# Hidden Markov Models

# +

# Exam 2 Review

Matt Gormley
Lecture 18
Mar. 25, 2022

1

# Reminders

- **Homework 6: Learning Theory / Generative Models**
  - **Out: Fri, Mar. 18**
  - **Due: Fri, Mar. 25 at 11:59pm**
  - **IMPORTANT: only 2 grace/late days permitted**
- **Exam 2 (Thu, Mar 3rd)**
  - **Thu, Mar. 31, 6:30pm – 8:30pm**
- **Practice for Exam 2**
  - **Practice problems released on course website**
    - **Out: Fri, Mar. 25**
  - **Mock Exam 2**
    - **Out: Fri, Mar. 25**
    - **Due Wed, Mar. 30 at 11:59pm**

# EXAM 2 LOGISTICS

# Exam 2

- **Time / Location**
  - **Time:** Thu, Mar. 31, 6:30pm – 8:30pm
  - **Location & Seats:** You have all been split across multiple rooms. Everyone has an assigned seat in one of these room. Please watch Piazza carefully for announcements.
- **Logistics**
  - Covered material: Lecture 8 – Lecture 17
  - Format of questions:
    - Multiple choice
    - True / False (with justification)
    - Derivations
    - Short answers
    - Interpreting figures
    - Implementing algorithms on paper
  - No electronic devices
  - You are allowed to **bring** one 8½ x 11 sheet of notes (front and back, handwritten with pen/pencil or tablet)

# Topics for Exam 1

- Foundations
  - Probability, Linear Algebra, Geometry, Calculus
  - Optimization
- Important Concepts
  - Overfitting
  - Experimental Design

- Classification
  - Decision Tree
  - KNN
  - Perceptron
- Regression
  - Linear Regression

# Topics for Exam 2

- **Classification**
  - Binary Logistic Regression
- **Important Concepts**
  - Stochastic Gradient Descent
  - Regularization
  - Feature Engineering
- **Feature Learning**
  - Neural Networks
  - Basic NN Architectures
  - Backpropagation

- **Learning Theory**
  - PAC Learning
- **Generative Models**
  - Generative vs. Discriminative
  - MLE / MAP
  - Naïve Bayes

- **Regression**
  - Linear Regression

# SAMPLE QUESTIONS

# Sample Questions

## 3.2 Logistic regression

Given a training set $\{(x_i, y_i), i = 1, \ldots, n\}$ where $x_i \in \mathbb{R}^d$ is a feature vector and $y_i \in \{0, 1\}$ is a binary label, we want to find the parameters $\hat{w}$ that maximize the likelihood for the training set, assuming a parametric model of the form

$$p(y = 1|x; w) = \frac{1}{1 + \exp(-w^T x)}.$$

The conditional log likelihood of the training set is

$$\ell(w) = \sum_{i=1}^{n} y_i \log p(y_i, |x_i; w) + (1 - y_i) \log(1 - p(y_i, |x_i; w)),$$

and the gradient is

$$\nabla \ell(w) = \sum_{i=1}^{n} (y_i - p(y_i|x_i; w)) x_i.$$

(b) [5 pts.] What is the form of the classifier output by logistic regression?

(c) [2 pts.] **Extra Credit:** Consider the case with binary features, i.e., $x \in \{0, 1\}^d \subset \mathbb{R}^d$, where feature $x_1$ is rare and happens to appear in the training set with only label 1. What is $\hat{w}_1$? Is the gradient ever zero for any finite $w$? Why is it important to include a regularization term to control the norm of $\hat{w}$?

# Samples Questions

## 2.1 Train and test errors

In this problem, we will see how you can debug a classifier by looking at its train and test errors. Consider a classifier trained till convergence on some training data $\mathcal{D}^{\text{train}}$, and tested on a separate test set $\mathcal{D}^{\text{test}}$. You look at the test error, and find that it is very high. You then compute the training error and find that it is close to 0.
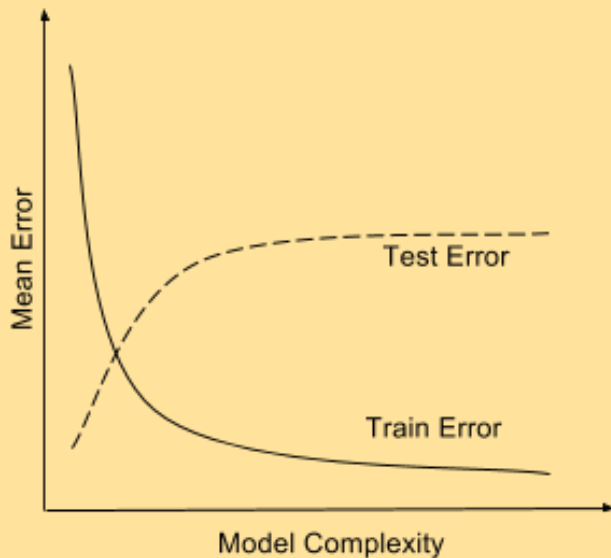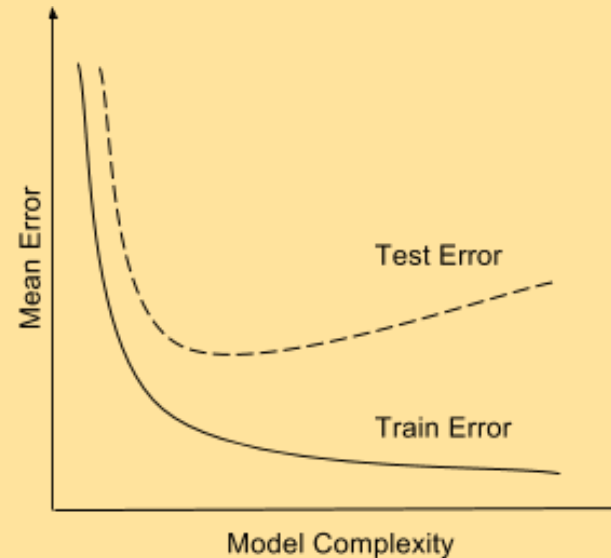
1. **[4 pts]** Which of the following is expected to help? Select all that apply.

   (a) Increase the training data size.

   (b) Decrease the training data size.

   (c) Increase model complexity (For example, if your classifier is an SVM, use a more complex kernel. Or if it is a decision tree, increase the depth).

   (d) Decrease model complexity.

   (e) Train on a combination of $\mathcal{D}^{\text{train}}$ and $\mathcal{D}^{\text{test}}$ and test on $\mathcal{D}^{\text{test}}$

   (f) Conclude that Machine Learning does not work.

# Samples Questions

## 2.1 Train and test errors

In this problem, we will see how you can debug a classifier by looking at its train and test errors. Consider a classifier trained till convergence on some training data $\mathcal{D}^{\text{train}}$, and tested on a separate test set $\mathcal{D}^{\text{test}}$. You look at the test error, and find that it is very high. You then compute the training error and find that it is close to 0.

4. **[1 pts]** Say you plot the train and test errors as a function of the model complexity. Which of the following two plots is your plot expected to look like?



(a)

(b)

# Sample Questions

## 5 Learning Theory [20 pts.]

(a) [3 pts.] **T or F**: It is possible to label 4 points in $\mathbb{R}^2$ in all possible $2^4$ ways via linear separators in $\mathbb{R}^2$.

(d) [3 pts.] **T or F**: The VC dimension of a hypothesis space with infinite size is also infinite.

# Sample Questions

**Neural Networks**

Can the neural network in Figure (b) correctly classify the dataset given in Figure (a)?



(a) The dataset with groups $S_1$, $S_2$, and $S_3$.

(b) The neural network architecture

# Sample Questions

**Neural Networks**

Apply the backpropagation algorithm to obtain the partial derivative of the mean-squared error of y with the true value y* with respect to the weight $w_{22}$ assuming a sigmoid nonlinear activation function for the hidden layer.



(b) The neural network architecture

# Sample Questions

## 1.2 Maximum Likelihood Estimation (MLE)

Assume we have a random sample that is Bernoulli distributed $X_1, \ldots, X_n \sim \text{Bernoulli}(\theta)$. We are going to derive the MLE for $\theta$. Recall that a Bernoulli random variable $X$ takes values in $\{0, 1\}$ and has probability mass function given by

$$P(X; \theta) = \theta^X (1 - \theta)^{1-X}.$$

(a) [2 pts.] Derive the likelihood, $L(\theta; X_1, \ldots, X_n)$.

(c) **Extra Credit:** [2 pts.] Derive the following formula for the MLE: $\hat{\theta} = \dfrac{1}{n} \left( \sum_{i=1}^{n} X_i \right)$.

# Sample Questions

## 1.3 MAP vs MLE

Answer each question with **T** or **F** and **provide a one sentence explanation of your answer**:

(a) [2 pts.] **T or F:** In the limit, as $n$ (the number of samples) increases, the MAP and MLE estimates become the same.

# Sample Questions

## 1.1 Naive Bayes

You are given a data set of 10,000 students with their sex, height, and hair color. You are trying to build a classifier to predict the sex of a student, so you randomly split the data into a training set and a testing set. Here are the specifications of the data set:

- sex $\in$ {male,female}

- height $\in$ [0,300] centimeters

- hair $\in$ {brown, black, blond, red, green}

- 3240 men in the data set

- 6760 women in the data set

Under the assumptions necessary for Naive Bayes (not the distributional assumptions you might naturally or intuitively make about the dataset) answer each question with **T** or **F** and **provide a one sentence explanation of your answer**:

(a) [2 pts.] **T or F:** As height is a continuous valued variable, Naive Bayes is not appropriate since it cannot handle continuous valued variables.

(c) [2 pts.] **T or F:** $P(\texttt{height}|\texttt{sex}, \texttt{hair}) = P(\texttt{height}|\texttt{sex})$.

# THE BIG PICTURE

# ML Big Picture

## Learning Paradigms:

*What data is available and when? What form of prediction?*

- supervised learning
- unsupervised learning
- semi-supervised learning
- reinforcement learning
- active learning
- imitation learning
- domain adaptation
- online learning
- density estimation
- recommender systems
- feature learning
- manifold learning
- dimensionality reduction
- ensemble learning
- distant supervision
- hyperparameter optimization

## Theoretical Foundations:

*What principles guide learning?*

- ❏ probabilistic
- ❏ information theoretic
- ❏ evolutionary search
- ❏ ML as optimization

## Problem Formulation:

*What is the structure of our output prediction?*

| | |
|---|---|
| boolean | Binary Classification |
| categorical | Multiclass Classification |
| ordinal | Ordinal Classification |
| real | Regression |
| ordering | Ranking |
| multiple discrete | Structured Prediction |
| multiple continuous | (e.g. dynamical systems) |
| both discrete & cont. | (e.g. mixed graphical models) |

## Facets of Building ML Systems:

*How to build systems that are robust, efficient, adaptive, effective?*

1. Data prep
2. Model selection
3. Training (optimization / search)
4. Hyperparameter tuning on validation data
5. (Blind) Assessment on test data

## Application Areas

*Key challenges?*
NLP, Speech, Computer Vision, Robotics, Medicine, Search

## Big Ideas in ML:

*Which are the ideas driving development of the field?*

- inductive bias
- generalization / overfitting
- bias-variance decomposition
- generative vs. discriminative
- deep nets, graphical models
- PAC learning
- distant rewards

# ML Big Picture

*Whiteboard*

- Decision Rules / Models
- Objective Functions
- Regularization
- Optimization

# MOTIVATION: STRUCTURED PREDICTION

# Structured Prediction

- Most of the models we've seen so far were for **classification**
  - Given observations: $\qquad \boldsymbol{x} = (x_1,\ x_2,\ ...,\ x_K)$
  - Predict a (binary) **label:** $\quad y$
- Many real-world problems require **structured prediction**
  - Given observations: $\qquad \boldsymbol{x} = (x_1,\ x_2,\ ...,\ x_K)$
  - Predict a **structure:** $\qquad \boldsymbol{y} = (y_1,\ y_2,\ ...,\ y_J)$
- Some *classification* problems benefit from **latent structure**

# Structured Prediction Examples

- **Examples of structured prediction**
  - Part-of-speech (POS) tagging
  - Handwriting recognition
  - Speech recognition
  - Word alignment
  - Congressional voting
- **Examples of latent structure**
  - Object recognition

# Dataset for Supervised
# Part-of-Speech (POS) Tagging

Data: $\mathcal{D} = \{\boldsymbol{x}^{(n)}, \boldsymbol{y}^{(n)}\}_{n=1}^{N}$

# Dataset for Supervised Handwriting Recognition

Data: $\mathcal{D} = \{\boldsymbol{x}^{(n)}, \boldsymbol{y}^{(n)}\}_{n=1}^{N}$



Sample 1: $\boldsymbol{y}^{(1)}$, $\boldsymbol{x}^{(1)}$

Sample 2: $\boldsymbol{y}^{(2)}$, $\boldsymbol{x}^{(2)}$

Sample 2: $\boldsymbol{y}^{(3)}$, $\boldsymbol{x}^{(3)}$

# Dataset for Supervised Phoneme (Speech) Recognition

Data: $\mathcal{D} = \{\boldsymbol{x}^{(n)}, \boldsymbol{y}^{(n)}\}_{n=1}^{N}$



Sample 1: h# dh ih s w uh z iy z iy $\quad\}\; y^{(1)}$

$\quad\}\; x^{(1)}$

Sample 2: f ao r ah s s h# $\quad\}\; y^{(2)}$

$\quad\}\; x^{(2)}$

Figures from (Jansen & Niyogi, 2013)

# Word Alignment / Phrase Extraction

- **Variables (boolean):**
  - For each (Chinese phrase, English phrase) pair, are they linked?

- **Interactions:**
  - Word fertilities
  - Few "jumps" (discontinuities)
  - Syntactic reorderings
  - "ITG contraint" on alignment
  - Phrases are disjoint (?)

**(Burkett & Klein, 2012)**

# Congressional Voting



- **Variables:**
  - Representative's vote
  - **Text of all speeches of a representative**
  - Local contexts of references between two representatives

- **Interactions:**
  - Words used by representative and their vote
  - Pairs of representatives and their local context

(Stoyanov & Eisner, 2012)

# Structured Prediction Examples

- **Examples of structured prediction**
  - Part-of-speech (POS) tagging
  - Handwriting recognition
  - Speech recognition
  - Word alignment
  - Congressional voting
- **Examples of latent structure**
  - Object recognition

# Case Study: Object Recognition

Data consists of images $x$ and labels $y$.



pigeon — $x^{(1)}$, $y^{(1)}$

rhinoceros — $x^{(2)}$, $y^{(2)}$

leopard — $x^{(3)}$, $y^{(3)}$

llama — $x^{(4)}$, $y^{(4)}$

# Case Study: Object Recognition

## Data consists of images $x$ and labels $y$.

- Preprocess data into "patches"

- Posit a latent labeling $z$ describing the object's parts (e.g. head, leg, tail, torso, grass)

- Define graphical model with these latent variables in mind

- $z$ is not observed at train or test time



leopard

# Case Study: Object Recognition

## Data consists of images $x$ and labels $y$.

- Preprocess data into "patches"

- Posit a latent labeling $z$ describing the object's parts (e.g. head, leg, tail, torso, grass)

- Define graphical model with these latent variables in mind

- $z$ is not observed at train or test time



leopard $Y$

# Case Study: Object Recognition

Data consists of images $x$ and labels $y$.

- Preprocess data into "patches"

- Posit a latent labeling $z$ describing the object's parts (e.g. head, leg, tail, torso, grass)

- Define graphical model with these latent variables in mind

- $z$ is not observed at train or test time



leopard

# Structured Prediction

## Preview of challenges to come...

- Consider the task of finding the **most probable assignment** to the output

| Classification | Structured Prediction |
|---|---|
| $\hat{y} = \underset{y}{\operatorname{argmax}}\, p(y\|\mathbf{x})$ | $\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}}\, p(\mathbf{y}\|\mathbf{x})$ |
| where $y \in \{+1, -1\}$ | where $\mathbf{y} \in \mathcal{Y}$ |
| | and $\|\mathcal{Y}\|$ is very large |

# Machine Learning

The **data** inspires the structures we want to predict

Our **model** defines a score for each structure

It also tells us what to optimize

**Inference** finds {best structure, marginals, partition function} for a new observation

(**Inference** is usually called as a subroutine in learning)

**Learning** tunes the parameters of the model



Domain Knowledge

Mathematical Modeling

ML

Combinatorial Optimization

Optimization

# Machine Learning

**Data**



**Model**



**Objective**



**Inference**



(**Inference** is usually called as a subroutine in learning)

**Learning**



38

# BACKGROUND

# Background: Chain Rule of Probability

For random variables $A$ and $B$:

$$P(A, B) = P(A|B)P(B)$$

For random variables $X_1, X_2, X_3, X_4$:

$$P(X_1, X_2, X_3, X_4) = P(X_1|X_2, X_3, X_4)$$
$$P(X_2|X_3, X_4)$$
$$P(X_3|X_4)$$
$$P(X_4)$$

# Background:
# Conditional Independence

Random variables $A$ and $B$ are conditionally independent given $C$ if:

$$P(A, B|C) = P(A|C)P(B|C) \qquad (1)$$

or equivalently:

$$P(A|B, C) = P(A|C) \qquad (2)$$

We write this as:

$$A \perp\!\!\!\perp B | C$$

Later we will also write: $I<A, \{C\}, B>$

# HIDDEN MARKOV MODEL (HMM)

# From Mixture Model to HMM



**"Naïve Bayes":**

$$P(\mathbf{X}, \mathbf{Y}) = \prod_{t=1}^{T} P(X_t | Y_t) p(Y_t)$$

**HMM:**

$$P(\mathbf{X}, \mathbf{Y}) = P(Y_1) \left( \prod_{t=1}^{T} P(X_t | Y_t) \right) \left( \prod_{t=2}^{T} p(Y_t | Y_{t-1}) \right)$$

44

# Markov Models

*Whiteboard*

- Example: Tunnel Closures
  [courtesy of Roni Rosenfeld]

- First-order Markov assumption

- Conditional independence assumptions

# Mixture Model for Time Series Data

We could treat each (tunnel state, travel time) pair as independent. This corresponds to a Naïve Bayes model with a single feature (travel time).

$$p(\text{O, S, S, O, C, 2m, 3m, 18m, 9m, 27m}) = (.8 * .2 * .1 * .03 * \dots)$$

# Hidden Markov Model

A Hidden Markov Model (HMM) provides a joint distribution over the the tunnel states / travel times with an assumption of dependence between adjacent tunnel states.

$$p(\text{O}, \text{S}, \text{S}, \text{O}, \text{C}, 2\text{m}, 3\text{m}, 18\text{m}, 9\text{m}, 27\text{m}) = (.8 * .08 * .2 * .7 * .03 * \ldots)$$

|   |    |
|---|----|
| O | .8 |
| S | .1 |
| C | .1 |

|   | O  | S   | C   |
|---|----|-----|-----|
| O | .9 | .08 | .02 |
| S | .2 | .7  | .1  |
| C | .9 | 0   | .1  |

|   | O  | S   | C   |
|---|----|-----|-----|
| O | .9 | .08 | .02 |
| S | .2 | .7  | .1  |
| C | .9 | 0   | .1  |

O → S → S → O → C

2m   3m   9m   27m

|   | 1min | 2min | 3min | … |
|---|------|------|------|---|
| O | .1   | .2   | .3   |   |
| S | .01  | .02  | .03  |   |
| C | 0    | 0    | 0    |   |

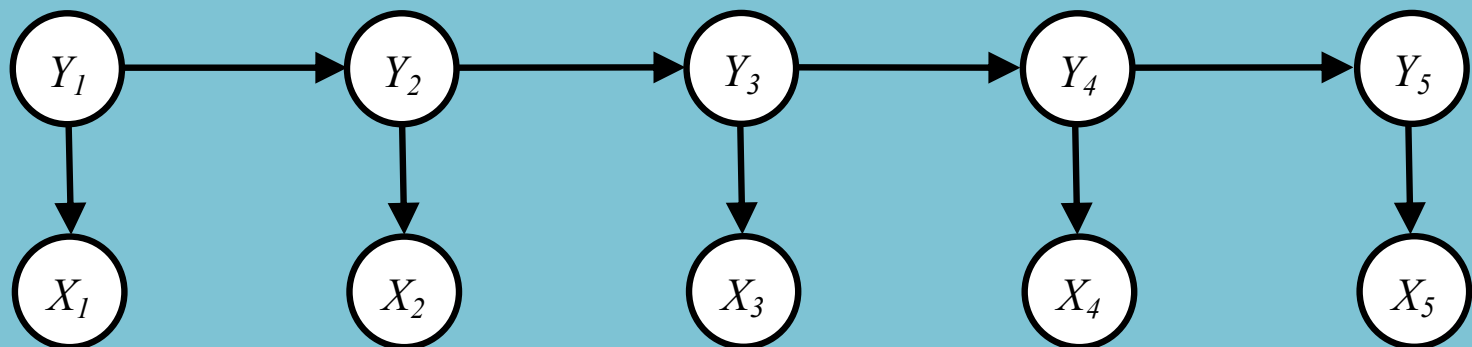|   | 1min | 2min | 3min | … |
|---|------|------|------|---|
| O | .1   | .2   | .3   |   |
| S | .01  | .02  | .03  |   |
| C | 0    | 0    | 0    |   |

# From Mixture Model to HMM



"Naïve Bayes":
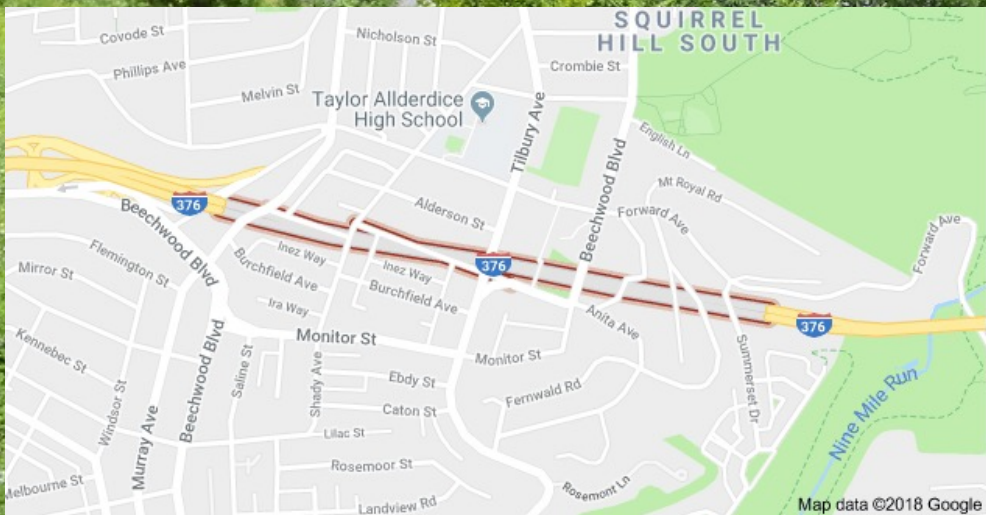
$$P(\mathbf{X}, \mathbf{Y}) = \prod_{t=1}^{T} P(X_t|Y_t)p(Y_t)$$

HMM:

$$P(\mathbf{X}, \mathbf{Y}) = P(Y_1) \left( \prod_{t=1}^{T} P(X_t|Y_t) \right) \left( \prod_{t=2}^{T} p(Y_t|Y_{t-1}) \right)$$
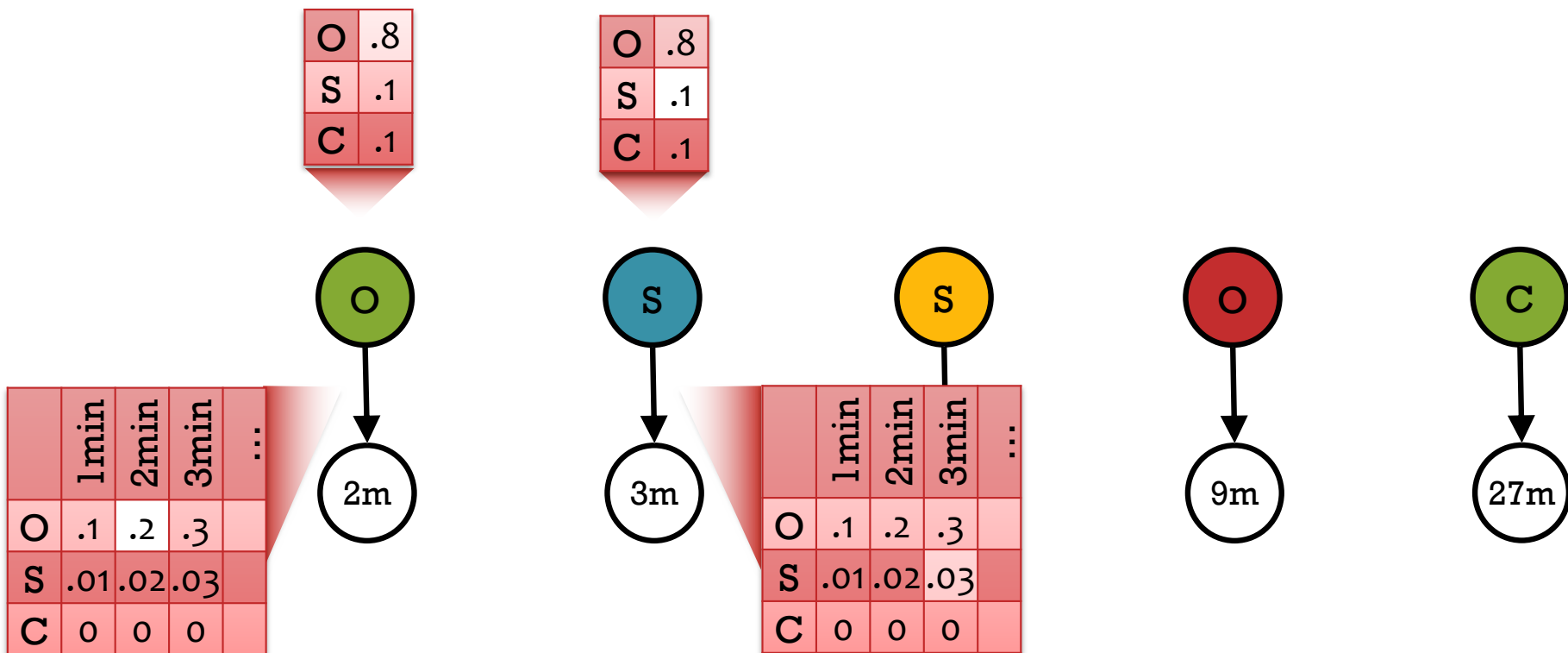
55

# From Mixture Model to HMM



"Naïve Bayes":

$$P(\mathbf{X}, \mathbf{Y}) = \prod_{t=1}^{T} P(X_t | Y_t) p(Y_t)$$

HMM:

$$P(\mathbf{X}, \mathbf{Y} | Y_0) = \prod_{t=1}^{T} P(X_t | Y_t) p(Y_t | Y_{t-1})$$

# SUPERVISED LEARNING FOR HMMS

# Recipe for Closed-form MLE

1. Assume data was generated i.i.d. from some model
   (i.e. write the generative story)
   $$x^{(i)} \sim p(x|\boldsymbol{\theta})$$

2. Write log-likelihood
   $$\ell(\boldsymbol{\theta}) = \log p(x^{(1)}|\boldsymbol{\theta}) + \ldots + \log p(x^{(N)}|\boldsymbol{\theta})$$

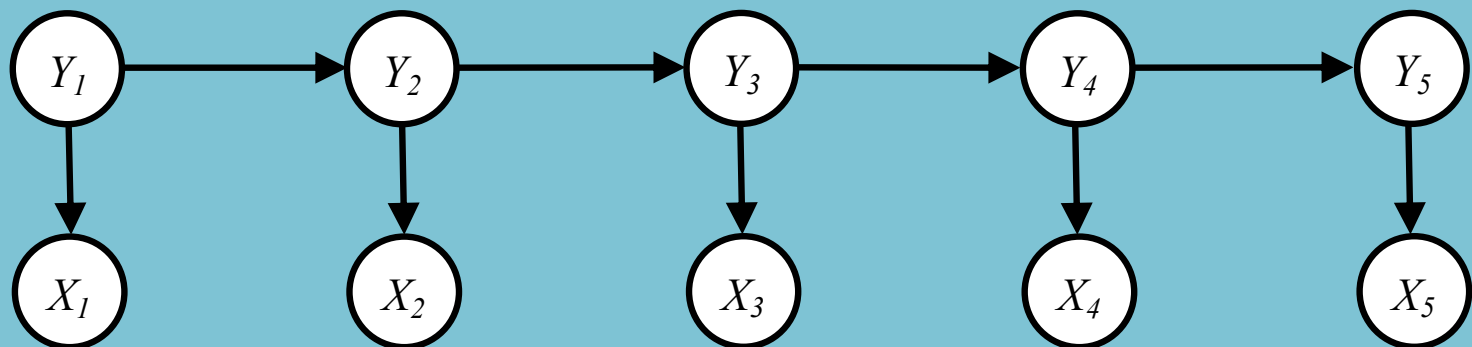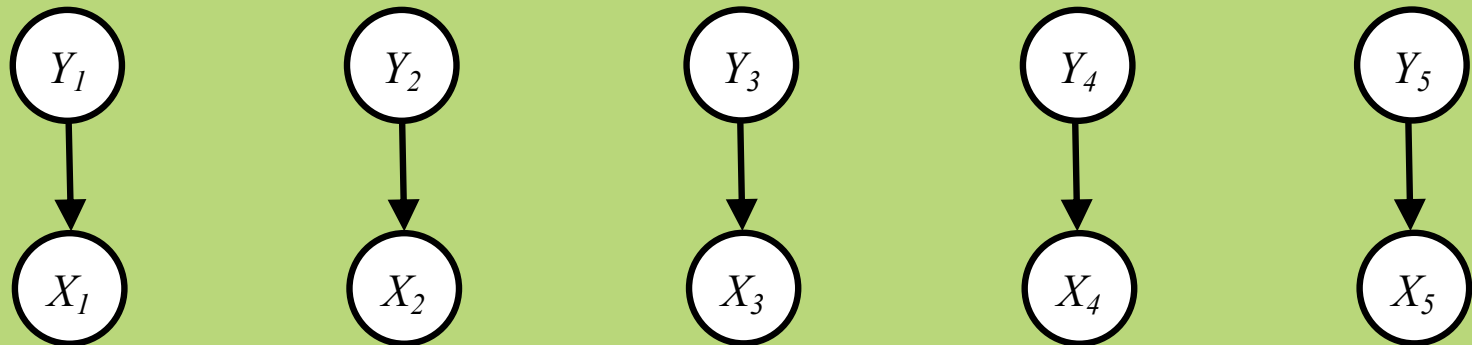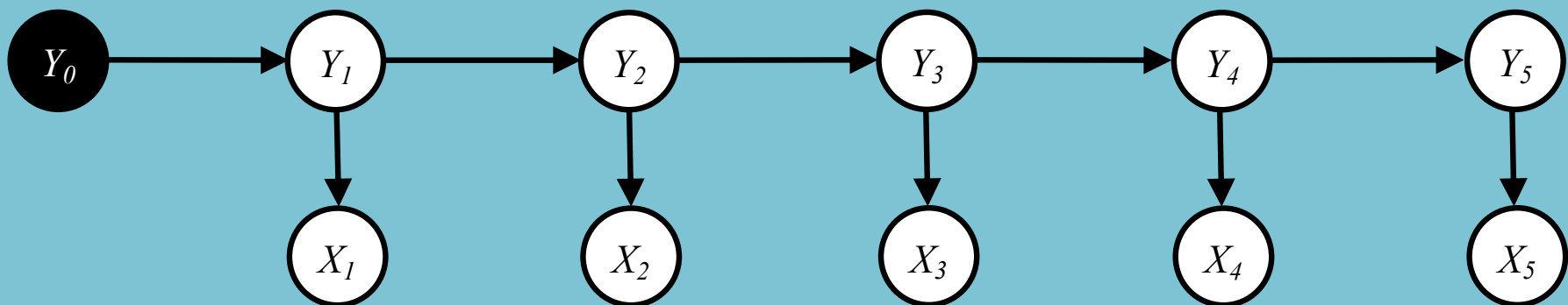3. Compute partial derivatives (i.e. gradient)
   $$\partial\ell(\boldsymbol{\theta})/\partial\theta_1 = \ldots$$
   $$\partial\ell(\boldsymbol{\theta})/\partial\theta_2 = \ldots$$
   $$\ldots$$
   $$\partial\ell(\boldsymbol{\theta})/\partial\theta_M = \ldots$$

4. Set derivatives to zero and solve for $\boldsymbol{\theta}$
   $$\partial\ell(\boldsymbol{\theta})/\partial\theta_m = 0 \text{ for all } m \in \{1, \ldots, M\}$$
   $$\boldsymbol{\theta}^{MLE} = \text{solution to system of } M \text{ equations and } M \text{ variables}$$

5. Compute the second derivative and check that $\ell(\boldsymbol{\theta})$ is concave down
   at $\boldsymbol{\theta}^{MLE}$

# MLE of Categorical Distribution

1. Suppose we have a **dataset** obtained by repeatedly rolling a $M$-sided (weighted) die $N$ times. That is, we have data

$$\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$$

   where $x^{(i)} \in \{1, \ldots, M\}$ and $x^{(i)} \sim \text{Categorical}(\boldsymbol{\phi})$.

2. A random variable is **Categorical** written $X \sim \text{Categorical}(\boldsymbol{\phi})$ iff

$$P(X = x) = p(x; \boldsymbol{\phi}) = \phi_x$$

   where $x \in \{1, \ldots, M\}$ and $\sum_{m=1}^{M} \phi_m = 1$. The **log-likelihood** of the data becomes:

$$\ell(\boldsymbol{\phi}) = \sum_{i=1}^{N} \log \phi_{x^{(i)}} \text{ s.t. } \sum_{m=1}^{M} \phi_m = 1$$

3. Solving this *constrained* optimization problem yields the **maximum likelihood estimator** (MLE):

$$\phi_m^{MLE} = \frac{N_{x=m}}{N} = \frac{\sum_{i=1}^{N} \mathbb{I}(x^{(i)} = m)}{N}$$

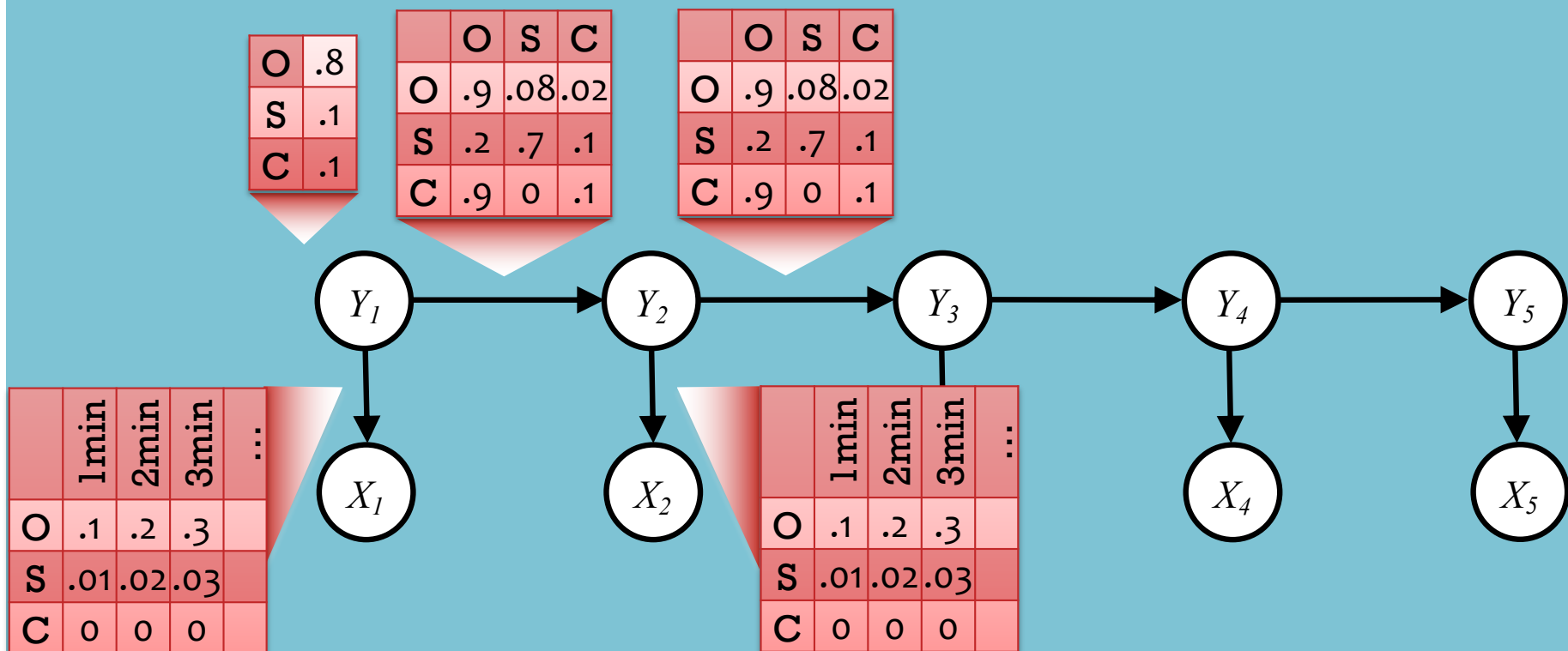# Hidden Markov Model

**HMM Parameters:**

Emission matrix, $\mathbf{A}$, where $P(X_t = k | Y_t = j) = A_{j,k}, \forall t, k$

Transition matrix, $\mathbf{B}$, where $P(Y_t = k | Y_{t-1} = j) = B_{j,k}, \forall t, k$

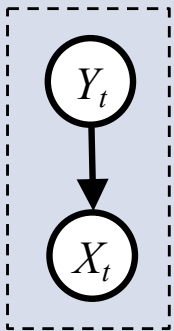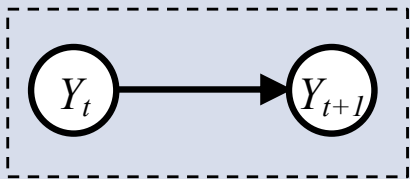Initial probs, $\mathbf{C}$, where $P(Y_1 = k) = C_k, \forall k$

# Training HMMs

*Whiteboard*

- (Supervised) Likelihood for an HMM
- Maximum Likelihood Estimation (MLE) for HMM

# Supervised Learning for HMMs

Learning an HMM decomposes into solving two (independent) Mixture Models



Data: $D = \{(\vec{x}^{(i)}, \vec{y}^{(i)})\}_{i=1}^{N}$  $\qquad \vec{x} = [x_1, \ldots, x_T]^T$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \vec{y} = [y_1, \ldots, y_T]^T$

Likelihood:

$$\ell(A,B,C) = \sum_{i=1}^{N} \log p(\vec{x}^{(i)}, \vec{y}^{(i)} \mid A, B, C)$$

$$= \sum_{i=1}^{N} \left[ \underbrace{\log p(y_1^{(i)} \mid C)}_{\text{initial}} + \underbrace{\left( \sum_{t=2}^{T} \log p(y_t^{(i)} \mid y_{t-1}^{(i)}, B) \right)}_{\text{transition}} + \underbrace{\left( \sum_{t=1}^{T} \log p(x_t^{(i)} \mid y_t^{(i)}, A) \right)}_{\text{emission}} \right]$$

MLE:

$$\hat{A}, \hat{B}, \hat{C} = \underset{A,B,C}{\text{argmax}} \; \ell(A,B,C)$$

$$\Rightarrow \hat{C} = \underset{C}{\text{argmax}} \sum_{i=1}^{N} \log p(y_1^{(i)} \mid C)$$

$$\hat{B} = \underset{B}{\text{argmax}} \sum_{i=1}^{N} \sum_{t=2}^{T} \log p(y_t^{(i)} \mid y_{t-1}^{(i)}, B)$$

$$\hat{A} = \underset{A}{\text{argmax}} \sum_{i=1}^{N} \sum_{t=1}^{T} \log p(x_t^{(i)} \mid y_t^{(i)}, A)$$

Can solve in closed form, which yields...

$$\hat{C}_k = \frac{\#(y_1^{(i)} = k)}{N} \qquad \forall i, k$$

$$\hat{B}_{jk} = \frac{\#(y_t^{(i)} = k \text{ and } y_{t-1}^{(i)} = j)}{\#(y_{t-1}^{(i)} = j)} \qquad \forall i, t>1, j, k$$

$$\hat{A}_{jk} = \frac{\#(x_t^{(i)} = k \text{ and } y_t^{(i)} = j)}{\#(y_t^{(i)} = j)} \qquad \forall i, t, j, k$$

63

# Hidden Markov Model

**HMM Parameters:**

Emission matrix, $\mathbf{A}$, where $P(X_t = k | Y_t = j) = A_{j,k}, \forall t, k$

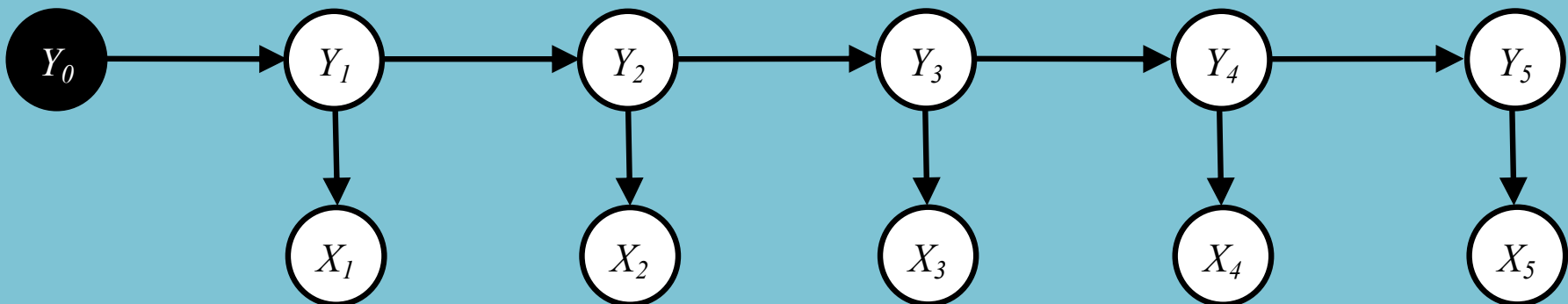Transition matrix, $\mathbf{B}$, where $P(Y_t = k | Y_{t-1} = j) = B_{j,k}, \forall t, k$

**Assumption:** $y_0 = \text{START}$

**Generative Story:**

$$Y_t \sim \text{Multinomial}(\mathbf{B}_{Y_{t-1}}) \ \forall t$$

$$X_t \sim \text{Multinomial}(\mathbf{A}_{Y_t}) \ \forall t$$

For notational convenience, we fold the *initial probabilities* **C** into the *transition matrix* **B** by our assumption.
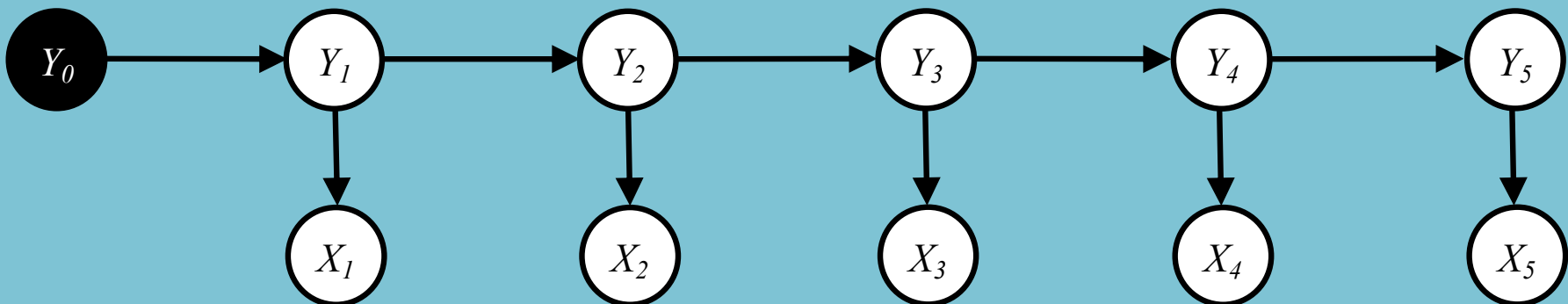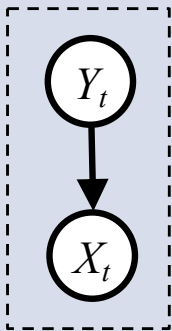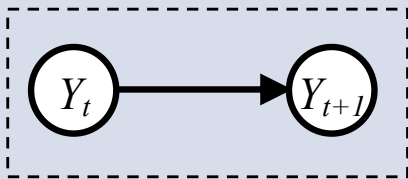


65

# Hidden Markov Model

**Joint Distribution:**

$y_0 = \text{START}$

$$p(\mathbf{x}, \mathbf{y}|y_0) = \prod_{t=1}^{T} p(x_t|y_t)p(y_t|y_{t-1})$$

$$= \prod_{t=1}^{T} A_{y_t,x_t} B_{y_{t-1},y_t}$$

# Supervised Learning for HMMs

Learning an HMM decomposes into solving two (independent) Mixture Models



$$D = \{(\vec{x}^{(i)}, \vec{y}^{(i)})\}_{i=1}^{N}$$

Likelihood :  $\ell(A,B) = \sum_{i=1}^{N} \log p(\vec{x}^{(i)}, \vec{y}^{(i)})$

$$= \sum_{i=1}^{N} \left[ \sum_{t=1}^{T} \log p(y_t^{(i)} | y_{t-1}^{(i)}, B) + \log p(x_t^{(i)} | y_t^{(i)}, A) \right]$$

MLE :  $\hat{A}, \hat{B} = \text{argmax } \ell(A,B)$

$$\hat{A} = \text{argmax} \sum_{i=1}^{N} \left[ \sum_{t=1}^{T} \log p(x_t^{(i)} | y_t^{(i)}, A) \right]$$

$$\hat{B} = \text{argmax} \sum_{i=1}^{N} \left[ \sum_{t=1}^{T} \log p(y_t^{(i)} | y_{t-1}^{(i)}, B) \right]$$

← can solve in closed form to set...

$$\hat{B}_{jk} = \frac{\#\left(y_t^{(i)} = k \text{ and } y_{t-1}^{(i)} = j\right)}{\#\left(y_{t-1}^{(i)} = j\right)}$$

$$\hat{A}_{jk} = \frac{\#\left(x_t^{(i)} = k \text{ and } y_t^{(i)} = j\right)}{\#\left(y_t^{(i)} = j\right)}$$

67