

Announcements

Assignments

- HW6
 - Due Today, 11:59 pm

Midterm 2

- Mon, 11/9, during lecture
- See Piazza for details on practice exam tomorrow

Schedule change

- Lecture on Friday instead of recitation

Plan

Last Time

- M(C)LE

$$\operatorname{argmax}_{\theta} p(y | \mathbf{x}, \theta)$$

- MAP

$$\operatorname{argmax}_{\theta} p(y | \mathbf{x}, \theta) p(\theta)$$

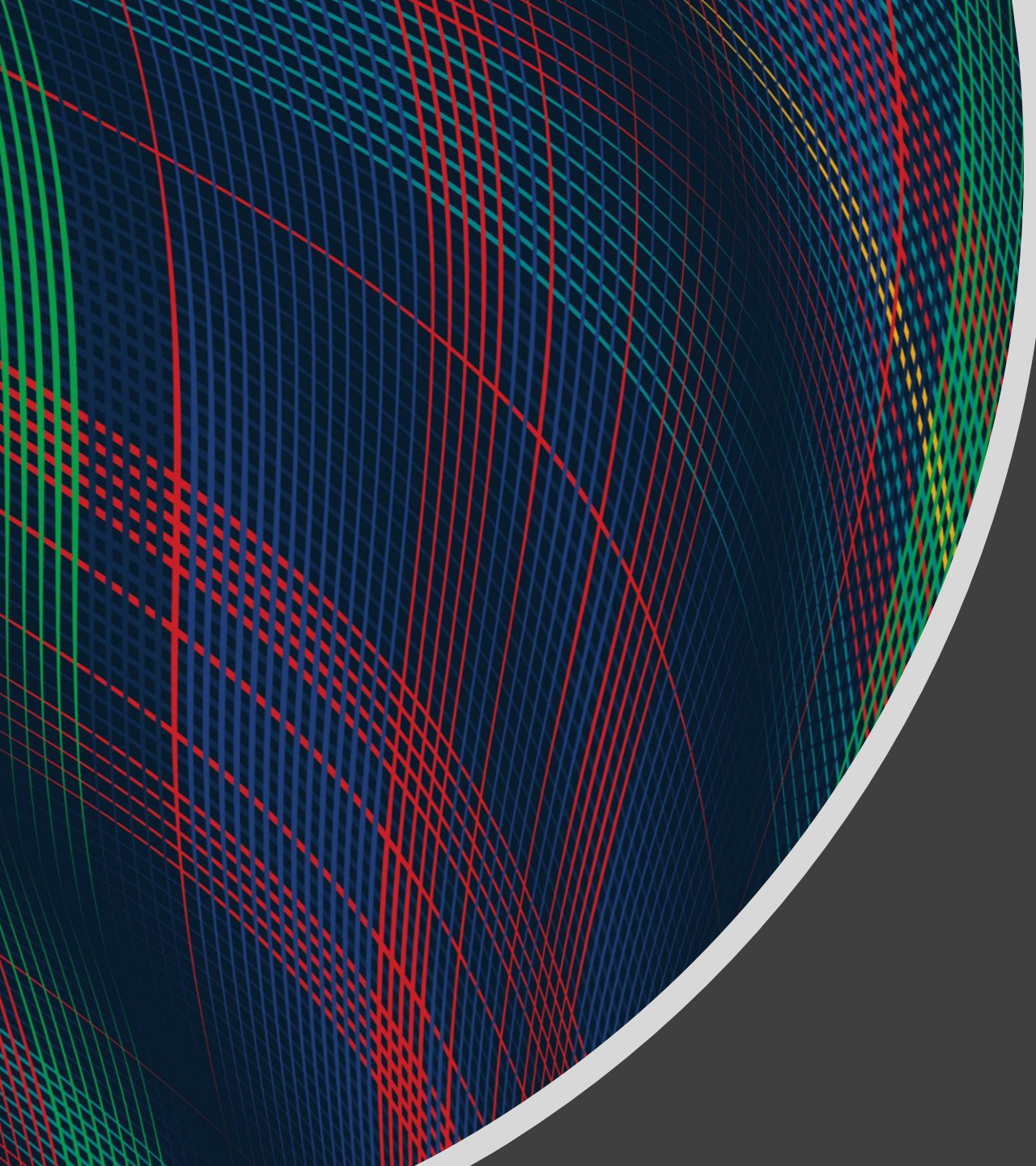
Today

- Generative models

$$\operatorname{argmax}_{\theta} p(\mathbf{x} | y, \theta) p(y | \theta)$$

- Naïve Bayes

$$\operatorname{argmax}_{\theta} \prod_{m=1}^M p(x_m | y, \theta) p(y | \theta)$$



Introduction to Machine Learning

Generative Models & Naïve Bayes

Instructor: Pat Virtue

Recall: Fisher Iris Dataset

https://en.wikipedia.org/wiki/Iris_flower_data_set



Recall: Fisher Iris Dataset

https://en.wikipedia.org/wiki/Iris_flower_data_set

Fisher (1936) used 150 measurements of flowers from 3 different species: Iris setosa (0), Iris virginica (1), Iris versicolor (2) collected by Anderson (1936)

Species	Sepal Length	Sepal Width	Petal Length	Petal Width
0	4.3	3.0	1.1	0.1
0	4.9	3.6	1.4	0.1
0	5.3	3.7	1.5	0.2
1	4.9	2.4	3.3	1.0
1	5.7	2.8	4.1	1.3
1	6.3	3.3	4.7	1.6
1	6.7	3.0	5.0	1.7

Modeling the Fisher Iris Dataset

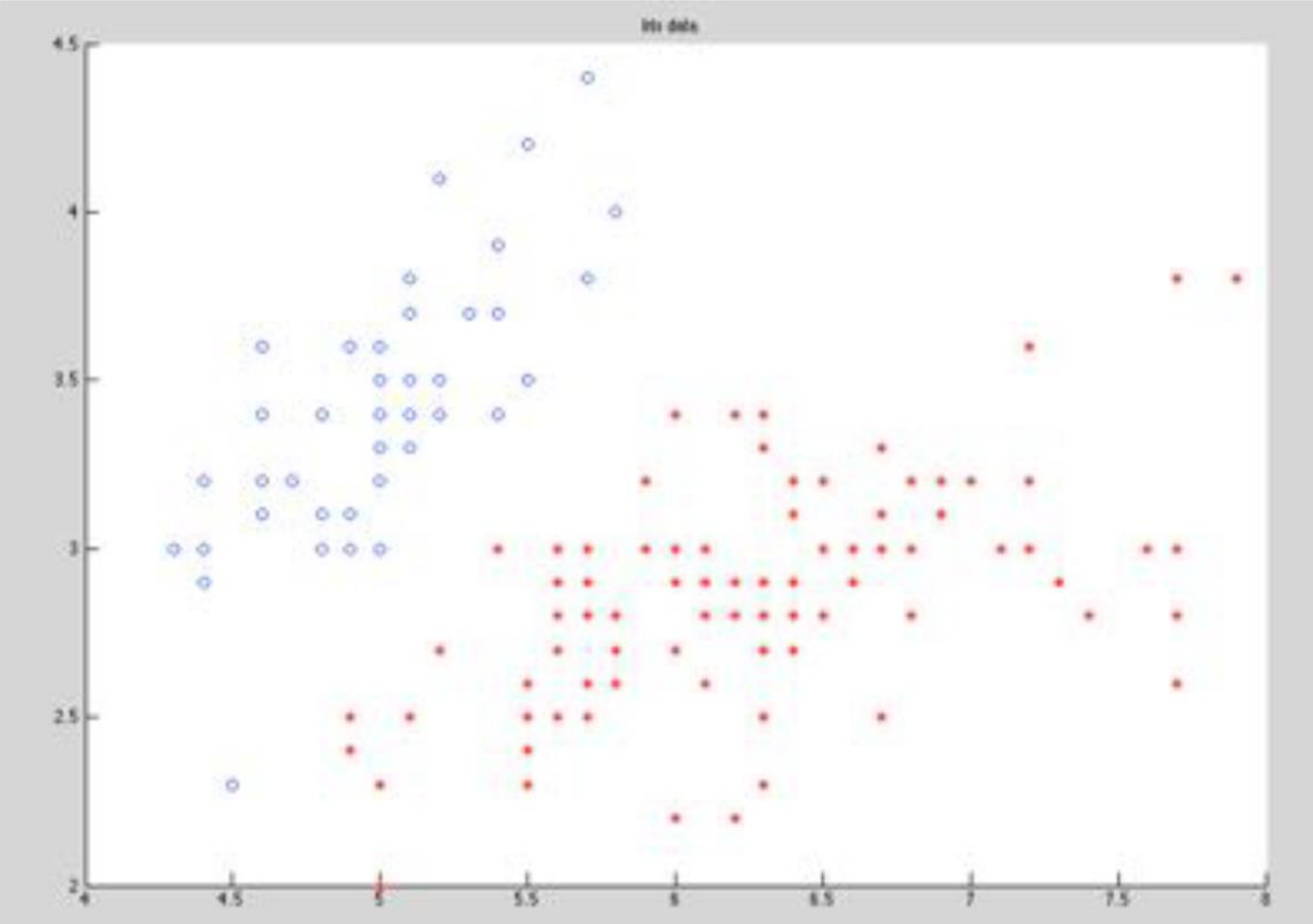


Image: CMU MLD, William Cohen

Modeling the Fisher Iris Dataset

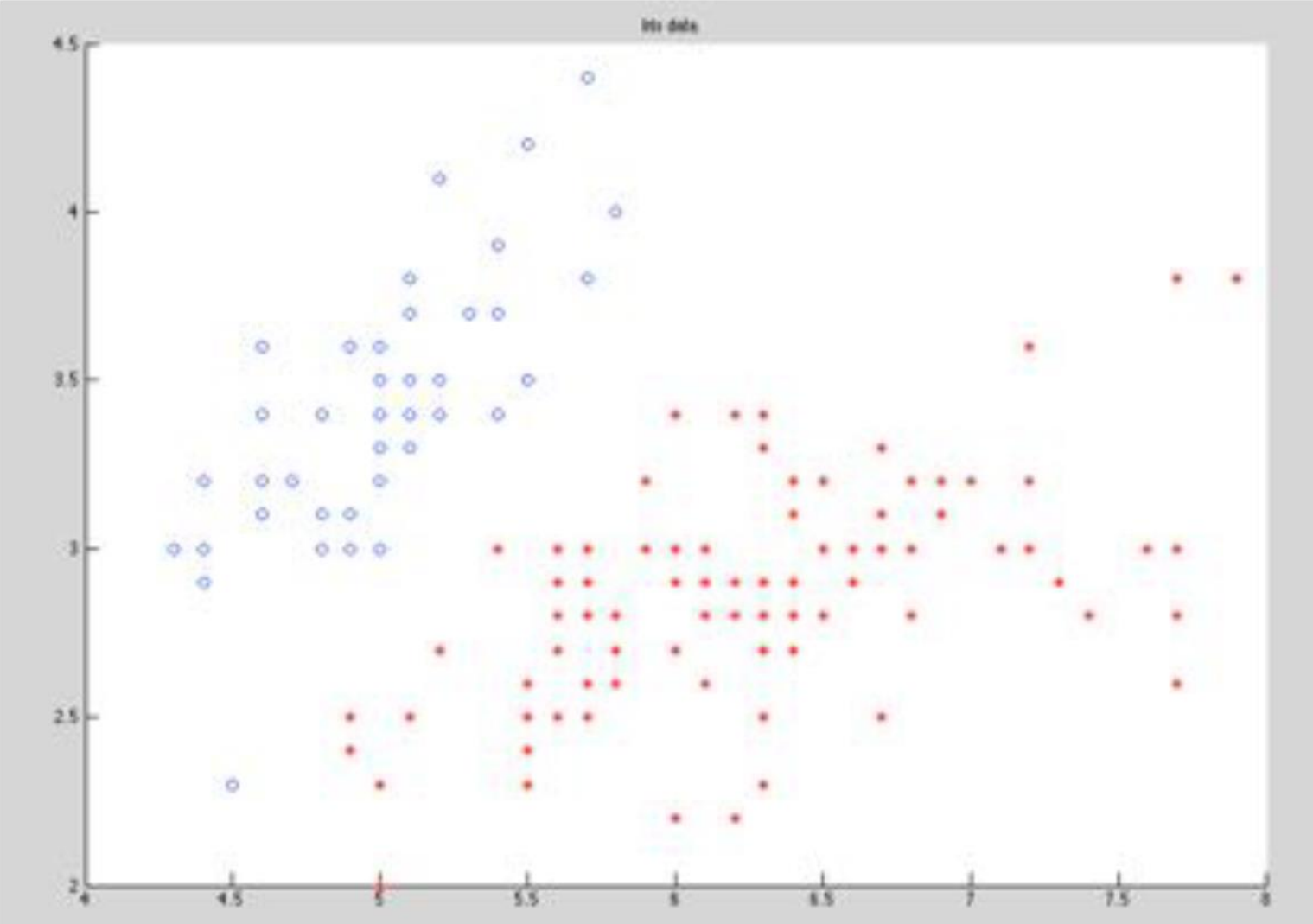
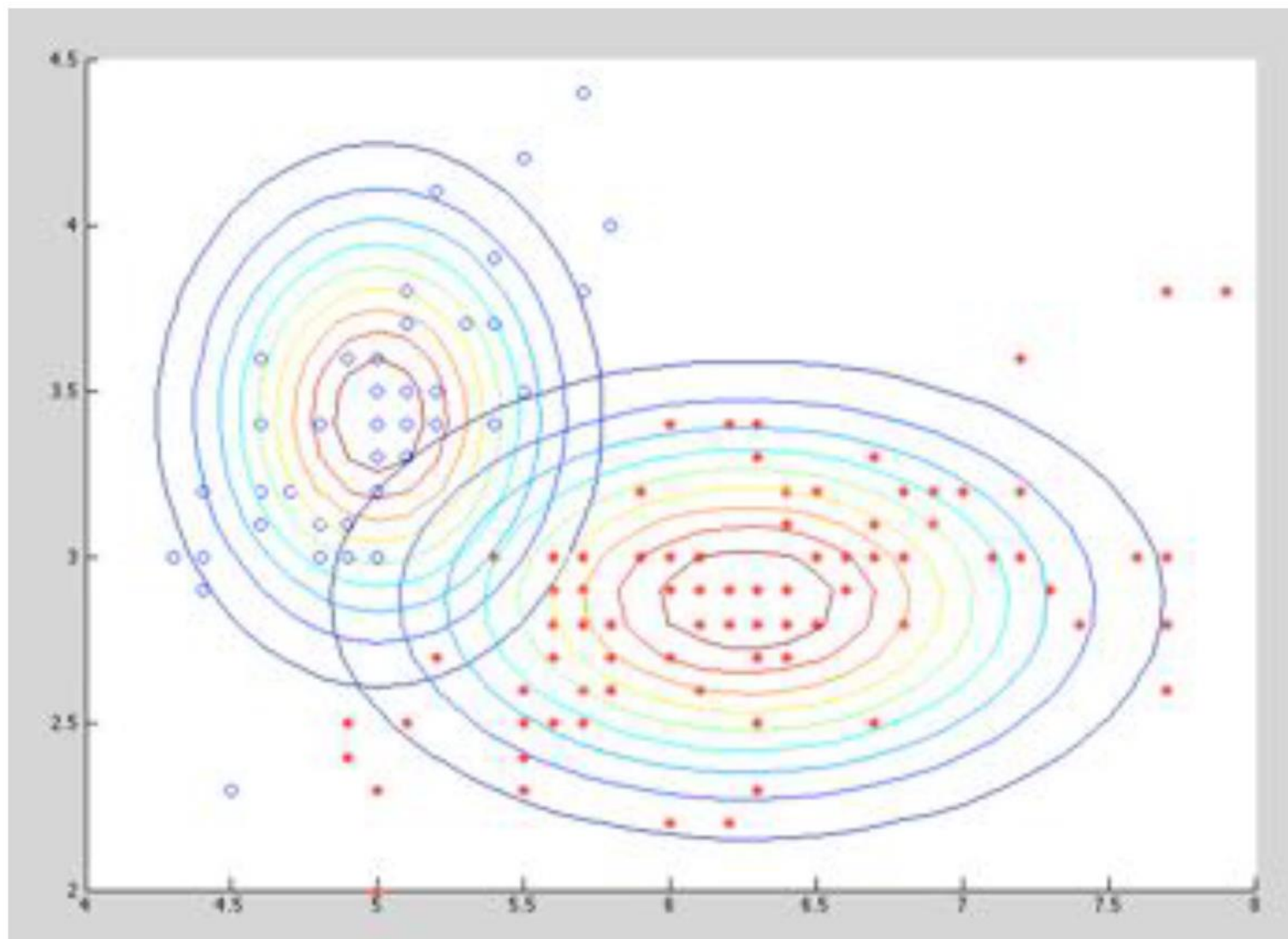


Image: CMU MLD, William Cohen

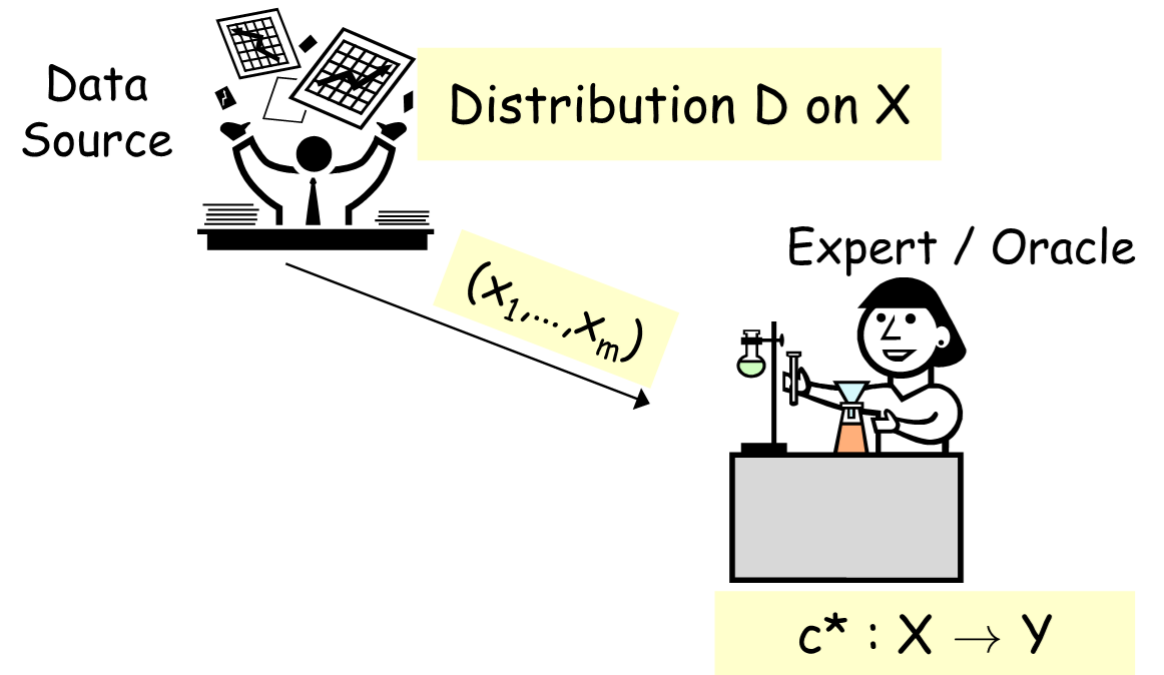
Modeling the Fisher Iris Dataset



Generative vs Discriminative Modeling

Discriminative: modeling $X \rightarrow Y$ directly

Generative: Stronger modeling assumptions about where data came from



Generative vs Discriminative Modeling

Discriminative: model $p(y | x, \theta)$ directly

- Learn parameters θ from data

Generative: model $p(y | \theta_{class})$ and $p(x | y, \theta_{class\ conditional})$

- Learn parameters θ_{class} and $\theta_{class\ conditional}$ from data
- Use Bayes rule to compute $p(y | x, \theta_{class}, \theta_{class\ conditional})$

$$p(y | x, \theta_{class}, \theta_{class\ conditional}) \propto p(x | y, \theta_{class\ conditional}) p(y | \theta_{class})$$

Generative Story

News article topic classification

- Document class: Business, Entertainment, Politics
- Words in the document

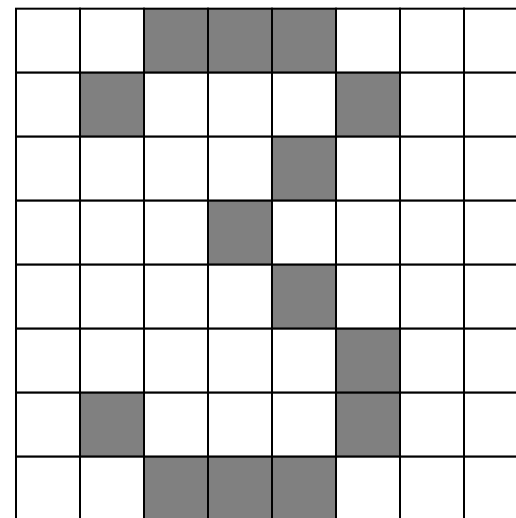
SPAM classification

- Document class: SPAM or not
- Words in the document

Generative Story

Hand-written digits

- Digit class: 0-9
- Pixels in images



Generative vs Discriminative Modeling

Discriminative: $p(y | x)$

Generative: $p(y | x) = \alpha p(x, y) = \alpha p(x | y) p(y)$

Assumptions vs Data

- Discriminative:
- Generative:

Quick Check

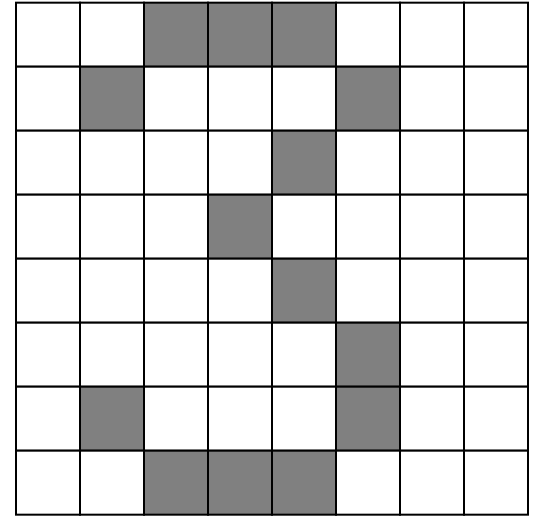
How many parameters?

- $P(Y)$, Y represents outcome of a 6-sided die roll

Multivariate Generative Models

Hand-written digits: How many parameters?

- $P(Y)$
- $P(\mathbf{X} \mid Y = 3)$
 $= p(X_1, X_2, \dots, X_{64} \mid Y = 3)$



Naïve Bayes assumption (bag of pixels)

- $P(Y)$
- $P(\mathbf{X} \mid Y = 3)$
 $= p(X_1 \mid Y = 3)p(X_2 \mid Y = 3) \dots p(X_{64} \mid Y = 3)$

Conditional Independence and Naïve Bayes

Independence

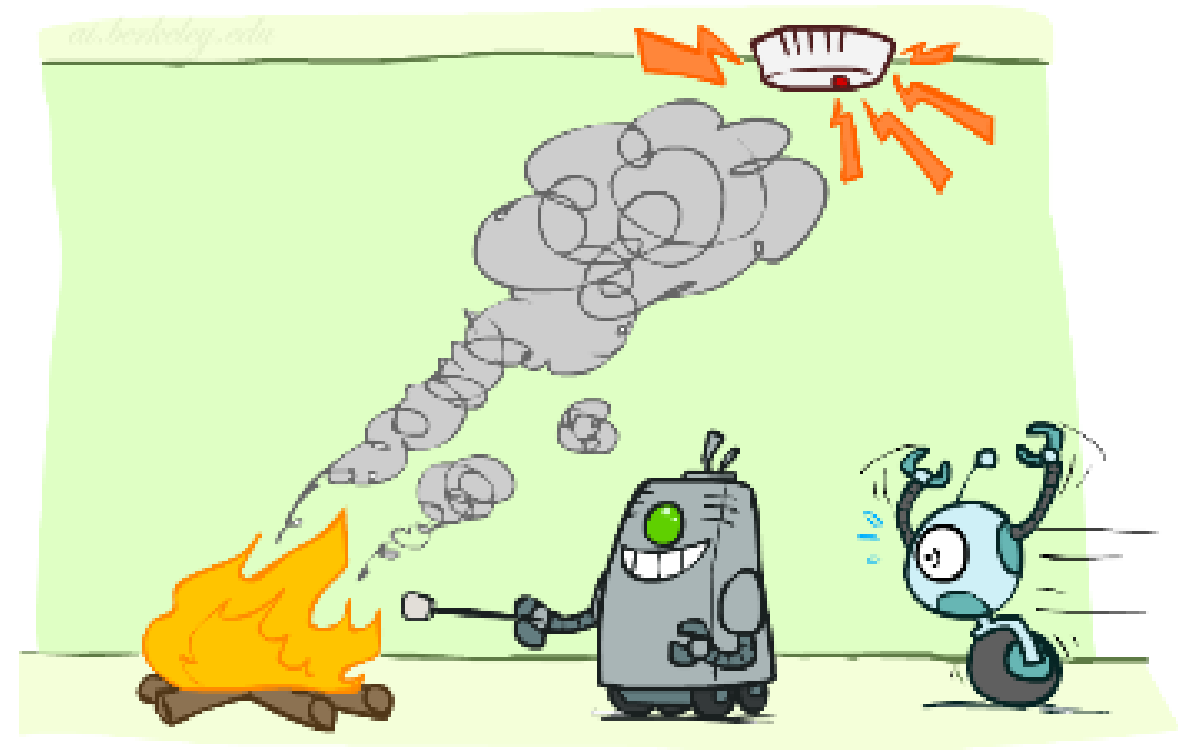
Conditional independence

Conditional Independence

Example: Fire, Smoke, Alarm

Fire and Alarm are independent given Smoke

$$P(A | S, F) = P(A | S)$$



Conditional Independence and Naïve Bayes

Independence

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

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

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

Conditional independence

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

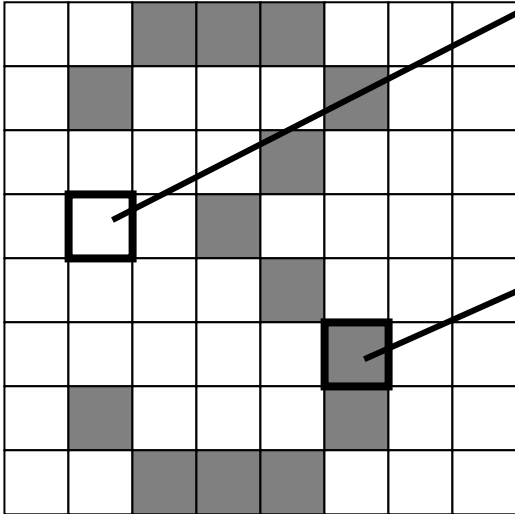
$$P(B | A, C) = P(B | C)$$

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

Naïve Bayes assumption

Naïve Bayes for Digits

y	$P(Y)$
1	0.1
2	0.1
3	0.1
4	0.1
5	0.1
6	0.1
7	0.1
8	0.1
9	0.1
0	0.1



y	$P(X_{3,1} = 1 y)$
1	0.01
2	0.05
3	0.05
4	0.30
5	0.80
6	0.90
7	0.05
8	0.60
9	0.50
0	0.80

y	$P(X_{5,5} = 1 y)$
1	0.05
2	0.01
3	0.90
4	0.80
5	0.90
6	0.90
7	0.25
8	0.85
9	0.60
0	0.80

SPAM Classification

Breakout room exercise

<https://tinyurl.com/301601spam>

SPAM classification

- Y : Binary random variable
Document is SPAM ($Y = 1$) or not ($Y = 0$)
- X_m : Binary random variable
Word m appears in document ($X_m = 1$) or not ($X_m = 0$)

Page 1: Estimate parameters from data

Page 2: Calculate probabilities of new sentence given page 1 parameters

SPAM Classification

Breakout room exercise

$P(Y = 0, X_1, \dots, X_M)$	$P(Y = 1, X_1, \dots, X_M)$

$P(Y = 0 X_1, \dots, X_M)$	$P(Y = 1 X_1, \dots, X_M)$

Reminder MLE for Bernoulli

Bernoulli distribution:

$$Y \sim \text{Bern}(\phi)$$

$$p(y) = \begin{cases} \phi, & y = 1 \\ 1 - \phi, & y = 0 \end{cases}$$

What is the log likelihood for three i.i.d. samples, given parameter ϕ ?

$$\mathcal{D} = \{y^{(1)} = 1, y^{(2)} = 1, y^{(3)} = 0\}$$

$$L(\phi) = \phi \cdot \phi \cdot (1 - \phi) \qquad = \prod_n \phi^{y^{(n)}} (1 - \phi)^{(1-y^{(n)})}$$

$$L(\phi) = \phi^2 \cdot (1 - \phi)^1 \qquad = \phi^{N_{y=1}} (1 - \phi)^{N_{y=0}}$$

Naïve Bayes MLE

$$L(\phi, \Theta) = p(\mathcal{D} \mid \phi, \Theta)$$

$$= \prod_{n=1}^N p(\mathcal{D}^{(n)} \mid \phi, \Theta) \quad \text{i.i.d assumption}$$

$$= \prod_{n=1}^N p(y^{(n)}, \mathbf{x}^{(n)} \mid \phi, \Theta)$$

$$= \prod_{n=1}^N p(y^{(n)} \mid \phi) p(\mathbf{x}^{(n)} \mid y^{(n)}, \Theta) \quad \text{Generative model}$$

$$= \prod_{n=1}^N p(y^{(n)} \mid \phi) p(x_1^{(n)}, x_2^{(n)}, \dots, x_M^{(n)} \mid y^{(n)}, \Theta)$$

$$= \prod_{n=1}^N p(y^{(n)} \mid \phi) \prod_{m=1}^M p(x_m^{(n)} \mid y^{(n)}, \theta_{m,y}) \quad \text{Naïve Bayes}$$

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

$$y^{(n)} \in \{0,1\}$$

$$\mathbf{x}^{(n)} \in \{0,1\}^M$$

$$\phi \in [0,1]$$

$$\Theta \in [0,1]^{M \times 2}$$

Naïve Bayes MLE

$$L(\phi, \Theta) = p(\mathcal{D} \mid \phi, \Theta)$$

$$= \prod_{n=1}^N p(\mathcal{D}^{(n)} \mid \phi, \Theta) \quad \text{i.i.d assumption}$$

$$= \prod_{n=1}^N p(y^{(n)}, \mathbf{x}^{(n)} \mid \phi, \Theta)$$

$$= \prod_{n=1}^N p(y^{(n)} \mid \phi) p(\mathbf{x}^{(n)} \mid y^{(n)}, \Theta) \quad \text{Generative model}$$

$$= \prod_{n=1}^N p(y^{(n)} \mid \phi) p(x_1^{(n)}, x_2^{(n)}, \dots, x_M^{(n)} \mid y^{(n)}, \Theta)$$

$$= \prod_{n=1}^N p(y^{(n)} \mid \phi) \prod_{m=1}^M p(x_m^{(n)} \mid y^{(n)}, \theta_{m,y}) \quad \text{Naïve Bayes}$$

$$= \prod_{n=1}^N \phi^{y^{(n)}} (1 - \phi)^{1-y^{(n)}} \prod_{m=1}^M \theta_{m,1}^{\mathbb{I}(y^{(n)}=1 \wedge x_m^{(n)}=1)} (1 - \theta_{m,1})^{\mathbb{I}(y^{(n)}=1 \wedge x_m^{(n)}=0)}$$
$$\theta_{m,0}^{\mathbb{I}(y^{(n)}=0 \wedge x_m^{(n)}=1)} (1 - \theta_{m,0})^{\mathbb{I}(y^{(n)}=0 \wedge x_m^{(n)}=0)}$$

$$= \phi^{N_{y=1}} (1 - \phi)^{N_{y=0}} \prod_{m=1}^M \theta_{m,1}^{N_{y=1, x_m=1}} (1 - \theta_{m,1})^{N_{y=1, x_m=0}} \theta_{m,0}^{N_{y=0, x_m=1}} (1 - \theta_{m,0})^{N_{y=0, x_m=0}}$$

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

$$y^{(n)} \in \{0,1\}$$

$$\mathbf{x}^{(n)} \in \{0,1\}^M$$

$$\phi \in [0,1]$$

$$\Theta \in [0,1]^{M \times 2}$$

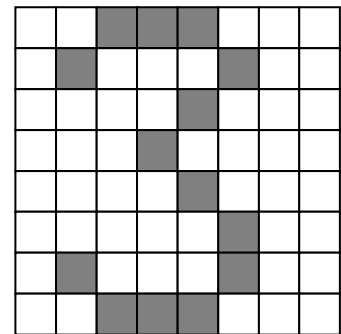
Generative Models

SPAM:

- Class distribution: $Y \sim \text{Bern}(\phi)$
- Class conditional distribution: $X_m \sim \text{Bern}(\theta_{m,y})$
- Naïve Bayes X_i conditionally independent X_j given Y for all $i \neq j$
$$p(X_i, X_j | Y) = p(X_i | Y) p(X_j | Y)$$

Digits:

- Class distribution: $Y \sim \text{Categorical}(\boldsymbol{\phi})$
- Class conditional distribution: $X_m \sim \text{Bern}(\theta_{m,y})$
- Naïve Bayes X_i conditionally independent X_j given Y for all $i \neq j$
$$p(X_i, X_j | Y) = p(X_i | Y) p(X_j | Y)$$



Generative Models with Continuous Features

Iris dataset:

- Class distribution: $Y \sim \text{Bern}(\phi)$
- Class conditional distribution: Multivariate Gaussian $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$
- Naïve Bayes assumption?

Piazza Poll 1

Iris dataset:

- Class distribution: $Y \sim \text{Bern}(\phi)$
- Class conditional distribution: $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$
- Naïve Bayes assumption?

Which of the following pairs of Gaussian class conditional distributions satisfy the Naïve Bayes assumptions? Select ALL that apply.

$$A. \quad \boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$B. \quad \boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}$$

$$C. \quad \boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}, \quad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$$

$$D. \quad \boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}, \quad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Piazza Poll 1

$$\begin{array}{ll} A. \boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, & \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \\ B. \boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, & \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix} \\ C. \boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}, & \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix} \\ D. \boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}, & \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \end{array}$$

Class-conditional Gaussian Distributions

Iris dataset:

- Class distribution: $Y \sim \text{Bern}(\phi)$ (or Categorical)
- Class conditional distribution: $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$

- Naïve Bayes assumption:

- Linear Decision Boundary:

- Quadratic Decision Boundary:

MLE vs MAP vs Generative vs Discriminative