

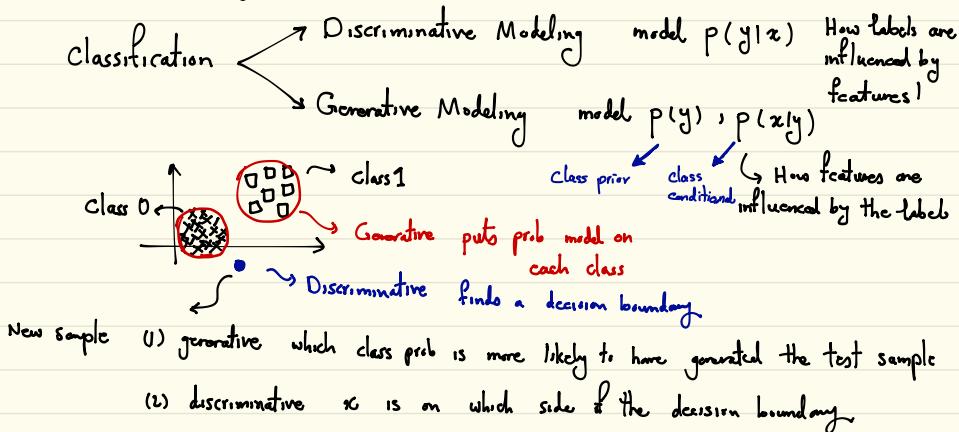
Generative Modeling and Naive Bayes :

DS 4400
Instructor: Ehsan Elhamifar

Given input features/attributes $\{x^1, x^n\}$ and discrete outputs/labels $\{y^1, y^n\}$,
find a model to predict y given x

$$x^i = \begin{pmatrix} x_1^i \\ x_2^i \end{pmatrix} \rightarrow \text{feature/attribute/variable}$$
$$y^i \in \{0, 1, \dots, L-1\}$$

We discussed two ways to address the problem



Generative Modeling specify $p(y), p(x|y)$ and using Bayes' rule

$$\underset{j \in \{0, 1, \dots, L-1\}}{\operatorname{argmax}} p(y=j | x) = \underset{j \in \{0, 1, \dots, L-1\}}{\operatorname{argmax}} p(y=j) p(x|y=j)$$

- * We need to learn parameters of generative model via maximum likelihood, given a training dataset $D = \{(x^1, y^1), \dots, (x^n, y^n)\}$

* what are parameters to learn?

$$p(y=j) \triangleq \theta_j^y$$

$$p(x=\bar{x} | y=j) \triangleq \theta_{\bar{x}|j}^{xy}$$

$$\rightarrow \Theta \triangleq \{\theta_j^y, \theta_{\bar{x}|j}^{xy}; \forall j, \forall \bar{x}\}$$

Ex) Spam detection

$$x_i = \begin{pmatrix} \text{'Free'} \\ \text{'CSPS'} \\ \text{'Call now'} \\ \text{'Nigerian'} \\ \vdots \end{pmatrix} \quad y \in \{\text{'spam'}, \text{'non-spam'}\}$$

Need to specify $\left\{ \theta_0^y = p(y=0) = p(\text{'spam' class}) \right.$

$$\theta_1^y = p(y=1) = p(\text{'not spam' class})$$

$$\theta_{\bar{x}|0}^{xy} = p(x=\bar{x} | y=0) = p(x=\bar{x} | \text{'spam'})$$

$$\theta_{\bar{x}|1}^{xy} = p(x=\bar{x} | y=1) = p(x=\bar{x} | \text{'not spam'})$$

$$* \text{ MLE} \quad L(\theta) = P_\theta(D) = P_\theta(x_1^y, \dots, x_N^y) \stackrel{iid}{=} \prod_{i=1}^N P_\theta(x_i^y, y_i) = \prod_{i=1}^N P_\theta(y_i^y) \prod_{i=1}^N P_\theta(x_i^y | y_i)$$

$$= \prod_{i=1}^N \prod_{j=0}^{L-1} \underbrace{P_\theta(y^y=j)}_{=\theta_j^y} \times \prod_{i=1}^N \prod_{j=0}^{L-1} \prod_{\bar{x}} \underbrace{P_\theta(x^i=\bar{x} | y^i=j)}_{\theta_{\bar{x}|j}^{xy}} \frac{1}{\# \text{Samples in } j}$$

$$\rightarrow l(\theta) = \log P_\theta(D) = \sum_{i=1}^N \sum_{j=0}^{L-1} 1(y^i=j) \log \theta_j^y + \sum_{i=1}^N \sum_{j=0}^{L-1} \sum_{\bar{x}} 1(x^i=\bar{x}, y^i=j) \log \theta_{\bar{x}|j}^{xy}$$

$$= \sum_{j=0}^{L-1} \left(\sum_{i=1}^N 1(y^i=j) \right) \log \theta_j^y + \sum_{j=0}^{L-1} \sum_{\bar{x}} \left(\sum_{i=1}^N 1(x^i=\bar{x}, y^i=j) \right) \log \theta_{\bar{x}|j}^{xy}$$

$$\text{Recall: } \sum_{j=0}^{L-1} \alpha_j \log \theta_j \quad \text{if } \theta_j > 0, \sum_{j=0}^{L-1} \theta_j = 1 \rightarrow \hat{\theta}_j = \frac{\alpha_j}{\sum_{l=0}^{L-1} \alpha_l} \Rightarrow \hat{\theta}_j^y = \frac{\sum_{i=1}^N 1(y^i=j)}{\sum_{i=1}^N \sum_{j=0}^{L-1} 1(y^i=j)} = \frac{\# \text{Samples in } j}{N}$$

ex) 1520 spam emails, 1933 non-spam emails

$$\theta_0^y = \frac{1520}{1520+1933}, \quad \theta_1^y = \frac{1933}{1520+1933}$$

$$\text{Similarly, } \hat{\theta}_{\bar{x}|j}^{xy} = \frac{\sum_{i=1}^N 1(x^i=\bar{x}, y^i=j)}{\sum_{\bar{x}} \sum_{i=1}^N 1(x^i=\bar{x}, y^i=j)} = \frac{\# \text{Samples in class } j \text{ where feature is } \bar{x}}{\# \text{Samples in class } j}$$

* How about specifying $\theta_{\bar{x}_j}^{x_j} \triangleq p(x=\bar{x}|y=j)$?

→ Similar to before, for each possible value of $x = \bar{x}$, we need to count the fraction of training samples in class j whose feature is \bar{x} !

Ex) $x = \begin{pmatrix} \text{'Free'} \\ \text{'CAPS'} \\ \text{'RBJ'} \end{pmatrix}$ for $\bar{x} = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$ count the fraction of emails in

class 'spam' where $x^i = \bar{x} \rightarrow \theta_{\bar{x}, i}^{x_i}$

class 'nonspam' " " " $\rightarrow \theta_{\bar{x}, i}^{x_i}$

* Challenge if x 's have K features $\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_K \end{pmatrix}$ where each feature takes

in discrete values \rightarrow we have m^K possible \bar{x} 's !!

(in the case of spam detection example, 2^K) \hookrightarrow Large !

We need to have a very large training set with multiple examples per each of m^K cases, to get a valid estimation of $\theta_{\bar{x}, j}^{x_j}$!!

$L \times m^K$ pars Impossible! In the case of limited training we do not have any example for which $x = \bar{x}$ for many \bar{x} 's!

ex) $\bar{x} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$ \rightsquigarrow have not seen this before, so estimate $\theta_{\bar{x}, 1}^{x_1} = \theta_{\bar{x}, 1}^{x_2} = \dots = 0$

but makes sense to have $\theta_{\bar{x}, 1}^{x_1} \approx 1$, $\theta_{\bar{x}, 1}^{x_2} \approx 0$!!

$$\arg \max_{j=1, \dots, L-1} p(y=j|x^t) = \arg \max_{j=1, \dots, L-1} p(x^t|y=j) p(y=j) = \arg \max_{j=1, \dots, L-1} \theta_{\bar{x}, j}^{x_j} \times \theta_j^y = 0, \forall j$$

← can't train

Drawbacks \rightarrow exponential number of parameters m^K in the number of features

\rightarrow many unseen examples among $m^K \rightarrow$ poor estimate

Fix "Naive Bayes" assumption

features are independent from each other given a class

$$\begin{aligned} p(x' | y=j) &= p(x'_1 | y=j) \times \dots \times p(x'_k | y=j) \\ &= \prod_{l=1}^k p(x'_{jl} | y=j) \quad \leadsto \text{Conditional Independence} \end{aligned}$$

Ex) In the case of Spam detection, if we know an email is a spam, the chance of seeing 'Free' does not change if we see words in 'CAPS'

Conditional Independence $p('Free'=1, 'CAPS'=1 | \text{Spam}) = p('Free'=1 | \text{Spam}) p('CAPS'=1 | \text{Spam})$

However, this independence does not hold if we do not condition

$$p('Free'=1, 'RB3'=1) \neq p('Free'=1) p('CAPS'=1)$$

If we see 'free', that increases the chance of email being 'spam' which consequently increases the chance of seeing 'CAPS'

So, we need to specify $p(x_n = \bar{x}_n | y=j) \triangleq \theta_{\bar{x}_n, j}^{x_n | y}$ for all $j = 1, \dots, L$
and all x_n values. If x_n takes m values $\rightarrow L \times m \times K$ parameters
 \Rightarrow significant decrease from $L \times m^K$ to $L \times m \times K$

Remember that for MLE:

$$P_\theta(D) = \prod_{i=1}^N P_\theta(y_i^i) \underbrace{\prod_{i=1}^N P_\theta(x_i^i | y_i^i)}_{\leftarrow \prod_{i=1}^N \prod_{l=1}^K P_\theta(x_i^i | y_i^i)} = \prod_{i=1}^N \prod_{l=1}^K \prod_{j=0}^{L-1} \underbrace{P(x_i^i = \bar{x}_i, y_i^i = j)}_{\stackrel{\text{def}}{=} \theta_{\bar{x}_i | j}^{x_i^i | y_i^i}}$$

$$1(x_i^i = \bar{x}_i, y_i^i = j)$$

$$\rightarrow l(\theta) = \dots + \sum_{l=1}^K \sum_{j=0}^{L-1} \sum_{\bar{x}_i} \left(\sum_{i=1}^N 1(x_i^i = \bar{x}_i, y_i^i = j) \right) \log \theta_{\bar{x}_i | j}^{x_i^i | y_i^i}$$

$$\rightarrow \hat{\theta}_{\bar{x}_i | j}^{x_i^i | y_i^i} = \frac{\sum_{i=1}^N 1(x_i^i = \bar{x}_i, y_i^i = j)}{\sum_{\bar{x}_i} \sum_{i=1}^N 1(x_i^i = \bar{x}_i, y_i^i = j)} = \frac{\# \text{Samples in class } j \text{ where } l_{\text{th}} \text{ feature is } \bar{x}_i}{\# \text{Samples in class } j}$$

$$\text{Ex) } \hat{\theta}_{111}^{x_1^1 | y_1^1} = \frac{\# \text{ emails in class spam that contain word 'Free'}}{\# \text{ emails in class spam}}$$

Ex) Assume in none of training examples, the word 'call now' appears, but words such as 'free', 'CAPS', ... appear $\Rightarrow \hat{\theta}_{\bar{x}_3 | j}^{x_3^3 | y_3^3} = 0, \forall \bar{x}_3 \in \{\dots\}, j \in \{\dots\}$

Assume in test email, we have word 'callnow' as well as 'free', 'CAPS', ...

Since $P(x_3^3 | y_3^3) = 0$

$$P(y_3^3 | x^t) \propto P(x^t | y_3^3) P(y_3^3) = P(x_1^1 | y_3^3) \cdots P(x_m^m | y_3^3) P(y_3^3) = 0$$

To overcome this zero prob effect, we assume all $\hat{\theta}_{\bar{x}_i | j}^{x_i^i | y_i^i} \neq 0$ by defining

$$\hat{\theta}_{\bar{x}_i | j}^{x_i^i | y_i^i} = \frac{\sum_{i=1}^N 1(x_i^i = \bar{x}_i, y_i^i = j) + t}{\sum_{\bar{x}_i} \sum_{i=1}^N 1(x_i^i = \bar{x}_i, y_i^i = j) + mt}$$

when x_i takes m values.

↪ Laplace smoothing

$\rightarrow t$ is a hyperparam which can be tuned using cross validation
 $\frac{1}{10^3}, \frac{1}{10^6}$

In the MLE framework,

$$\theta_{\bar{x}|j}^{x_k|y} = \frac{\sum_{i=1}^n 1(x_k^i = \bar{x}, y^i=j)}{\sum_{i=1}^n \sum_{l=1}^m 1(x_k^i = l, y^i=j)}$$

↳ fraction of emails in class j , whose k -th feature is \bar{x}

* Using Naive Bayes, we reduced the # parameters from exponential to linear in the # features. However, it still can be the case that we do not observe any training data for some $\theta_{\bar{x}|j}^{x_k|y}$

Ex) Bag of Words model for Documents

$$\text{Doc } i = \begin{pmatrix} 105 & \rightarrow a \\ 67 & \rightarrow \text{car} \\ 28 & \rightarrow \text{rp} \\ 23 & \rightarrow \text{cor} \\ 0 & \rightarrow \text{truck} \end{pmatrix}$$

each feature can take length of document # of values
 $\{1, 2, \dots, |\text{Doc}.x|\}$
 $\rightarrow |\text{Doc}|^K$ possible feature combinations / # poss
 $\xrightarrow{\text{NB}} K |\text{Doc}|$ possible parameters

Car vs politics classification of documents assume in none of training samples in class car, there is word 'truck' $\rightarrow p(\text{'truck'} \neq 0 | y = \text{'car'}) = 0$
 Now, in your test if word 'truck' appear in a document of class 'car', then $p(x | \text{'car'}) = \prod_k p(x_k | \text{'car'}) = 0 \quad ||| \rightarrow$ You do not choose class 'car' for true class!

Solution Laplace Smoothing
 in the past

assume we have seen each of the words/features t times

$$\theta_{\bar{x}|j}^{x_k|y} = \frac{\sum_{i=1}^n 1(x_k^i = \bar{x}, y^i=j) + t}{\sum_{i=1}^n \sum_{l=1}^m 1(x_k^i = l, y^i=j) + mt}$$

Example: Spam detection with $x = \begin{pmatrix} 1 & \text{freq} \\ 1 & \text{CAPS} \\ 1 & \text{call num} \end{pmatrix}$ Spam: $y=1$, non-spam: $y=0$

$$D = \left\{ \left(\left(\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, 1 \right), 1 \right), \left(\left(\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, 1 \right), 1 \right), \left(\left(\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, 1 \right), 0 \right), \left(\left(\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, 0 \right), 0 \right) \right\}$$

$$\theta_0^y = \frac{3}{7}, \quad \theta_1^y = \frac{4}{7}$$

$$\theta_{1,0}^{y=1} = \frac{2}{3}, \quad \theta_{1,0}^{y=0} = \frac{1}{3}$$

$$\theta_{0,0}^{y=1} = \frac{3}{3}, \quad \theta_{0,0}^{y=0} = 0$$

$$\theta_{0,1}^{y=1} = \frac{2}{3}, \quad \theta_{0,1}^{y=0} = \frac{1}{3}$$

$$\theta_{0,1}^{y=1} = \frac{2}{4}$$

$$\theta_{0,1}^{y=0} = 1$$

$$\theta_{1,1}^{y=1} = \frac{3}{4}$$

$$\theta_{1,1}^{y=0} = \frac{1}{4}$$

$$x^n = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \Rightarrow \underset{j=0,1}{\operatorname{argmax}} p(y=j|x^n) = \underset{j=0,1}{\operatorname{argmax}} p(x^n|y=j) p(y=j) = \underset{j=0,1}{\operatorname{argmax}} p(x^n|y=j) p(x^n|y=j) p(x^n|y=j)$$

$$\underset{j=0,1}{\operatorname{argmax}} \theta_{1,0}^{y=1} \theta_{1,0}^{y=0} \theta_{0,0}^{y=1} \theta_{0,0}^{y=0}$$

$$j=0 \rightarrow \frac{1}{3} \times \cancel{\frac{\epsilon}{3}} \times \frac{1}{3} \times \frac{3}{7}$$

$$j=1 \rightarrow \frac{2}{4} \times \cancel{\frac{\epsilon}{4}} \times \frac{1}{4} \times \frac{4}{7}$$

with $\epsilon \Rightarrow$ cannot distinguish

$$\text{with } \epsilon \quad \frac{1}{21} < \frac{1}{4} \Rightarrow j^* = 1$$

$$x^n = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \Rightarrow \underset{j=0,1}{\operatorname{argmax}} p(y=j|x^n) = \underset{j=0,1}{\operatorname{argmax}} p(x^n|y=j) p(y=j) = \underset{j=0,1}{\operatorname{argmax}} p(x^n=1|y=j) p(x^n=0|y=j)$$

$$= \underset{j=0,1}{\operatorname{argmax}} \theta_{1,0}^{y=1} \theta_{0,0}^{y=1} \theta_{1,0}^{y=0} \theta_{0,0}^{y=0} \theta_j^y$$

$$j=0 \rightarrow \frac{1}{3} \times 1 \times \frac{1}{3} \times \frac{3}{7} = \frac{1}{21}$$

$$j=1 \rightarrow \frac{2}{4} \times 1 \times \frac{1}{4} \times \frac{4}{7} = \frac{1}{14}$$

$\rightarrow j^* = 1$

Gaussian Naive Bayes

The previous examples, each feature takes a discrete set of values binary or m values
what if features are continuous $x = \begin{pmatrix} \text{time} \\ \text{location} \end{pmatrix}$? ex in weather classification

Model $p(x_k | y=j)$ using a continuous distribution

$$x_k | y=j \sim \mathcal{N}(\mu_{kj}, \sigma_{kj}^2)$$

$$\rightarrow \hat{\mu}_{kj} = \frac{\sum_{i=1}^N x_k^i \mathbb{1}(y^i=j)}{\sum_{i=1}^N \mathbb{1}(y^i=j)}$$

$$\rightarrow \sigma_{kj}^2 = \frac{\sum_{i=1}^N (x_k^i - \hat{\mu}_{kj})^2 \mathbb{1}(y^i=j)}{\sum_{i=1}^N \mathbb{1}(y^i=j)}$$