

9. Naive Bayes

- Naive Bayes
- Generative vs discriminative models

Many slides borrowed from Prof. Minh Hoai Nguyen's previous course offering

Bayes classifier is computationally expensive in general

$$\hat{y}_{\text{Bayes}}(x) = \underset{\hat{y} \in \{1, \dots, K\}}{\operatorname{argmax}} \Pr(y(x) = \hat{y}|x)$$

Bayes' rule

$$= \underset{\hat{y} \in \{1, \dots, K\}}{\operatorname{argmax}} \frac{\Pr(x|y(x) = \hat{y})}{\underbrace{\Pr(x|y(x) = \hat{y})}_{\text{Conditional likelihood}}} \underbrace{\Pr(y(x) = \hat{y})}_{\text{prior}}$$

- Assume K classes, D attributes, each taking N values.
- Sampling complexity for $y|x$ is $O(KN^D)$
- Problem: no modeling assumption on \hat{y} = high sampling complexity

Naive Bayes

Assumption necessary

Assumption: each attribute is independent

given label

$$\Pr(x|y) = \prod_k \Pr(x_k|y)$$

Example: Will it rain today?

- is it cold?
- is it cloudy?
- are my joints aching?
- did it rain yesterday?



• is temperature low?

Is this assumption reasonable? (Given it will rain, are cold, cloudy independent?)

Naive Bayes

Assumption: each attribute is independent

$$\Pr(x|y) = \prod_k \Pr(x_k|y)$$

~~Example: Will it rain today?~~

- is it cold?
- is it cloudy?
- are my joints aching?
- did it rain yesterday?

Is this assumption reasonable? (Given it will rain, are cold, cloudy independent?)

Maybe not, but it's useful and works fine in practice

“All models are wrong, but some are useful” – George Box

Naive Bayes

$$\hat{y}_{\text{Bayes}}(x)$$

$$= \\ \stackrel{\text{Bayes' rule}}{=}$$

Naive Bayes' assumption

From $O(KN^D)$ to $O(KDN)$

NP

IR

$$\operatorname{argmax}_{\hat{y} \in \{1, \dots, K\}} \Pr(y(x) = \hat{y} | x)$$

$$\operatorname{argmax}_{\hat{y} \in \{1, \dots, K\}} \Pr(x|y(x) = \hat{y}) \Pr(y(x) = \hat{y})$$

$$\operatorname{argmax}_{\hat{y} \in \{1, \dots, K\}} \prod_{i=1}^D \Pr(x_i|y(x) = \hat{y}) \Pr(y(x) = \hat{y})$$

humid cold

Example: Bayes classifier, no structural assumption

$$\hat{y}_{\text{Bayes}}(x) = \underset{\hat{y} \in \{1, \dots, K\}}{\operatorname{argmax}} \Pr(x|y(x) = \hat{y}) \Pr(y(x) = \hat{y})$$

- Will it rain?
- Collect historical data

rain	yes	no	yes	no	no	yes
A	1	1	-1	1	1	-1
B	99	41	4	18	3	12
C	◊	♡	♠	♣	♣	♡

To see every value of $\Pr(x_A, x_B, x_C | y(x) = \hat{y})$ at least once, require

$$\underbrace{2}_{\# \text{ classes}} \times \underbrace{2}_{\# \text{ values for A}} \times \underbrace{100}_{\# \text{ values for B}} \times \underbrace{4}_{\# \text{ values for C}} \text{ observations}$$

$$f(x) = \sum_i d_i f_i(x)$$

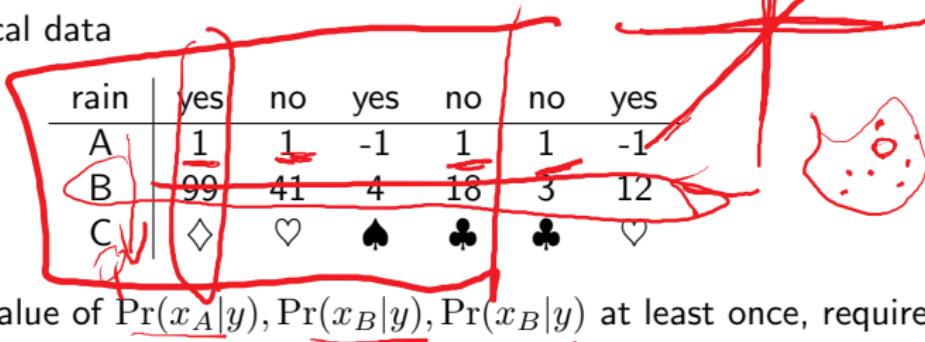
Naive Bayes classifier

B follows $\frac{L_1}{100}$

$$\hat{y}_{\text{Bayes}}(x) = \underset{\hat{y} \in \{1, \dots, K\}}{\operatorname{argmax}} \prod_{i=1}^D \Pr(x_i | y(x) = \hat{y}) \Pr(y(x) = \hat{y})$$

- Will it rain?
- Collect historical data

$$\begin{aligned} & \text{rain} | A=1 \\ &= \begin{cases} \frac{1}{4} \text{ yes} \\ \frac{3}{4} \text{ no} \end{cases} \end{aligned}$$



rain	yes	no	yes	no	no	yes
A	1	1	-1	1	1	-1
B	99	41	4	18	3	12
C	◊	♥	♠	♣	♣	♥

To see every value of $\Pr(x_A|y)$, $\Pr(x_B|y)$, $\Pr(x_C|y)$ at least once, require

$$\underbrace{2}_{\# \text{ classes}} \times \left(\underbrace{2}_{\# \text{ values for A}} + \underbrace{100}_{\# \text{ values for B}} + \underbrace{4}_{\# \text{ values for C}} \right) \text{ observations}$$

logistic regression

generative

discriminative

linear regression

$y \sim N(\theta^T x)$?

nearest neighbors

?

SVM
margin

?

Generative vs discriminative models

Bayes

naive Bayes

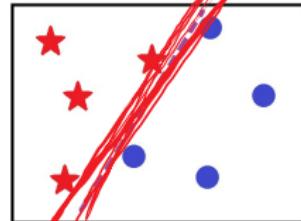
c^T

decision trees

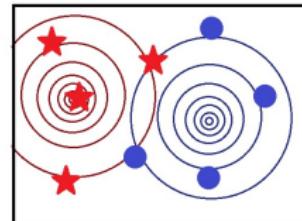
✓

Discriminative vs generative models

- Discriminative models focus on task
 - Goal: find the margin that minimizes misclassification rate
 - e.g.: logistic regression, regression trees, K-nearest neighbors
- Generative models focus on interpretability
 - First build a probabilistic model, then find MLE / MAP
 - e.g.: linear regression, naive Bayes, Gaussian mixture models

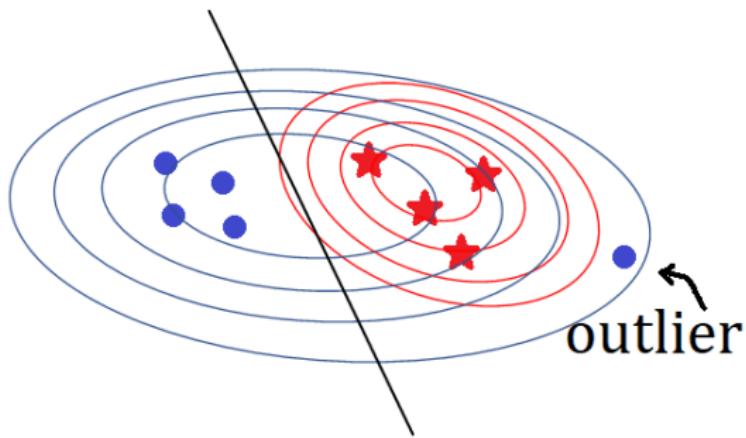


discriminative



generative

Generative models may not be robust



Generative vs discriminative models

- In classification, the goal is to produce a rule $\hat{y}(x) \approx y(x) \in \{1, \dots, K\}$
- A **generative classifier** aims to model the entire universe
 - Assume a model $\Pr(x|y)$ and prior $\Pr(y)$
 - Produces $\hat{y}(x) = \operatorname{argmax}_{y \in \{1, \dots, K\}} \Pr(x|y)\Pr(y)$
 - Requires a lot of training samples
 - Interpretable
- A **discriminative classifier** models only the decision rule
 - Directly models $f(x, k) = \Pr(y = k|x)$ for all k
 - Produces $\hat{y}(x) = \operatorname{argmax}_{k \in \{1, \dots, K\}} f(x, k)$
 - Requires much less training samples
 - Hard to interpret

Example: Logistic regression for binary classification

$$\theta^* := \operatorname{argmax}_{\theta} \sum_{i=1}^m \log(y_i x_i^T \theta), \quad y(x) = \operatorname{sign}(x^T \theta^*)$$

Generative or discriminative?

Example: Logistic regression for binary classification

$$\theta^* := \operatorname{argmax}_{\theta} \sum_{i=1}^m \log(y_i x_i^T \theta), \quad y(x) = \operatorname{sign}(x^T \theta^*)$$

Generative or discriminative? Ans: Discriminative!

Remember modeling assumption:

$$\log \left(\frac{\Pr(y(x) = 1)}{\Pr(y(x) = -1)} \right) = x^T \theta$$

No assumptions on $\Pr(x|y)$!

Example: Linear regression

$$\theta^* := \operatorname{argmin}_{\theta} \sum_{i=1}^m (x_i^T \theta - y_i)^2 + \|\theta\|_2^2, \quad y(x) = \operatorname{sign}(x^T \theta^*)$$

Generative or discriminative?

Example: Linear regression

$$\theta^* := \underset{\theta}{\operatorname{argmin}} \underbrace{\sum_{i=1}^m (x_i^T \theta - y_i)^2 + \|\theta\|_2^2}_{f(\theta)}, \quad y(x) = \mathbf{sign}(x^T \theta^*)$$

Generative or discriminative? Ans: Generative! <-- depending on interpretation
Under Gaussian modeling assumption:

$$y \sim \mathcal{N}(\theta^T x, 1), \quad \theta \sim \mathcal{N}(0, I) \Rightarrow f(\theta) \propto \Pr(\theta | \{x_i, y_i\})$$

Summary

- Up until now, we only considered discriminative classifiers
 - Logistic regression
 - Thresholded linear regression
- Naive Bayes: first generative approach
- Key assumption: each features independent given label
 - Reasonable?
 - Maybe not accurate, but effective
 - Considerably reduces complexity