

10-315 Introduction to ML

Probabilistic Models: Generative and Naïve Bayes

Instructor: Pat Virtue

Let's say we gave an exam that had 3 different versions where the probability of different exams was modeled as a **categorical** random variable Y and the scores

from each exam were modeled as $X_{Y=k} \sim \mathcal{N}(\mu_k, \sigma_k^2)$. $\mathcal{D} = \{x^{(i)}, y^{(i)}\}_{i=1}^N$

How many parameters are in this **generative** model?

- A. 2
- B. 3
- C. 6
- D. 9
- E. 3N

Let's say we gave an exam that had 3 different versions where the probability of different exams Y given the scores from the exam x was modeled as a multiclass

logistic regression $p(Y_k = 1 \mid x)$ trained on dataset $\mathcal{D} = \{x^{(i)}, y^{(i)}\}_{i=1}^N$

How many parameters are in this **discriminative** model?

- A. 2
- B. 3
- C. 6
- D. 9
- E. 3N

Which of the following probability models will allow us to sample new data,

 $\mathcal{D} = \left\{x_1^{(i)}, x_2^{(i)}, y^{(i)}\right\}_{i=1}^N$ where x_1, x_2, y are length, weight, and species (cat or dog), respectively? Assume that we have learned all of the corresponding parameters

A) Logistic regression

B) Bernoulli + Gaussian

$$p(Y = cat \mid x_1, x_2, w_1, w_2, b)$$

$$p(Y = cat | \phi)$$

$$p(x_1 | cat, \mu_{cat,1}, \sigma_{cat,1})$$

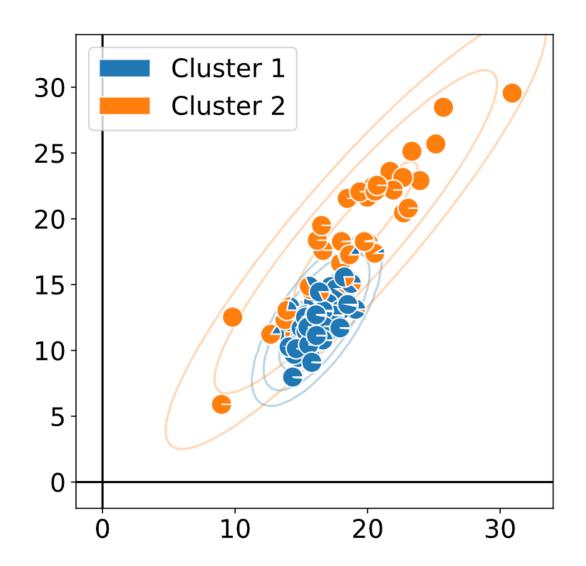
$$p(x_2 | cat, \mu_{cat,2}, \sigma_{cat,2})$$

$$p(x_1 | dog, \mu_{dog,1}, \sigma_{dog,1})$$

$$p(x_2 | dog, \mu_{dog,2}, \sigma_{dog,2})$$

C) None of the above

Sampling Cats and Dogs



B) Bernoulli + Gaussian

$$p(Y = cat \mid \phi)$$

$$p(x_1 \mid cat, \mu_{cat,1}, \sigma_{cat,1})$$

$$p(x_2 \mid cat, \mu_{cat,2}, \sigma_{cat,2})$$

$$p(x_1 \mid dog, \mu_{dog,1}, \sigma_{dog,1})$$

$$p(x_2 \mid dog, \mu_{dog,2}, \sigma_{dog,2})$$

Generative Stories

Generative vs Discriminative Modeling

Discriminative: $p(y \mid x)$

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

Model assumptions vs Data

- Discriminative:
- Generative:

Generative Story Examples

Spam Generation

Discriminative

Generative

Generative + Naive Bayes

Generative Story Examples

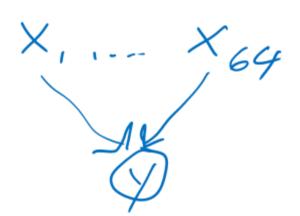
Hand-written digits

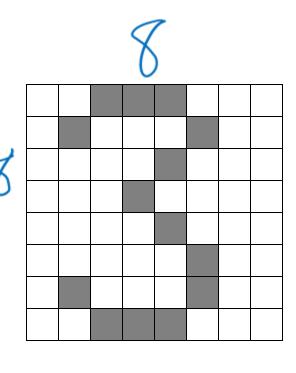
y: Digit class: 0-9

x: Pixels in images

X1 ... 764

Discriminative





Generative

Generative + Naive Bayes

P(X|Y=3)

X, X62

Discriminative and Generative

Two Applications of Bayes Rule

$$p(a \mid b) = \frac{p(b \mid a) p(a)}{p(b)}$$

$$p(\theta \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \theta) p(\theta)}{p(\mathcal{D})} \qquad p(y \mid x) = \frac{p(x \mid y) p(y)}{p(x)}$$

Bayes Rule

Terminology

Posterior Likelihood Prior

$$p(\theta \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \theta) p(\theta)}{p(\mathcal{D})}$$

Bayes Rule

Terminology

Posterior Likelihood Prior

Conditional likelihood Posterior Class conditional Class prior

$$p(\theta \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \theta) p(\theta)}{p(\mathcal{D})} \qquad p(y \mid x) = \frac{p(x \mid y) p(y)}{p(x)}$$

Bayes Rule

Where did the parameters go?!?

$$p(y \mid x) = \frac{p(x \mid y) p(y)}{p(x)}$$

Optimization: Generative vs Discriminative

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

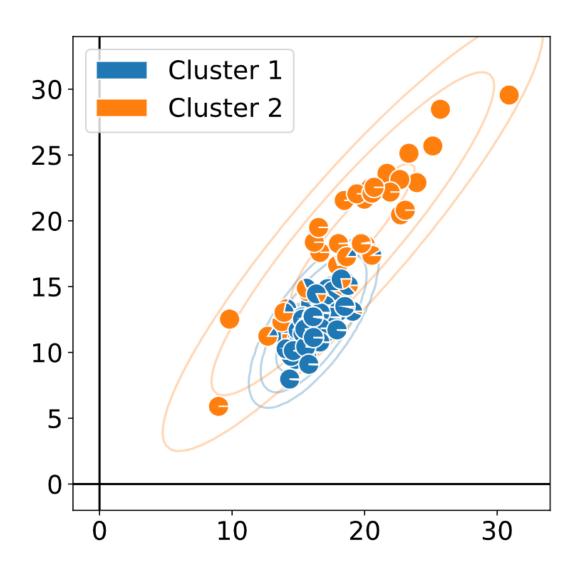
• Learn parameters θ from data

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

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

```
p(y \mid x, \theta_{class}, \theta_{class\ conditional}) \propto p(x \mid y, \theta_{class\ conditional}) p(y \mid \theta_{class})
```

Estimating Cats and Dogs



Bernoulli + Gaussian

$$p(Y = cat \mid \phi)$$

Reminder: Likelihood

Simple example with three i.i.d. samples of $Y \sim Bernoulli(\phi)$

$$\mathcal{D} = \left\{ y^{(i)} \right\}_{i=1}^3$$

Likelihood

$$\begin{split} &p(\mathcal{D};\theta) \\ &= p\big(y^{(1)},y^{(2)},y^{(3)};\phi^{(1)},\phi^{(2)},\phi^{(3)}\big) \text{ Expand notation} \\ &= p\big(y^{(1)},y^{(2)},y^{(3)};\phi\big) & \text{Identically distributed} \\ &= p\big(y^{(1)};\phi\big)\,p\big(y^{(2)};\phi\big)\,p\big(y^{(3)};\phi\big) & \text{Independent} \\ &= \prod_{i=1}^3 p\big(y^{(i)};\phi\big) \end{split}$$

Generative MLE

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

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

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

Generative model

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

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

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

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

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

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

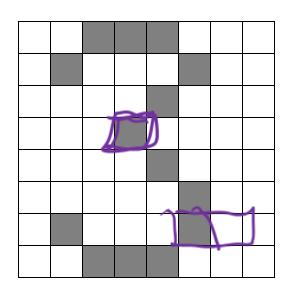
$$\mathbf{\Theta} \in [0,1]^{M \times 2}$$

Naïve Bayes

Multivariate Generative Models

Hand-written digits: How many parameters?

- \blacksquare P(Y)
- P(X | Y = 3)= $p(X_1, X_2, ... X_{64} | Y = 3)$



Naïve Bayes assumption (bag of pixels)

- \blacksquare P(Y)
- P(X | Y = 3) $= p(X_1 | Y = 3)p(X_2 | Y = 3) ... p(X_{64} | Y = 3)$

Conditional Independence and Naïve Bayes

Independence

$$P(A \mid B) = P(A)$$

 $P(B \mid A) = P(B)$

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

Conditional independence

$$A \coprod B$$

$$P(A \mid B,C) = P(A \mid C) \qquad P(A,B \mid C) = P(A \mid C)P(B \mid C)$$

$$P(B \mid A,C) = P(B \mid C)$$

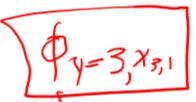
$$P(A, B \mid \underline{C}) = P(A \mid \underline{C})P(B \mid \underline{C})$$

Naïve Bayes assumption

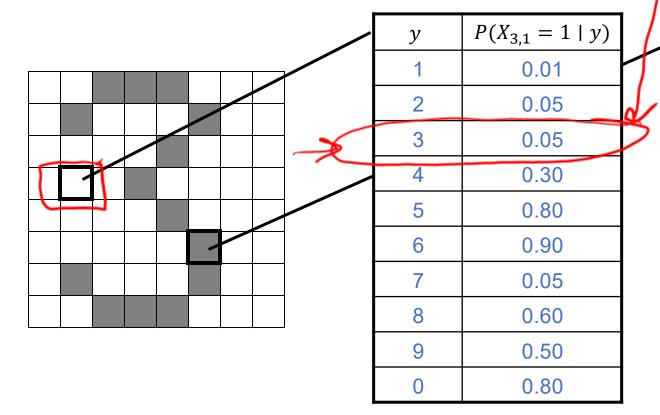
$$\forall_{j \neq k}$$
 $x_i \perp \perp x_j \mid Y$

$$\frac{1}{\sqrt{y}} = 1$$

Naïve Bayes for Digits



у	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



у	$P(X_{5,5} = 1 \mid y)$
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

Recitation Exercise

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

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

Naïve Bayes MLE

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

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

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

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

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

$$= \prod_{n=1}^N p(\mathbf{y}^{(n)} \mid \mathbf{\phi}) \prod_{m=1}^M p(\mathbf{x}_m^{(n)} \mid \mathbf{y}^{(n)}, \theta_{m,y}) \quad \text{Na\"ive Bayes}$$

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

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

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

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

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

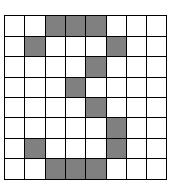
Generative Models

SPAM:

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

Digits:

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



Generative Models with Continuous Features

Iris dataset:

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

Iris dataset:

- Class distribution: $Y \sim Bern(\phi)$
- Class conditional distribution: $X \sim \mathcal{N}(\mu_y, \Sigma_y)$
- Naïve Bayes assumption?

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

A.
$$\mu_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \Sigma_{y=0} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mu_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \Sigma_{y=1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

B. $\mu_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \Sigma_{y=0} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mu_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \Sigma_{y=1} = \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}$

C. $\mu_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \Sigma_{y=0} = \begin{bmatrix} 1 & -1 \\ 1 & 2 \end{bmatrix}, \quad \mu_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \Sigma_{y=1} = \begin{bmatrix} 1 & -1 \\ 1 & 2 \end{bmatrix}$

D. $\mu_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \Sigma_{y=0} = \begin{bmatrix} 1 & -1 \\ 1 & 2 \end{bmatrix}, \quad \mu_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \Sigma_{y=1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

Class-conditional Gaussian Distributions

Iris dataset:

- Class distribution: $Y \sim Bern(\phi)$ (or Categorical)
- lacktriangle Class conditional distribution: $m{X} \sim \mathcal{N}(m{\mu}_y, m{\Sigma}_y)$
- Naïve Bayes assumption:
- Linear Decision Boundary:
- Quadradic Decision Boundary:

Generative + MAP

MLE vs MAP vs Generative vs Discriminative

Maximum likelihood estimation

Maximum a posteriori estimation

Descriminative

Models only the conditional likelihood, $\prod p(y \mid x, \theta)$

 $p(\mathcal{D}^* \mid \theta)$ * conditional likelihood $= \prod_{i=1}^{N} p(y^{(i)} \mid x^{(i)}, \theta)$

 $p(\theta) p(\mathcal{D}^* \mid \theta)$ * conditional likelihood $= p(\theta) \prod_{i=1}^{N} p(y^{(i)} \mid x^{(i)}, \theta)$

Generative

Models the full joint likelihood, $\prod p(x, y \mid \theta)$

 $p(\mathcal{D} \mid \theta)$ * actual likelihood $= \prod_{i=1}^{N} p(x^{(i)}, y^{(i)} \mid \theta)$ $= \prod_{i=1}^{N} p(x^{(i)} \mid y^{(i)}, \theta) p(y^{(i)} \mid \theta)$ classconditional prior

 $p(\theta) p(\mathcal{D} \mid \theta)$ * actual likelihood $= p(\theta) \prod_{i=1}^{N} p(x^{(i)}, y^{(i)} \mid \theta)$ $= p(\theta) \prod_{i=1}^{N} p(x^{(i)} \mid y^{(i)}, \theta) p(y^{(i)} \mid \theta)$ = class- class conditional prior

Reminder: Prior Distributions for MAP

If the prior $p(\theta)$ is uniform, then MLE and MAP are the same!

$$p(\mathcal{D} \mid \phi) p(\phi) = \prod_{i} p(y^{(i)} \mid \phi) p(\phi) = \phi^{N_{y=1}} (1 - \phi)^{N_{y=0}} p(\phi)$$

Conjugate priors: when the prior and the posterior distributions are in the same family

Bernoulli likelihood with a **Beta prior** has **Beta posterior**

Categorical likelihood with a <u>Dirichlet prior</u> has <u>Dirichlet posterior</u>

Gaussian likelihood with a Gaussian prior has Gaussian posterior

https://www.desmos.com/calculator/kr7m2m6cf7