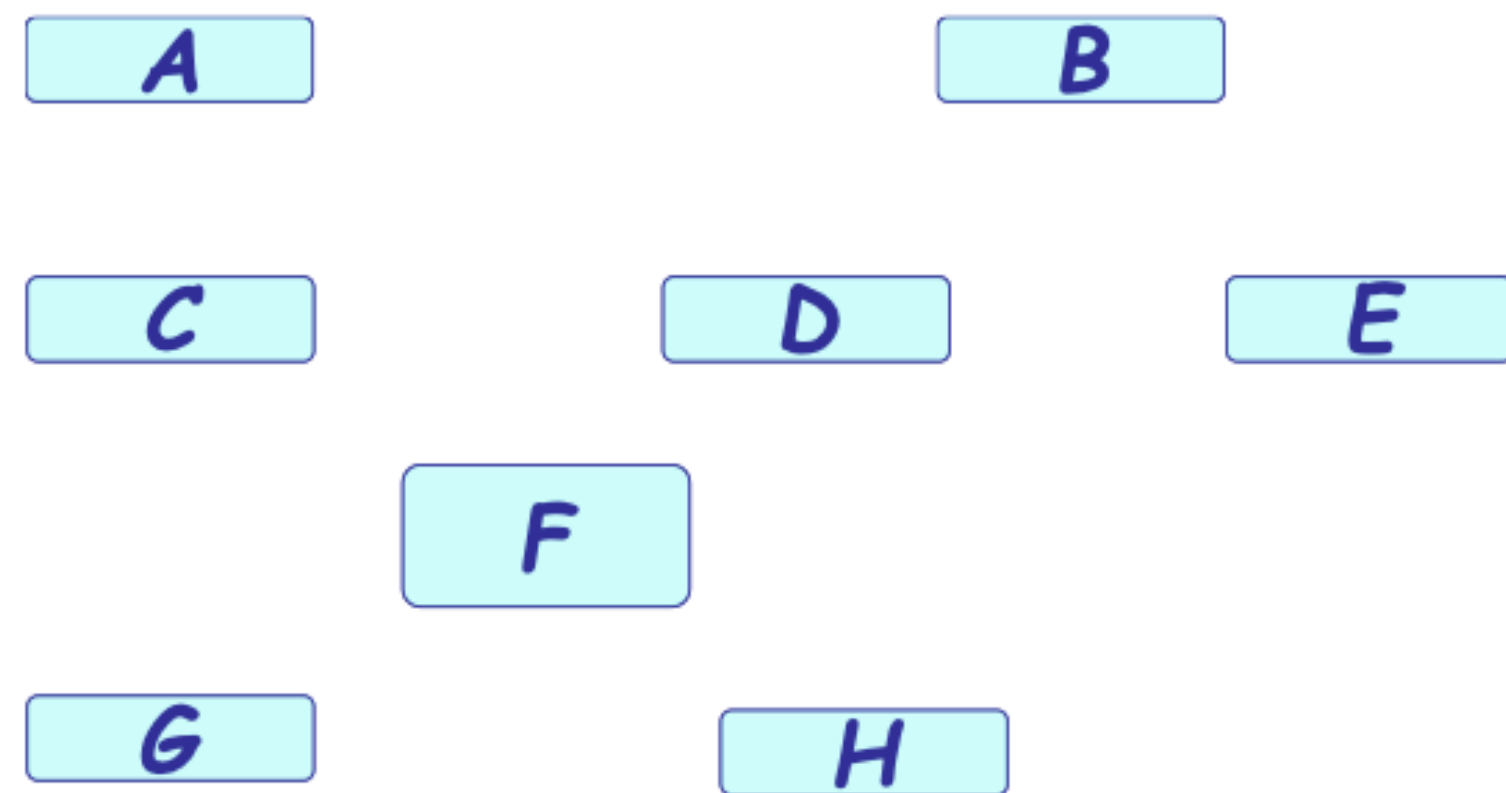
Junxian He  
Nov 5, 2024

# Announcements

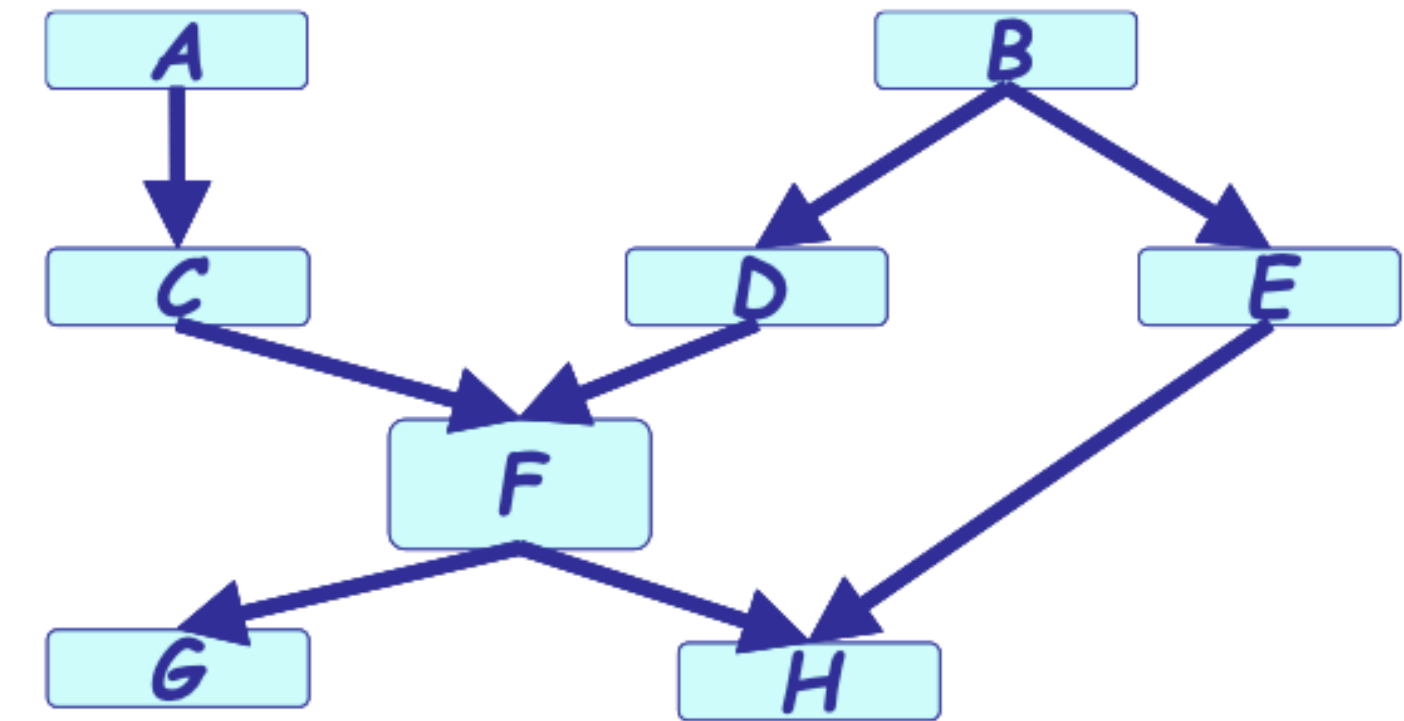
- We have a makeup lecture this Thursday on Nov 7, 7pm-820pm, at Room 2303 after we finish HMM. Attendance is not required, zoom recording will be released

# Recap: Probabilistic Graphical Models

It is a smart way to **write/specify/compose/design** exponentially-large probability distributions without paying an exponential cost, and at the same time endow the distributions with **structured semantics**



$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$



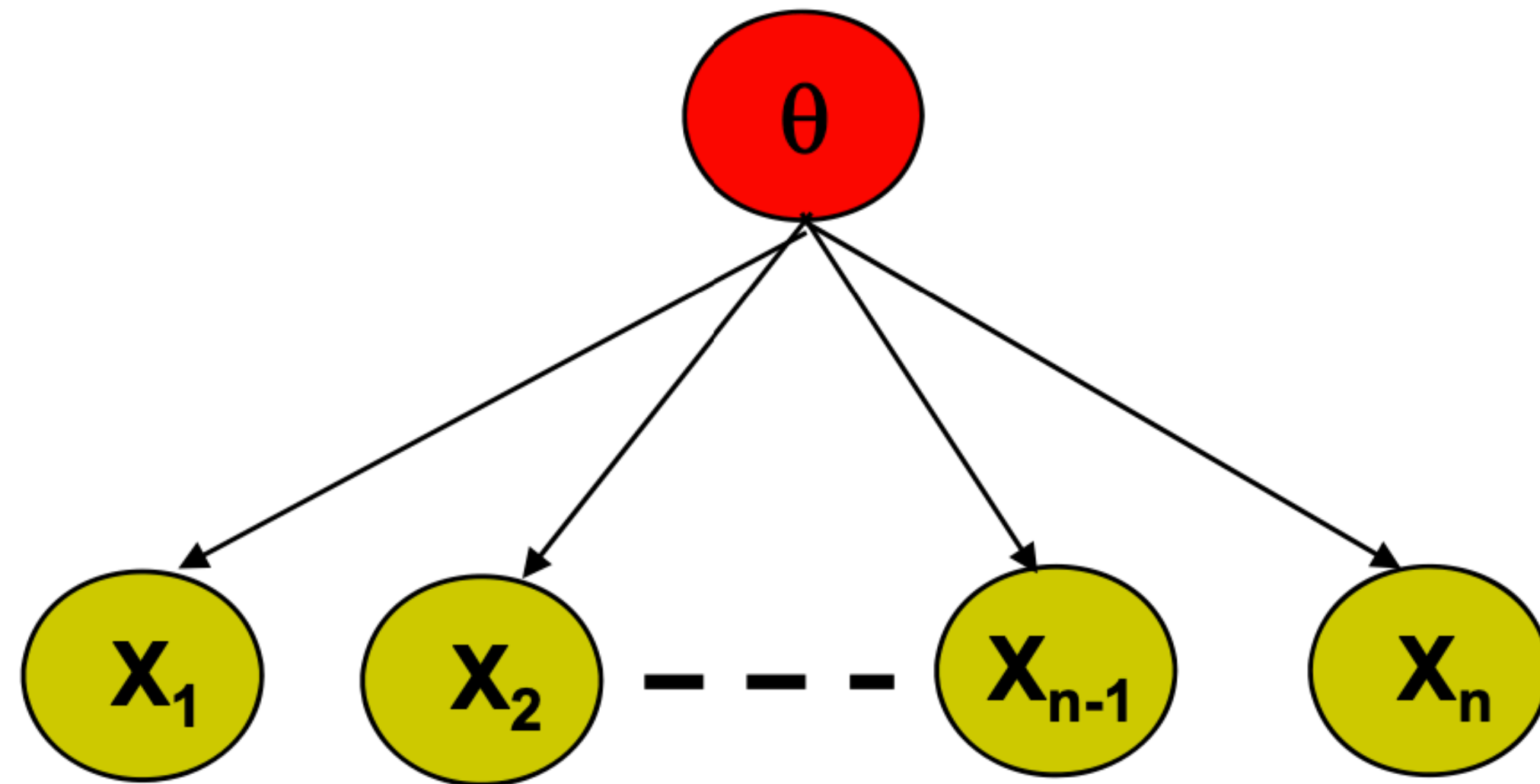
$$P(X_{1:8}) = P(X_1)P(X_2)P(X_3 | X_1 X_2)P(X_4 | X_2)P(X_5 | X_2) \\ P(X_6 | X_3, X_4)P(X_7 | X_6)P(X_8 | X_5, X_6)$$

**More formal definition:**

It refers to a family of distributions on a set of random variables that are compatible with all the probabilistic independence propositions encoded by a graph that connects these variables

**Probabilistic Graphical Model is a  
graphical language to express  
conditional independence**

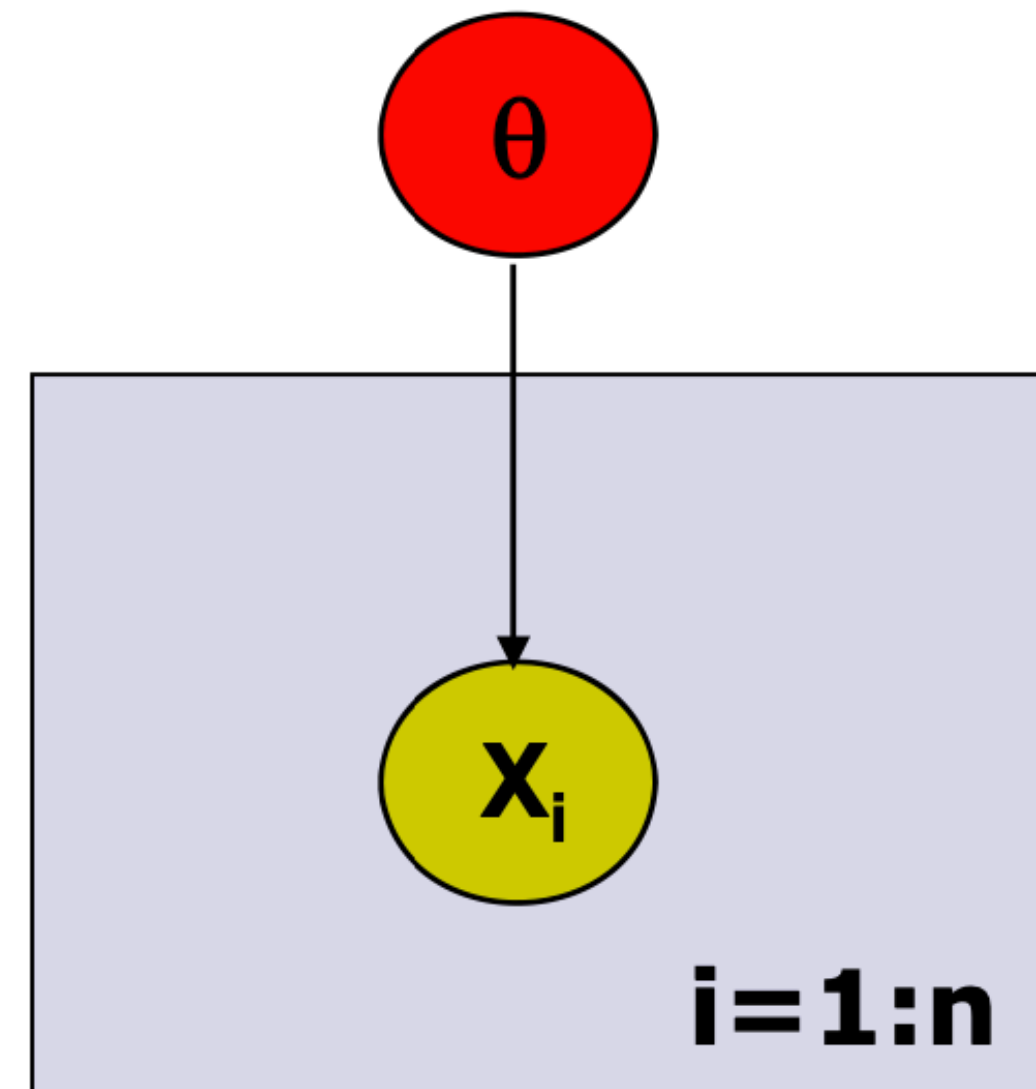
# Conditionally Independent Observations



Model parameters

Data  $\{X_1, X_2, \dots, X_n\}$

# “Plate” Notation



**Model parameters**

**Data =  $\{x_1, \dots, x_n\}$**

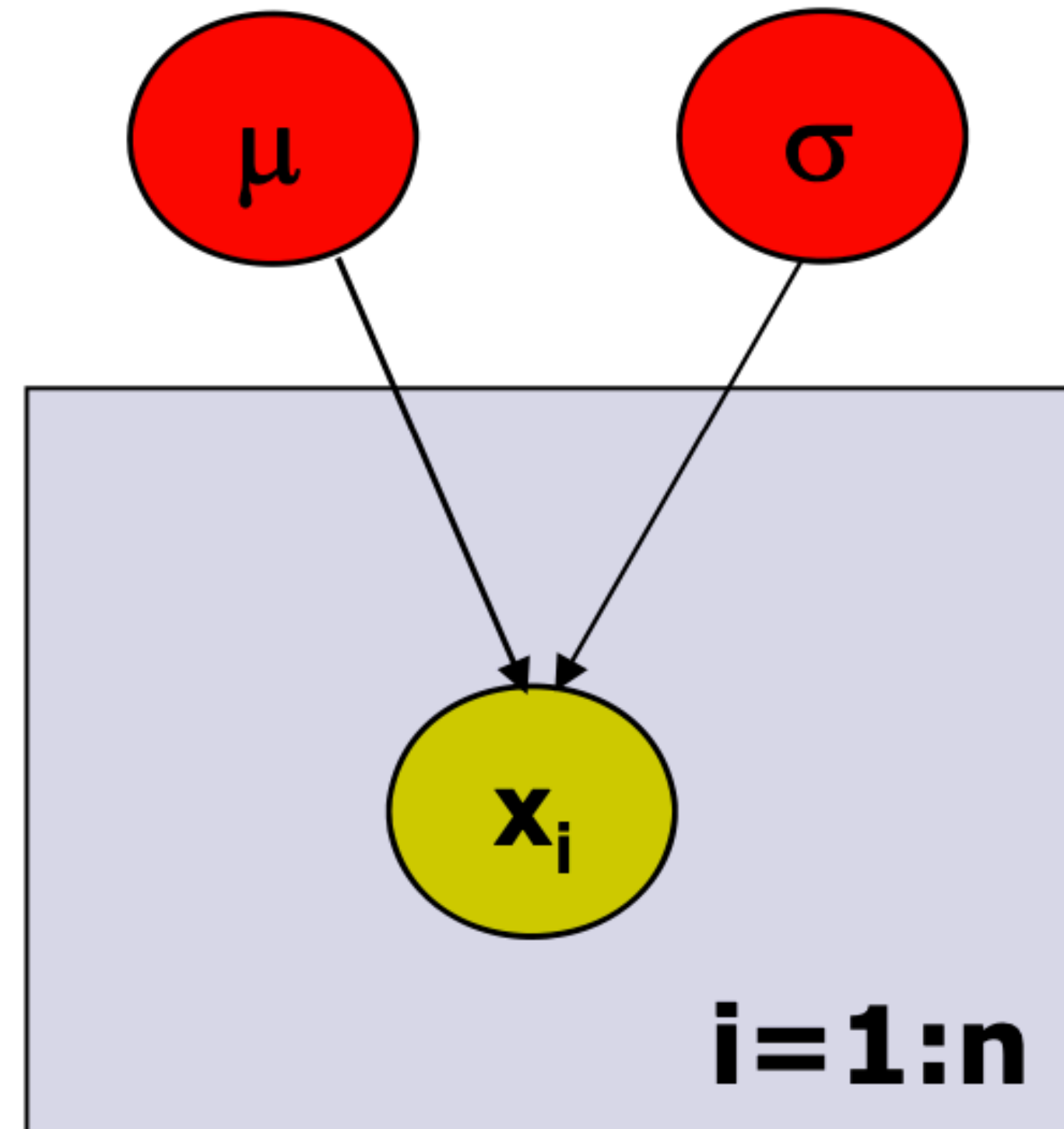
**variables within a plate are replicated  
in a conditionally independent manner**

# Example: Gaussian Model

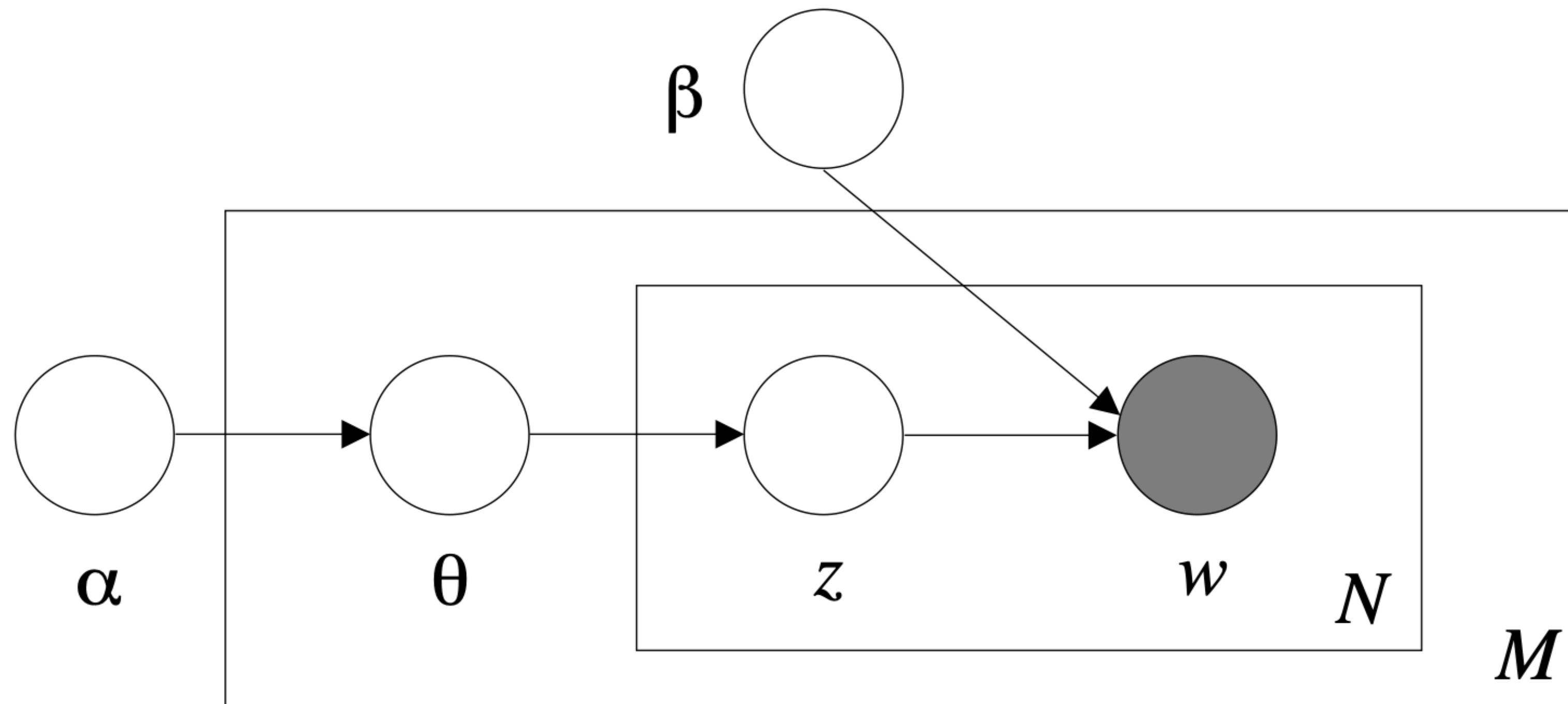
Generative model:

$$\begin{aligned} p(x_1, \dots, x_n \mid \mu, \sigma) &= \prod p(x_i \mid \mu, \sigma) \\ &= p(\text{data} \mid \text{parameters}) \\ &= p(D \mid \theta) \end{aligned}$$

where  $\theta = \{\mu, \sigma\}$



# Observed Variable and Latent Variable Notations



We typically use gray variables to denote observed variables



# Gaussian Mixture Model / Gaussian Discriminative Analysis in PGMs

# Inference and Learning

Query a node (random variable) in the graph

- Task 1: How do we answer **queries** about  $P$ ?
  - We use **inference** as a name for the process of computing answers to such queries
- Task 2: How do we estimate a **plausible model**  $M$  from data  $D$ ?
  - i. We use **learning** as a name for the process of obtaining point estimate of  $M$ .

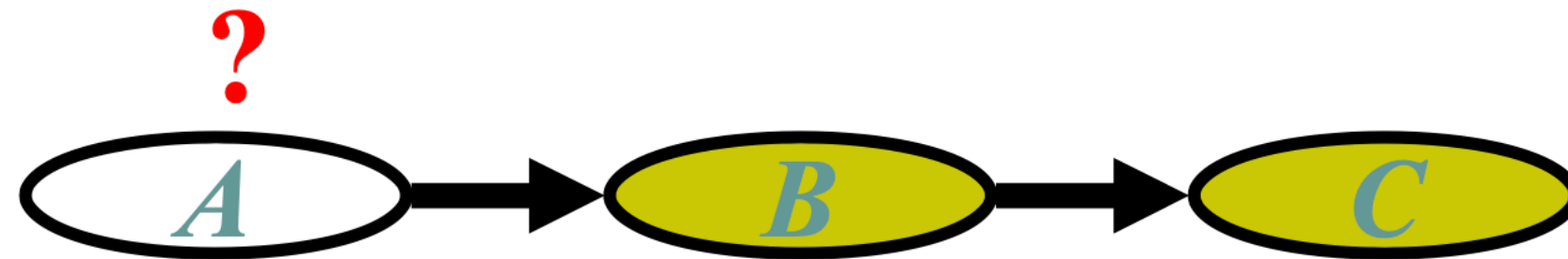
# Examples

- **Prediction:** what is the probability of an outcome given the starting condition



- the query node is a descendent of the evidence

- **Diagnosis:** what is the probability of disease/fault given symptoms



- the query node an ancestor of the evidence

In practice, the observed variable is often the data that is on the leaf nodes

# How to Learn the Parameters

1. When  $\theta$  is the parameter and does not have prior  $\rightarrow$  MLE

$$p(x, z; \theta)$$

2. When we add the prior over  $\theta \rightarrow$  MAP (Bayesian)

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

# How to do MLE on Latent Variable Models?

Expectation Maximization!

The E-step computes the posterior distribution  $p(z|x)$

This process is referred to as inference

# Approaches to Inference

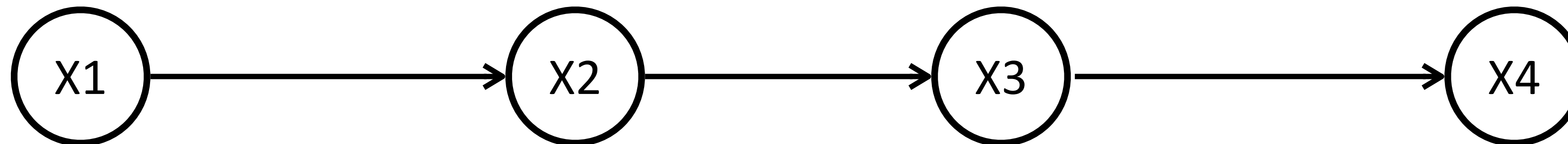
- Exact inference algorithms
    - The elimination algorithm
    - Belief propagation
    - The junction tree algorithms (but will not cover in detail here)
  - Approximate inference techniques
    - Variational algorithms
    - Stochastic simulation / sampling methods
    - Markov chain Monte Carlo methods
- Variational Autoencoders

# 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?

# Hidden Markov Models



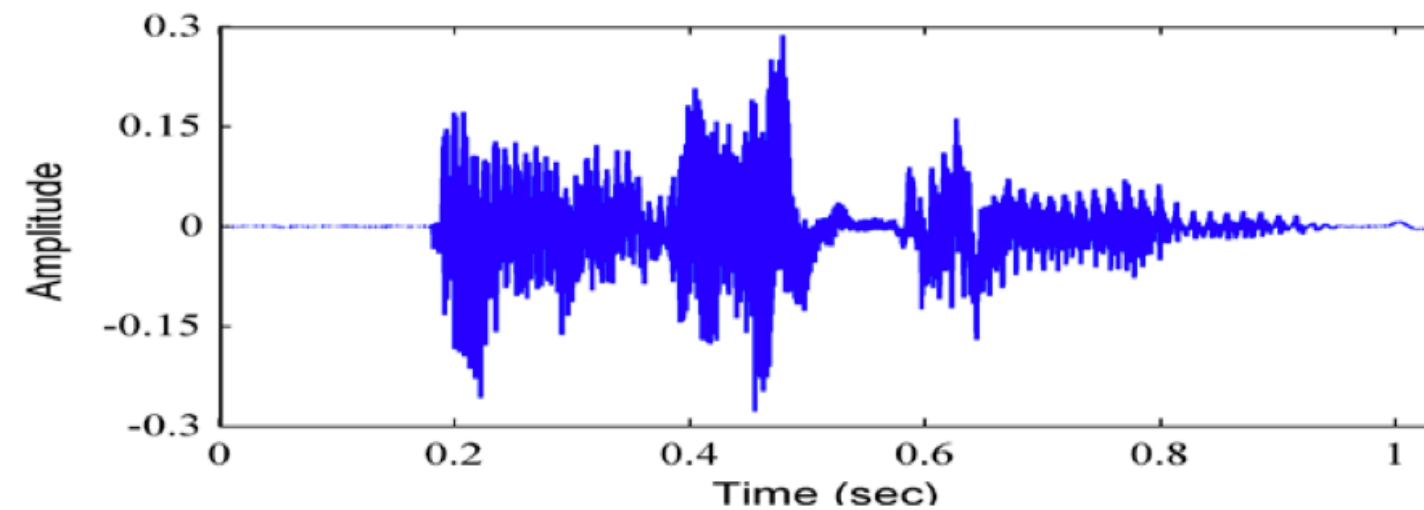
# i.i.d to sequential data

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

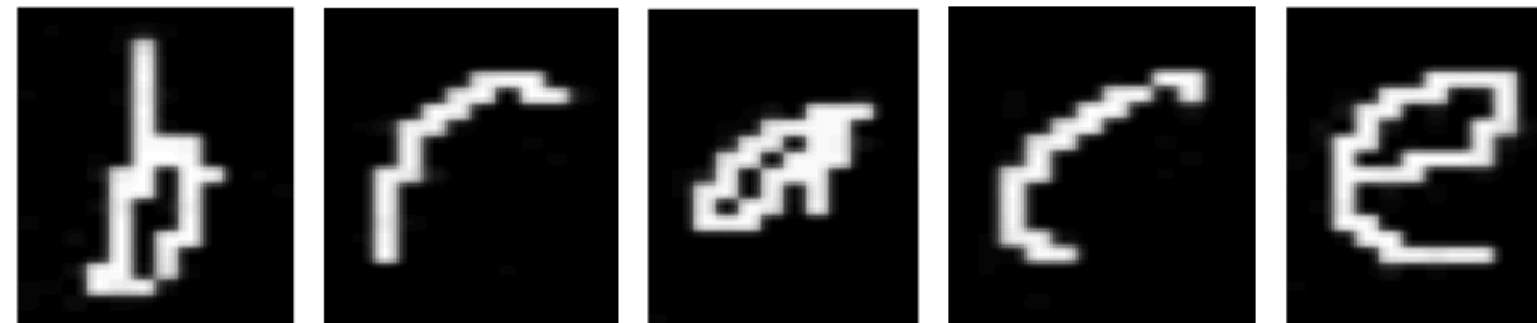
❑ Sequential (non i.i.d.) data

– Time-series data

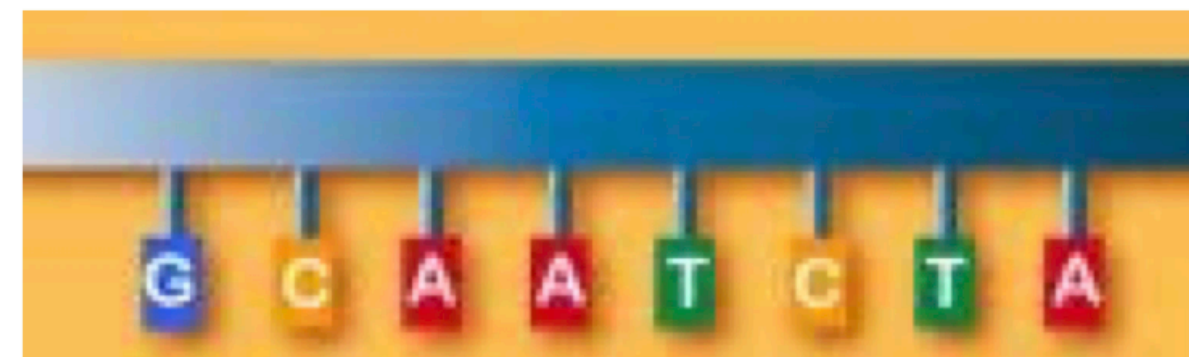
E.g. Speech



– Characters in a sentence



– Base pairs along a DNA strand



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

# Markov Models

□ Joint distribution of  $n$  arbitrary random variables

$$\begin{aligned} p(\mathbf{X}) &= p(X_1, X_2, \dots, X_n) \\ &= p(X_1)p(X_2|X_1)p(X_3|X_2, X_1) \dots p(X_n|X_{n-1}, \dots, X_1) \\ &= \prod_{i=1}^n p(X_i|X_{i-1}, \dots, X_1) \quad \text{Chain rule} \end{aligned}$$

□ Markov Assumption ( $m^{\text{th}}$  order)

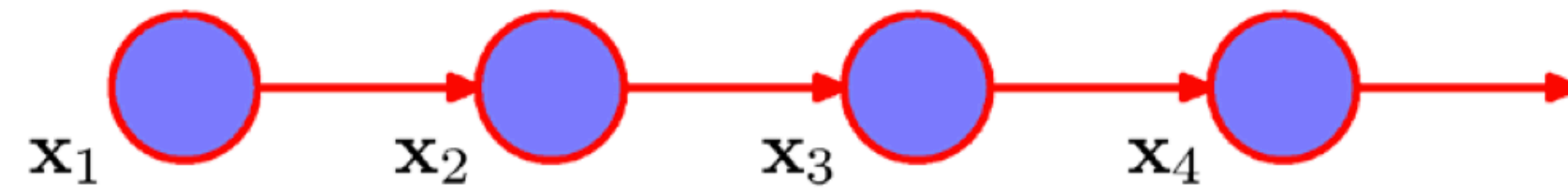
$$p(\mathbf{X}) = \prod_{i=1}^n p(X_i|X_{i-1}, \dots, X_{i-m})$$

Current observation only depends on past  $m$  observations

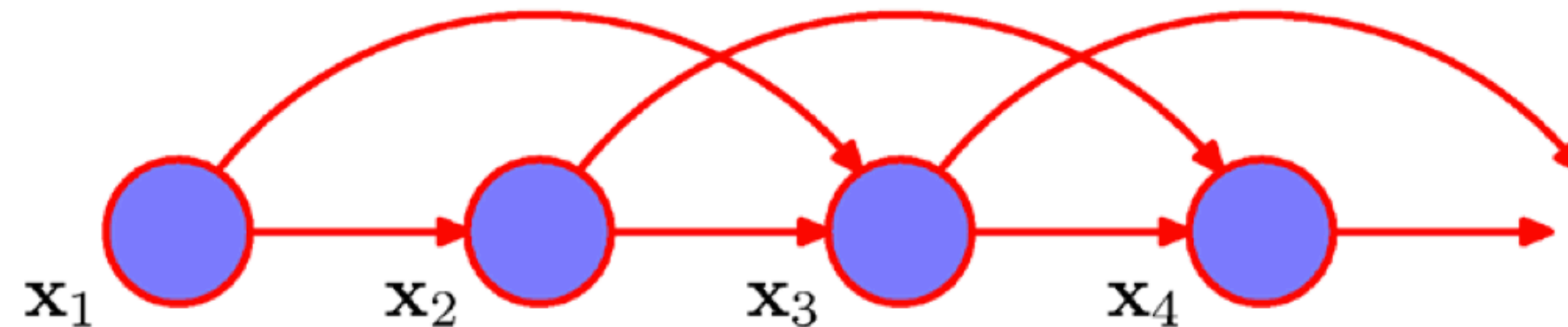
# Markov Models

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



# Markov Models

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

## □ Markov Assumption

# parameters in  
stationary model  
K-ary variables

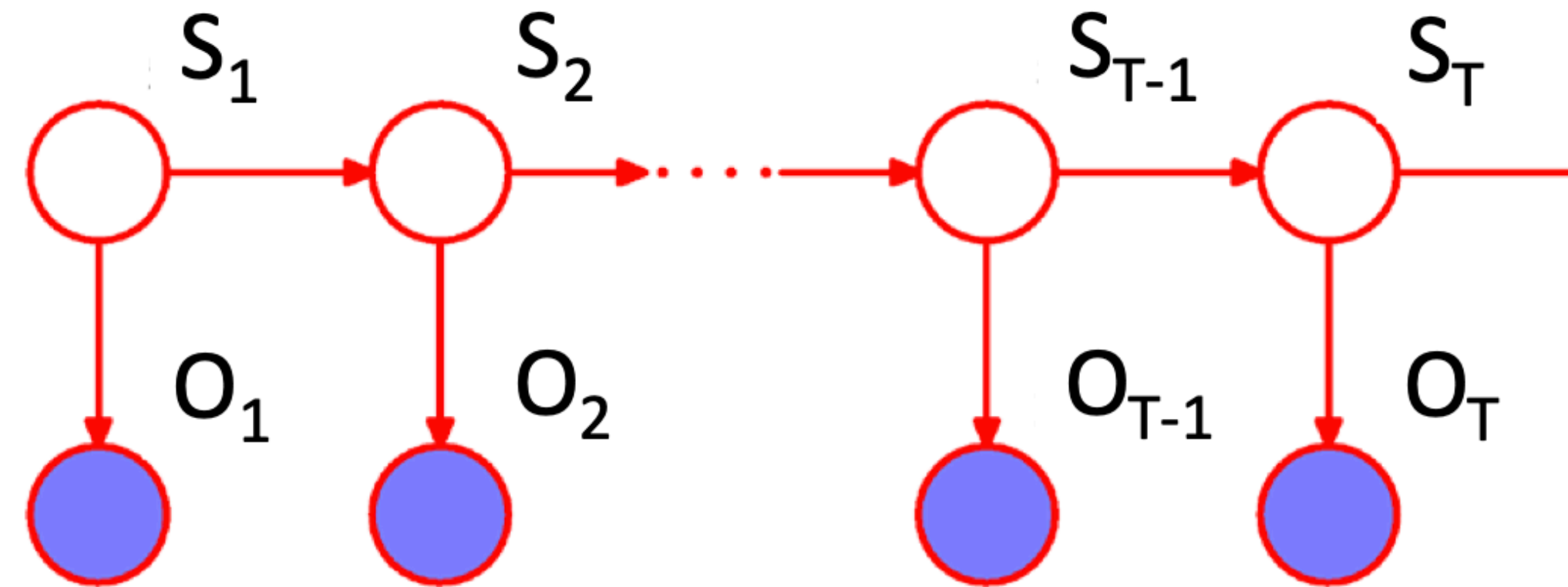
1<sup>st</sup> order  $p(\mathbf{X}) = \prod_{i=1}^n p(X_i | X_{i-1})$   $O(K^2)$

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

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

# Hidden Markov Models



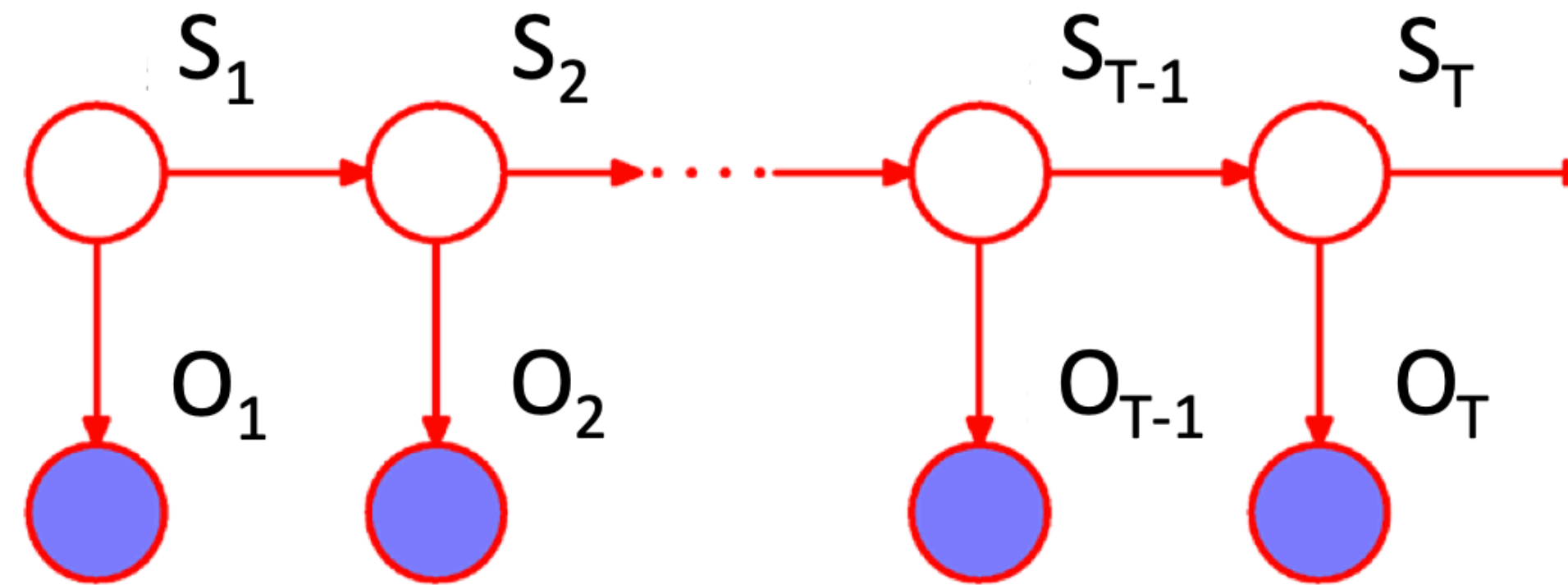
Observation space

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

Hidden states

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

# Hidden Markov Models



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

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?

# Hidden Markov Models

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

Initial probabilities

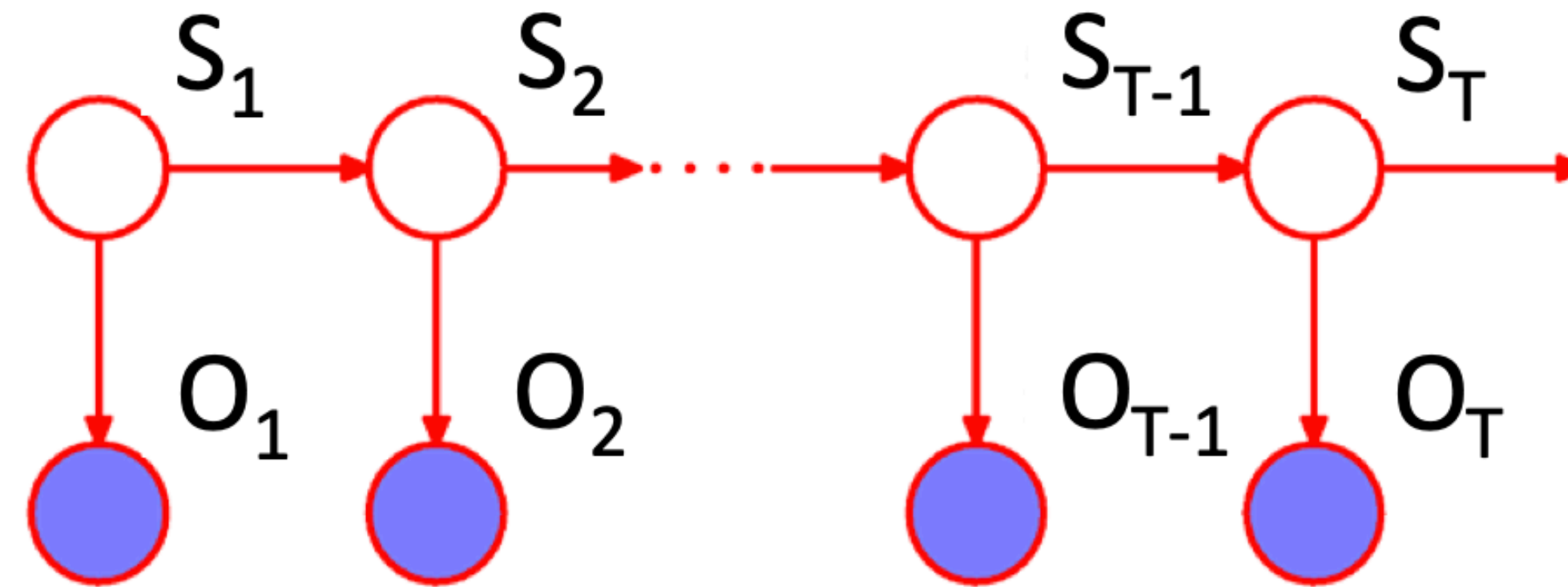
$$p(S_1 = i) = \pi_i$$

Transition probabilities

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

Emission probabilities

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



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

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

# 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





# 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

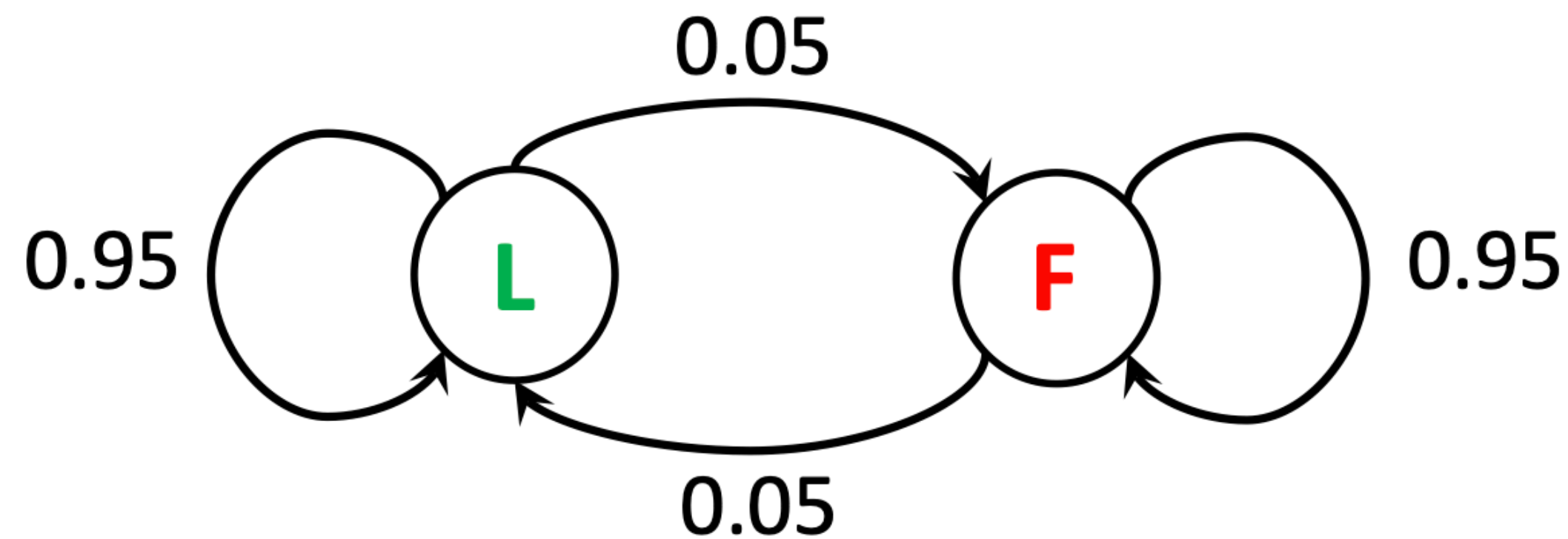
This is the decoding question in HMMs

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

# State Space Representation

- Switch between **F** and **L** with 5% probability



## HMM Parameters

Initial probs

$$P(S_1 = \mathbf{L}) = 0.5 = P(S_1 = \mathbf{F})$$

Transition probs

$$P(S_t = \mathbf{L}/\mathbf{F} | S_{t-1} = \mathbf{L}/\mathbf{F}) = 0.95$$

$$P(S_t = \mathbf{F}/\mathbf{L} | S_{t-1} = \mathbf{L}/\mathbf{F}) = 0.05$$

Emission probabilities

$$P(O_t = y | S_t = \mathbf{F}) = 1/6 \quad y = 1, 2, 3, 4, 5, 6$$

$$P(O_t = y | S_t = \mathbf{L}) = 1/10 \quad y = 1, 2, 3, 4, 5$$

$$= 1/2 \quad y = 6$$

# 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

- **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)**

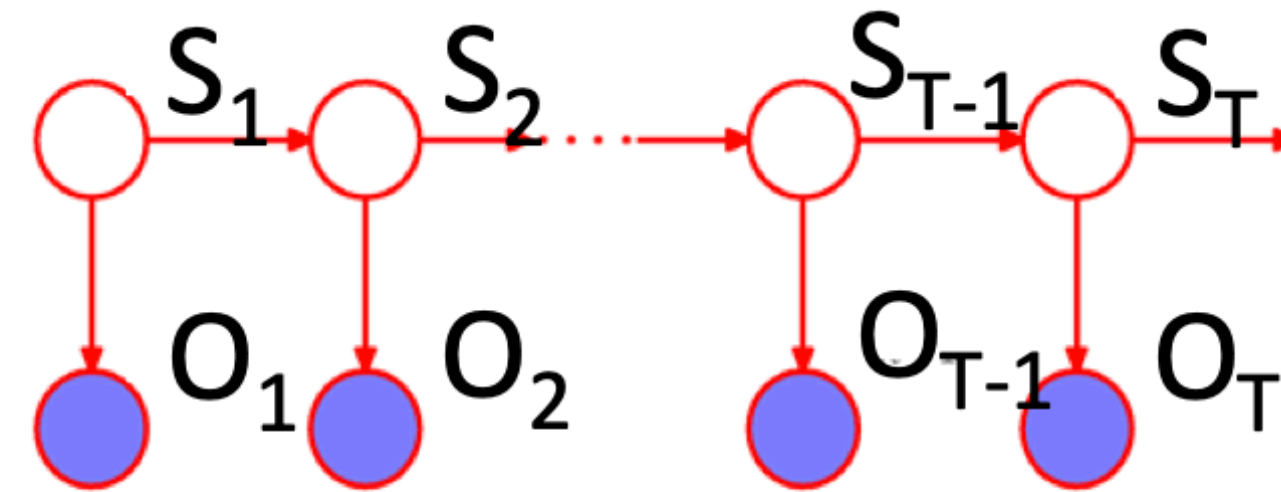
# Evaluation Problem

- 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

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

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



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

# Forward Probability

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

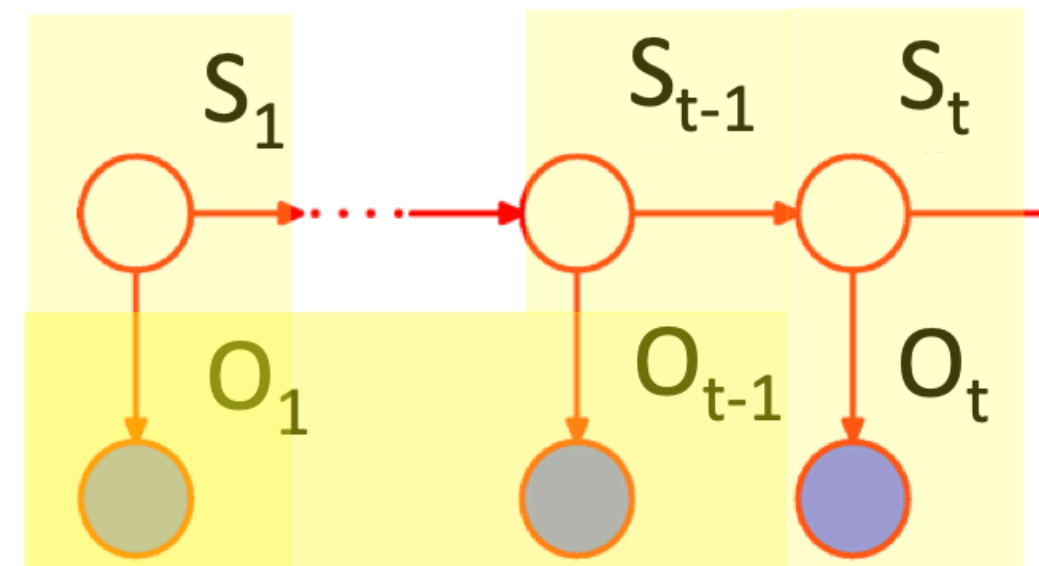
Compute forward probability  $\alpha_t^k$  recursively over  $t$

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

Introduce  $S_{t-1}$

Chain rule

Markov assumption



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

# 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?

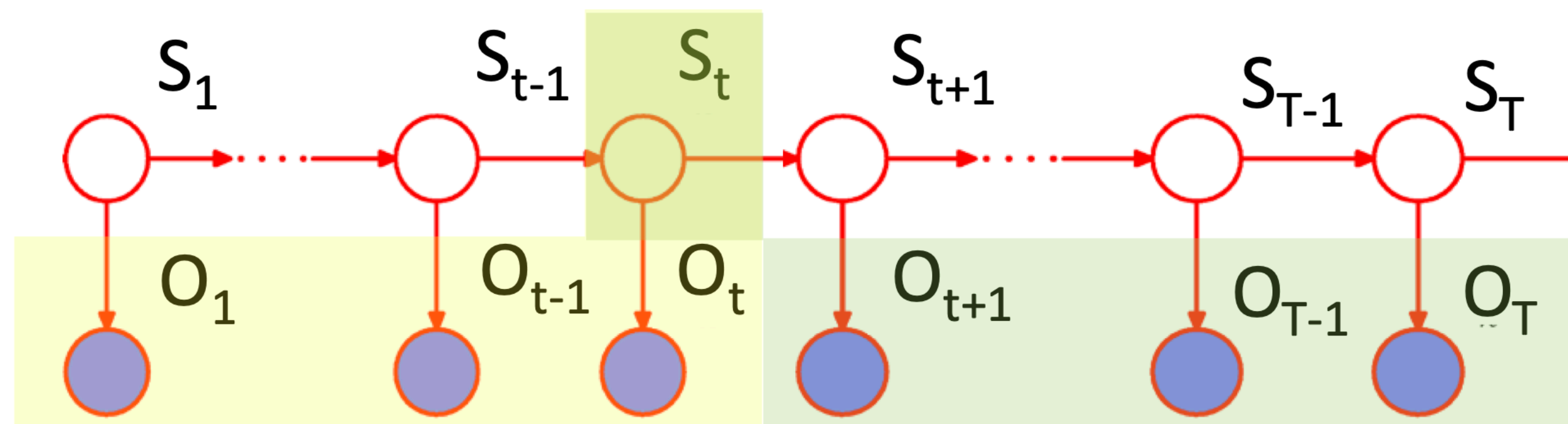
# Decoding Problem 1

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

$$\begin{aligned}
 p(S_t = k, \{O_t\}_{t=1}^T) &= p(O_1, \dots, O_t, S_t = k, O_{t+1}, \dots, O_T) \\
 &= \underbrace{p(O_1, \dots, O_t, S_t = k)}_{\alpha_t^k} \underbrace{p(O_{t+1}, \dots, O_T | S_t = k)}_{\beta_t^k}
 \end{aligned}$$

Compute recursively





# Forward-Backward Algorithm

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

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

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

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

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

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

# 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?

# 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

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

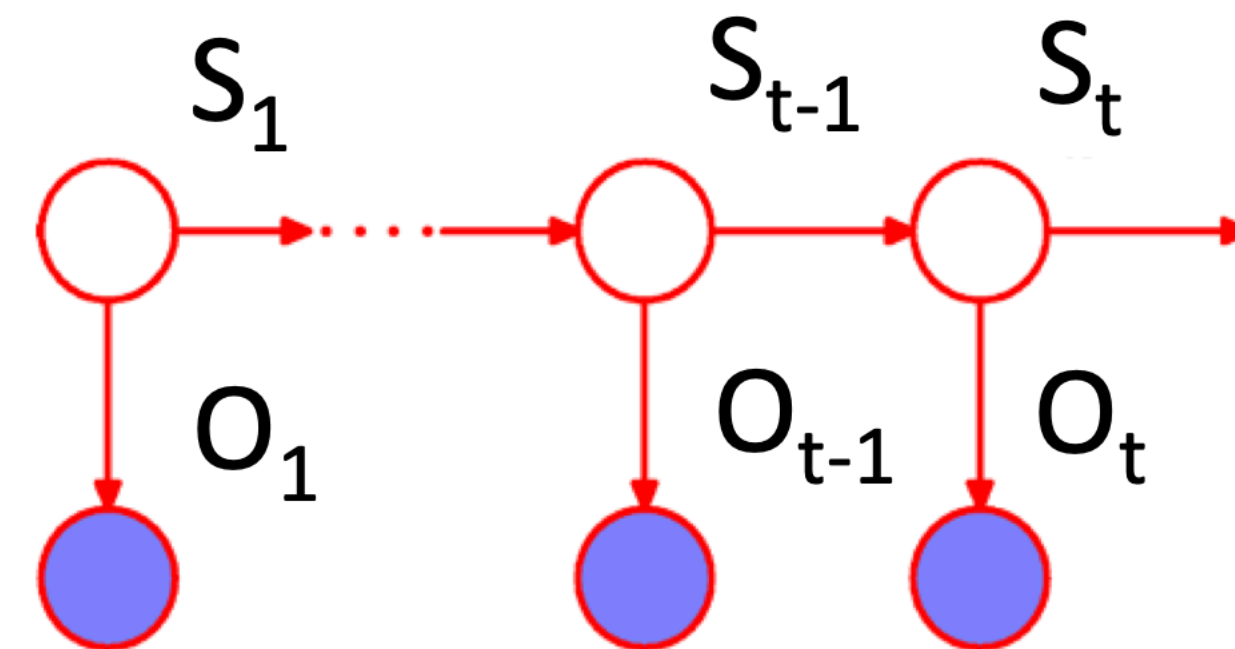
Compute probability  $V_t^k$  recursively over t

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

·  
·  
·

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

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

- 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

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

You will derive the EM  
in your HW

**Thank You!**  
**Q & A**