# Hidden Markov Models

# +

# Midterm Exam 2 Review

Matt Gormley
Lecture 19
Mar. 27, 2020

# Reminders

- **Homework 6: Learning Theory / Generative Models**
  - **Out: Fri, Mar 20**
  - **Due: Fri, Mar 27 at 11:59pm**
- **Practice Problems for Exam 2**
  - **Out: Fri, Mar 20**
- **Midterm Exam 2**
  - **Thu, Apr 2 – evening exam, details announced on Piazza**

- **Today's In-Class Poll**
  - **http://poll.mlcourse.org**

# MIDTERM EXAM LOGISTICS

# Midterm Exam

- **Time / Location**
  - **Time:** Evening Exam
    **Thu, Apr. 2 at 6:00pm – 9:00pm**
  - **Location:** We will contact you with additional details about how to join the appropriate Zoom meeting.
  - **Seats:** There will be **assigned Zoom rooms**. Please arrive online early.
  - Please watch Piazza carefully for announcements.
- **Logistics**
  - Covered material: Lecture 9 – Lecture 18 (95%), Lecture 1 – 8 (5%)
  - 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)

# Midterm Exam

- **How to Prepare**
  - Attend the midterm review lecture (right now!)
  - Review prior year's exam and solutions (we'll post them)
  - Review this year's homework problems
  - Consider whether you have achieved the "learning objectives" for each lecture / section

# Midterm Exam

- **Advice (for during the exam)**
  - Solve the easy problems first
    (e.g. multiple choice before derivations)
    - if a problem seems extremely complicated you're likely missing something
  - Don't leave any answer blank!
  - If you make an assumption, write it down *unresolved*
  - If you look at a question and don't know the answer:
    - we probably haven't told you the answer
    - but we've told you enough to work it out
    - imagine arguing for some answer and see if you like it

# Topics for Midterm 1

- Foundations
  - Probability, Linear Algebra, Geometry, Calculus
  - Optimization

- Important Concepts
  - Overfitting
  - Experimental Design

- Classification
  - Decision Tree
  - KNN
  - Perceptron

- Regression
  - Linear Regression

# Topics for Midterm 2

- Classification
  - Binary Logistic Regression
  - Multinomial 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

# 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?

$$h(\vec{x}) = \underset{y \in \{0,1\}}{\arg\max}\ p(y|x) = \begin{cases} 1 & \text{if } p(y=1|\vec{x}) \geq 0.5 \\ 0 & \text{otherwise} \end{cases} \quad / \text{ linear decision boundaries}$$

(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.

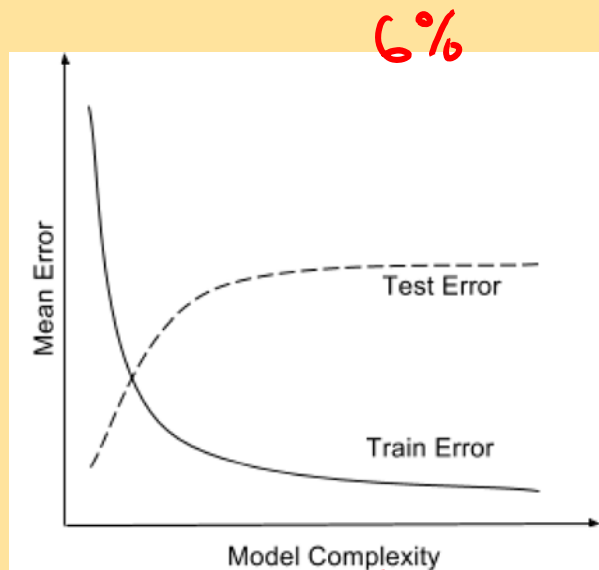1. **[4 pts]** Which of the following is expected to help? Select all that apply.

   ✓ (a) Increase the training data size. **80%**

   (b) Decrease the training data size. **20%**

   (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). **7%**

   ✓ (d) Decrease model complexity. **90%**

   (e) Train on a combination of $\mathcal{D}^{\text{train}}$ and $\mathcal{D}^{\text{test}}$ and test on $\mathcal{D}^{\text{test}}$ **11%**

   (f) ~~Conclude that Machine Learning does not work.~~ **Calamity**
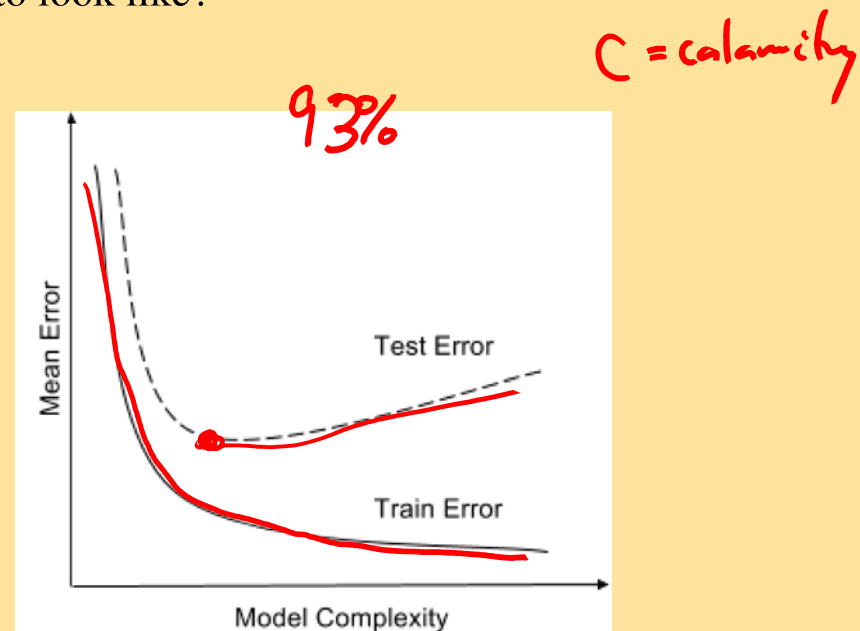
# Samples Questions

## 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?

C = calamity



6%

Mean Error

Test Error

Train Error

Model Complexity

(a) A

93%

Mean Error

Test Error

Train Error

Model Complexity

(b) B

12

# 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 concept class with infinite size is also infinite.

$c^* \in \mathcal{H}$

$D$

(f) [3 pts.] **T or F**: Given a realizable concept class and a set of training instances, a consistent learner will output a concept that achieves 0 error on the training instances.

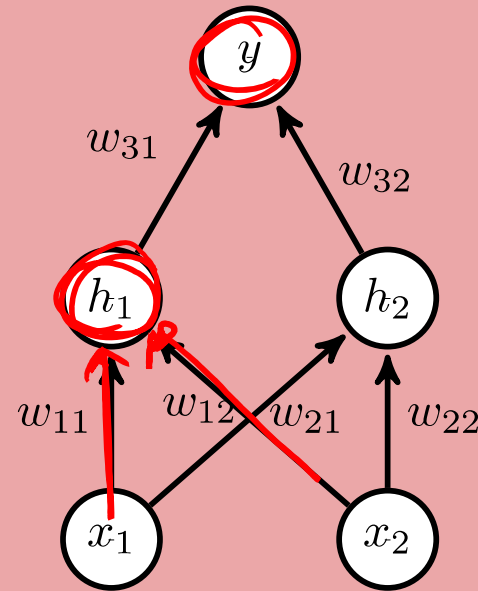will find a consistent hypothesis if it exists

hypothesis

# Sample Questions

**Neural Networks**

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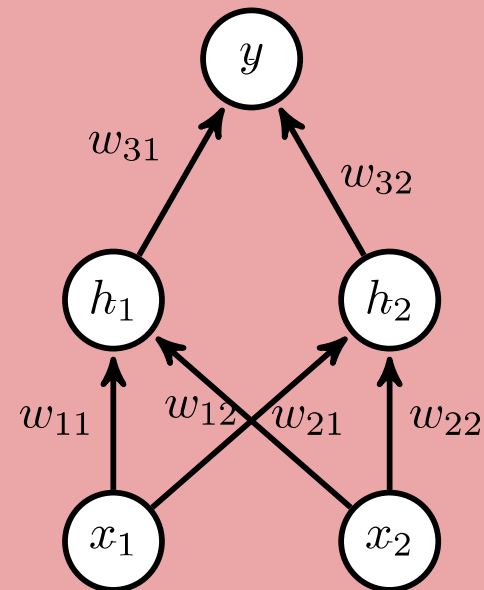


(a) The dataset with groups $S_1$, $S_2$, and $S_3$.

(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.

*Gaussian NB*

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

# HIDDEN MARKOV MODEL (HMM)

# HMM Outline

- **Motivation**
  - Time Series Data
- **Hidden Markov Model (HMM)**
  - Example: Squirrel Hill Tunnel Closures [courtesy of Roni Rosenfeld]
  - Background: Markov Models
  - From Mixture Model to HMM
  - History of HMMs
  - Higher-order HMMs
- **Training HMMs**
  - (Supervised) Likelihood for HMM
  - Maximum Likelihood Estimation (MLE) for HMM
  - EM for HMM (aka. Baum-Welch algorithm)
- **Forward-Backward Algorithm**
  - Three Inference Problems for HMM
  - Great Ideas in ML: Message Passing
  - Example: Forward-Backward on 3-word Sentence
  - Derivation of Forward Algorithm
  - Forward-Backward Algorithm
  - Viterbi algorithm

# Markov Models

*Whiteboard*

- Example: Tunnel Closures
  [courtesy of Roni Rosenfeld]

- First-order Markov assumption
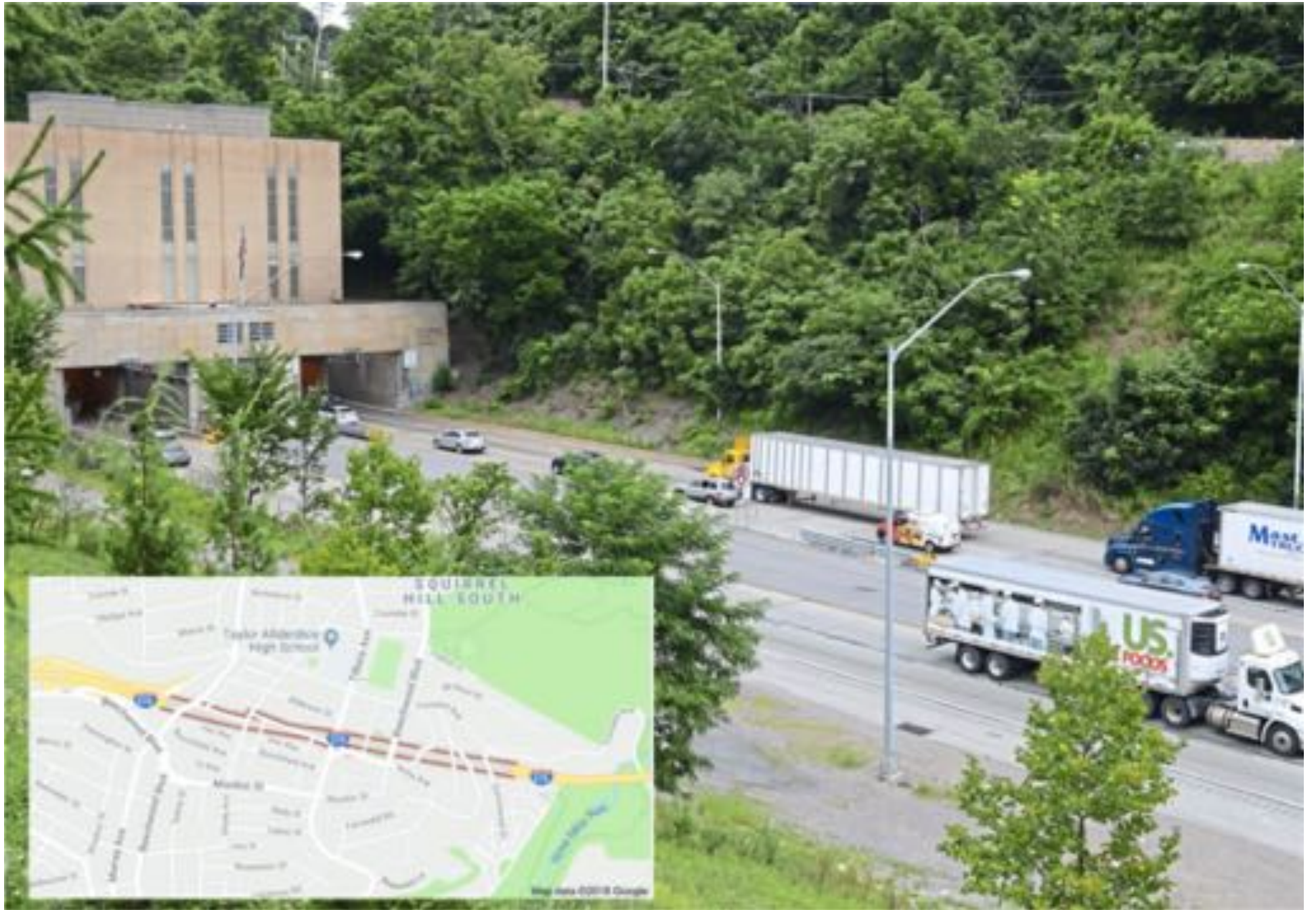
- Conditional independence assumptions
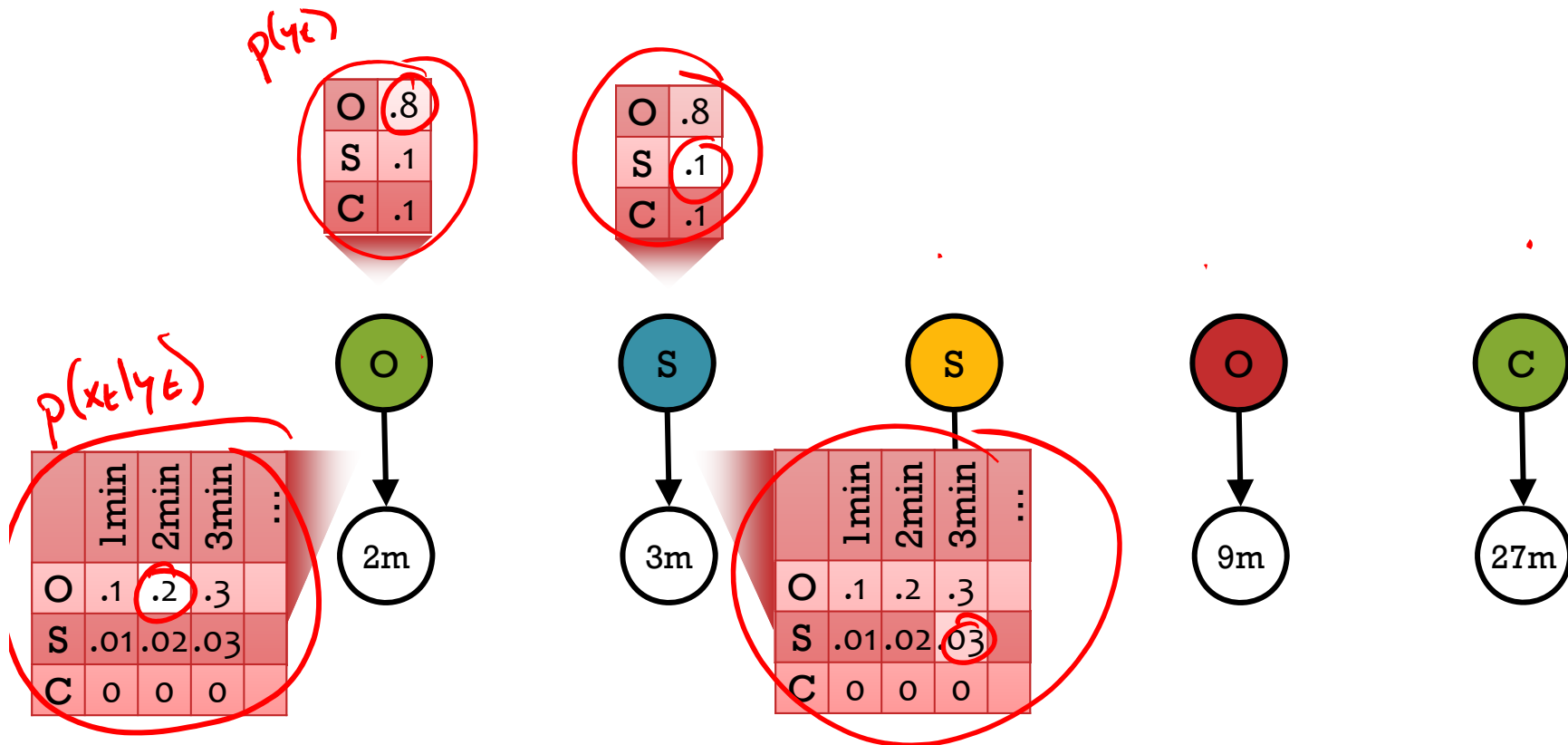
# Totoro's Tunnel

# 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).
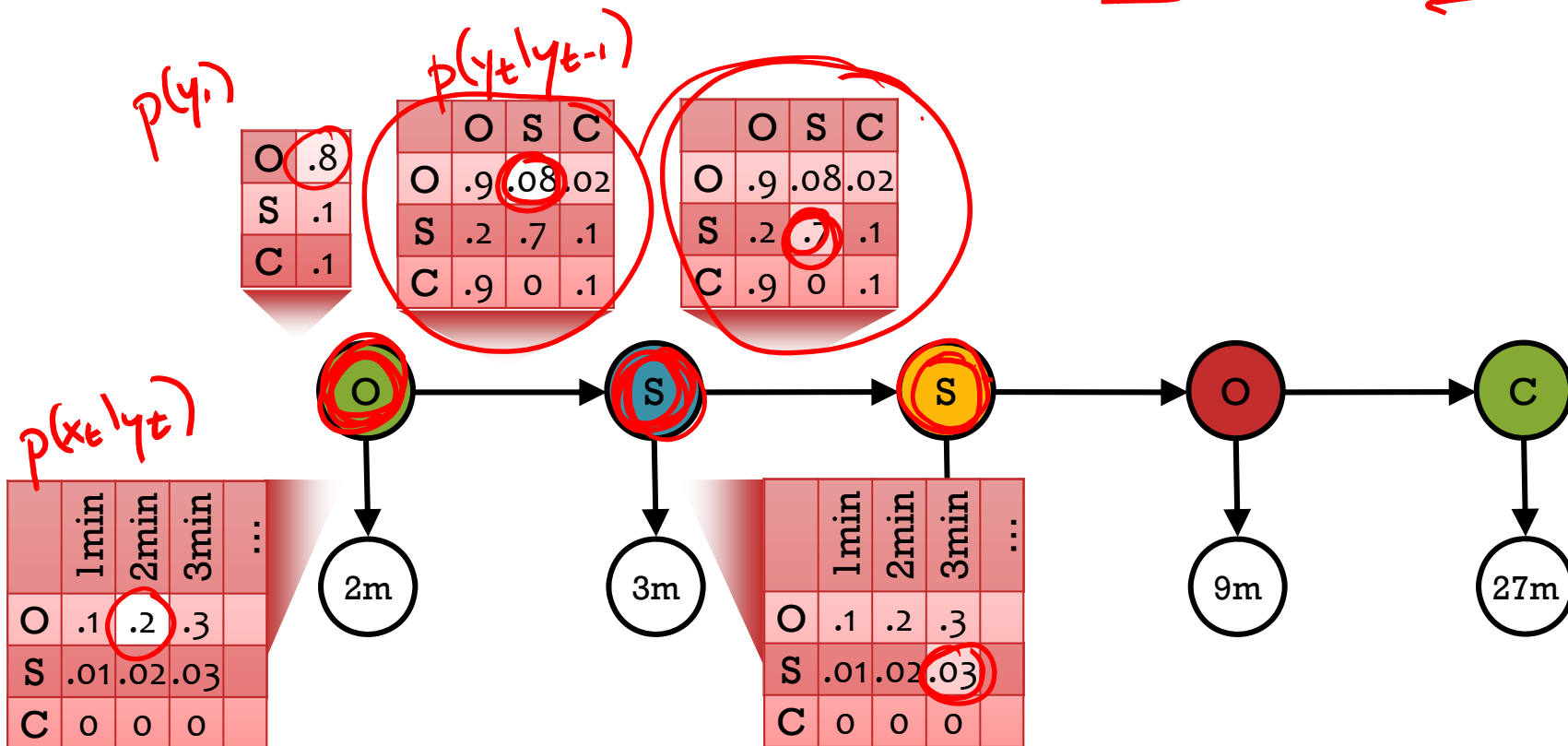
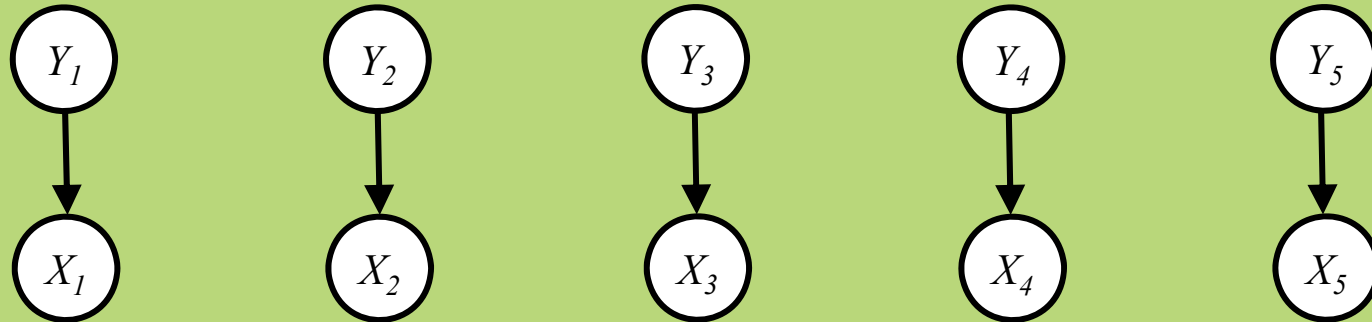$$p(O, S, S, O, C, 2m, 3m, 18m, 9m, 27m) \quad = \quad (.8 * .2 * .1 * .03 * \ldots)$$

$p(y_t)$

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

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

$p(x_t | y_t)$

| | 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 | |

O → 2m  
S → 3m  
S → ...  
O → 9m  
C → 27m

# 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.

$$p(O, S, S, O, C, 2m, 3m, 18m, 9m, 27m) = (.8 * .08 * .2 * .7 * .03 * \ldots)$$
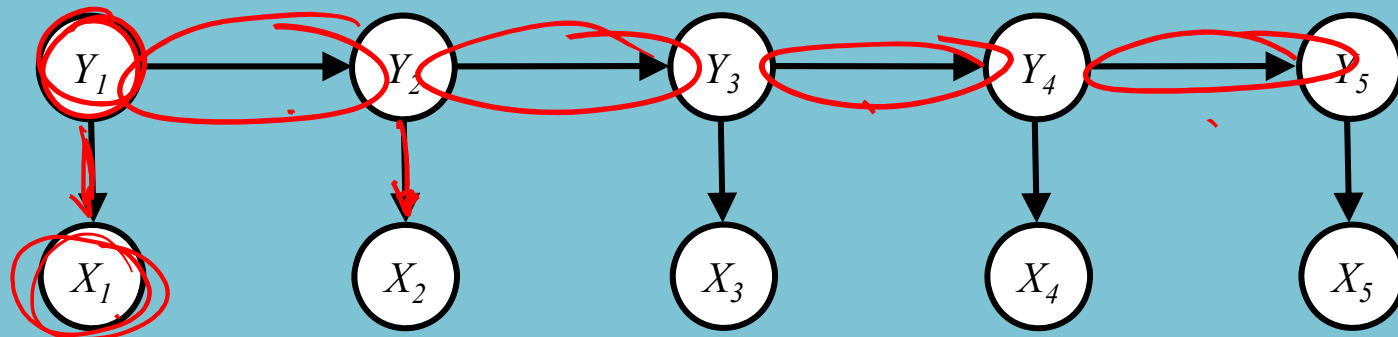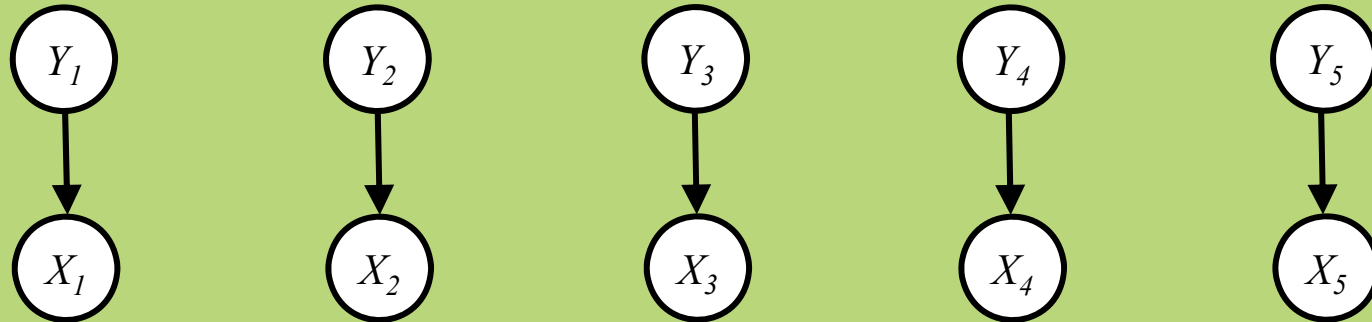
# From Mixture Model to HMM



"Naïve Bayes":

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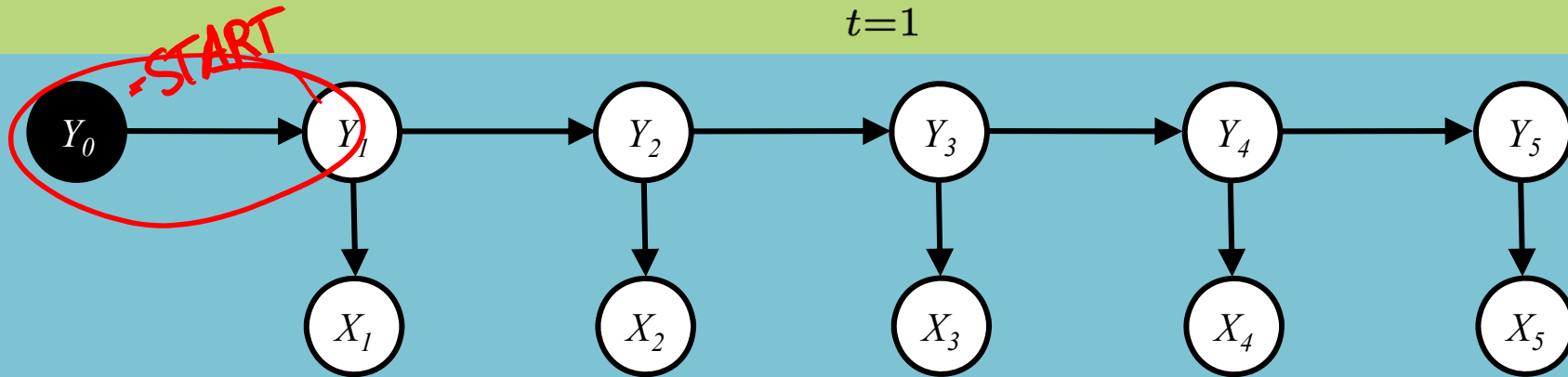
HMM:

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

31

# From Mixture Model to HMM



"Naïve Bayes":

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

START

HMM:

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

$p(y_1 | y_0) \neq p(y_1)$

32

# SUPERVISED LEARNING FOR HMMS

# Recipe for Closed-form MLE

$(x^{(i)}, \vec{y}^{(i)})$

1. Assume data was generated i.i.d. from some model
   (i.e. write the generative story)

   HMM

   $$x^{(i)} \sim p(x|\boldsymbol{\theta})$$

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 $\boldsymbol{\theta}$

   $$\partial\ell(\boldsymbol{\theta})/\partial\theta_m = 0 \text{ for all } m \in \{1, \ldots, M\}$$

   $\boldsymbol{\theta}^{MLE}$ = solution to system of M equations and M 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 $x^{(i)} \sim \text{Categorical}(\boldsymbol{\phi})$. *vector*

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

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

$\phi_m > 0 \ \forall m$

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):

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

39
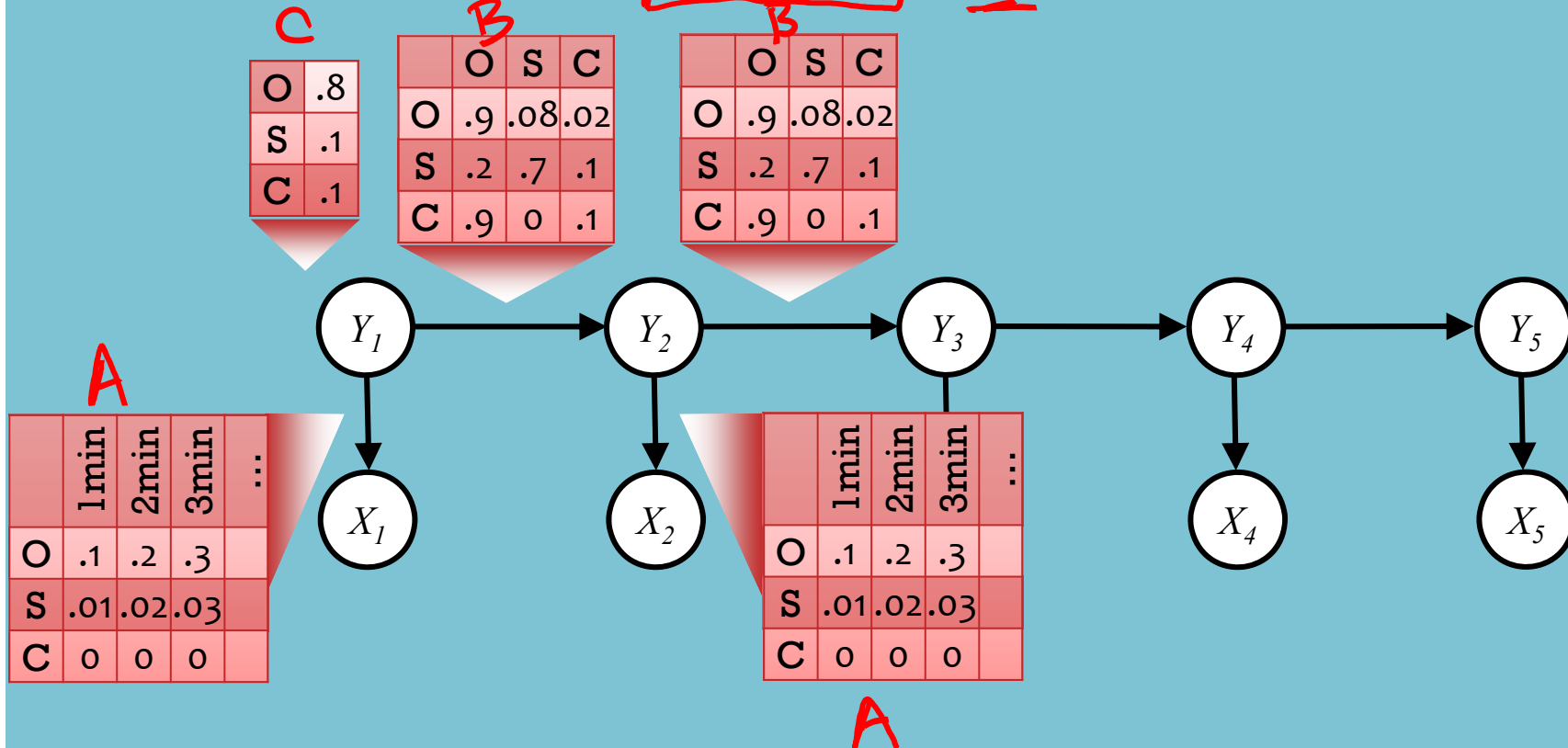
# Hidden Markov Model

**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 $P(Y_t = k | Y_{t-1} = j) = B_{j,k}, \forall t, k$

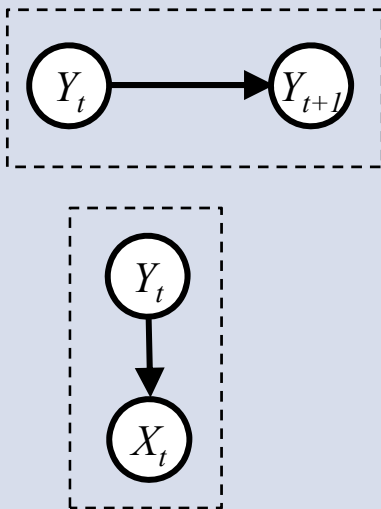Initial probs, $\mathbf{C}$, where $P(Y_1 = k) = C_k, \forall k$

# 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



$$\text{Data :} \quad D = \{(\vec{x}^{(i)}, \vec{y}^{(i)})\}_{i=1}^{N} \qquad \vec{x} = [x_1, \ldots, x_T]^T$$
$$\vec{y} = [y_1, \ldots, y_T]^T$$

**Likelihood :**

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

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

**MLE :**

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

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

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

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

Can solve in closed form, which yields...

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

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

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

43

# Hidden Markov Model

**HMM Parameters:**

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

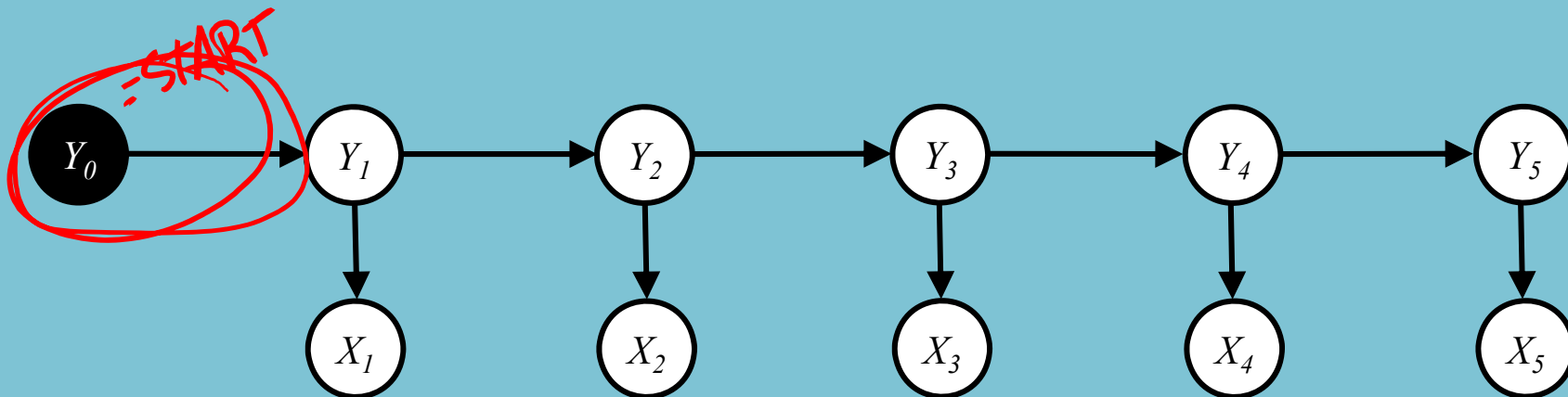Transition matrix, $\mathbf{B}$, where $P(Y_t = k | Y_{t-1} = j) = B_{j,k}, \forall t, k$

**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* $\mathbf{C}$ into the *transition matrix* $\mathbf{B}$ by our assumption.



44

# Hidden Markov Model

**Joint Distribution:**

$$y_0 = \text{START}$$

$$p(\mathbf{x}, \mathbf{y} | y_0) = \prod_{t=1}^{T} p(x_t | y_t) p(y_t | y_{t-1})$$
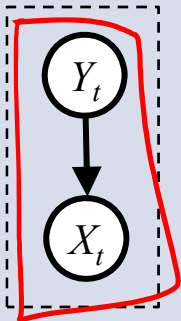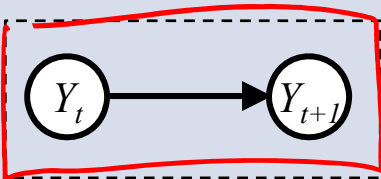
*C is in here, B is larger*

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

# Supervised Learning for HMMs

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



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

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

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

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

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

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

can solve in closed form to get...

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

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

# Unsupervised Learning for HMMs

- Unlike **discriminative** models p(y|x), **generative** models p(x,y) can maximize the likelihood of the data D = {x$^{(1)}$, x$^{(2)}$, …, x$^{(N)}$} where we don't observe any y's.
- This **unsupervised learning** setting can be achieved by finding parameters that maximize the **marginal likelihood**
- We optimize using the **Expectation-Maximization** algorithm

Since we don't observe **y**, we define the marginal probability:

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{Y}} p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{y}) \tag{1}$$

The log-likelihood of the data is thus:

$$\ell(\boldsymbol{\theta}) = \log \prod_{i=1}^{N} p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})$$

$$= \sum_{i=1}^{N} \log \sum_{\mathbf{y} \in \mathcal{Y}} p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \mathbf{y}) \tag{3}$$
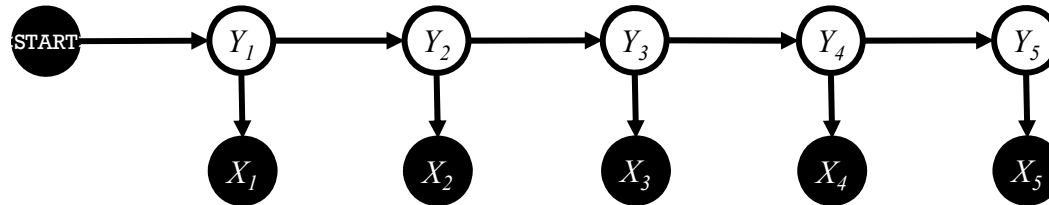
Beyond the scope of today's lecture!

# HMMs: History

- Markov chains: Andrey Markov (1906)
  - Random walks and Brownian motion
- Used in Shannon's work on information theory (1948)
- Baum-Welsh learning algorithm: late 60's, early 70's.
  - Used mainly for speech in 60s-70s.
- Late 80's and 90's: David Haussler  (major player in learning theory in 80's) began to use HMMs for modeling biological sequences
- Mid-late 1990's: Dayne Freitag/Andrew McCallum
  - Freitag thesis with Tom Mitchell on IE from Web using logic programs, grammar induction, etc.
  - McCallum:  multinomial Naïve Bayes for text
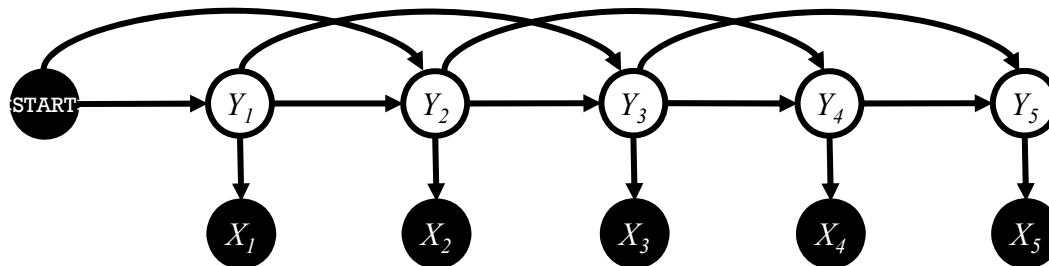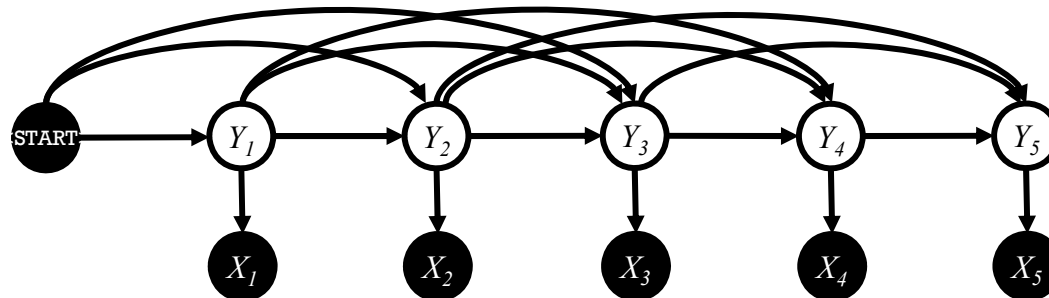  - With McCallum, IE using HMMs on CORA
- …

# Higher-order HMMs

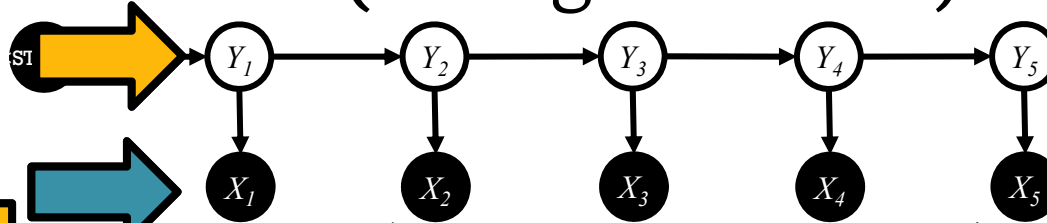- 1st-order HMM (i.e. bigram HMM)



- 2nd-order HMM (i.e. trigram HMM)



- 3rd-order HMM

# Higher-order HMMs

- 1$^{st}$-order HMM (i.e. bigram HMM)



2$^{nd}$-order HMM (i.e. trigram HMM)



3$^{rd}$-order HMM



Hidden States, **y**

Observa-tions, **x**