

DATA 5610

①

- Alpaydm Chp 5: Multivariate Methods
 - 5.1 Multivariate Data
 - 5.2 Parameter Estimation
 - 5.3 Estimation of Missing Values
 - 5.4 Multivariate Normal Distribution
 - 5.5 Multivariate Classification
 - 5.6 Tuning complexity
 - 5.7 Discrete Features
- ISLR 2
 - Chp 4: Classification
- ~~Willett~~ Willett Chp 5 Naive Bayes Classifier

Naïve Bayes Classifier

models the joint pdf $P(x, y)$

- good intro to "generative" classification

- Setup:

• Given Data = $[(x^{(1)}, y), \dots, (x^{(n)}, y)]$

$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)}) \quad \text{so } x^{(i)} \in \mathbb{R}^d, \quad y_i \in \mathcal{Y}$$

where $\mathcal{Y} = \{1, \dots, m\}$

• Assume a family of distributions P_θ s.t.

$$P_\theta(x, y) = P_\theta(x|y) P(y)$$

$$= P_\theta(x_1|y) \dots P_\theta(x_d|y) P_\theta(y)$$

* This key assumption is why it is called "Naïve"

- Let $(x^{(1)}, y_1), \dots, (x^{(n)}, y_n) \sim P_\theta$ iid for some θ

(If $(x, y) \sim P_\theta$, then x_1, \dots, x_d are independent given Y)

("Conditional Independ. Assumption")

→ "Independent Feature Model"

• x_1, \dots, x_d the "features"

- Goal: For some new $x \in \mathbb{R}^d$, predict its y (classify)

- Algorithm:

• Estimate θ (MLE, MAP, etc.)
from Data

• Compute $\hat{y} = \underset{y \in Y}{\operatorname{argmax}} \underset{\hat{\theta}}{P(y | x)}$
↑ estimated

$$\hat{y} = \underset{(y)}{\operatorname{argmax}} \left[\frac{P_{\hat{\theta}}(x|y) P_{\hat{\theta}}(y)}{P_{\hat{\theta}}(x)} \right] \quad \text{using Bayes' Rule}$$

— Does not depend on denom., so

$$\hat{y} = \underset{y}{\operatorname{argmax}} \left[P_{\hat{\theta}}(x|y) P_{\hat{\theta}}(y) \right] \quad \text{will select the same } \hat{y}$$

$$= \underset{y}{\operatorname{argmax}} P_{\hat{\theta}}(x_1|y) \dots P_{\hat{\theta}}(x_d|y) P_{\hat{\theta}}(y)$$

— The "Bayes" part of NBC is b/c it is a Bayes' Estimator

What is it used for?

⑤

- Peter Norvig: often start with NBC and don't need to improve on it
- classification, text, sentiment
- NLP: Natural Language Processing (getting the computer to "understand" humans)
- Email spam detection (your email prolly uses NBC)
- "Good" or "bad" news
- Predict which direction tweets on Twitter will influence election or referendum
- Determine if the tweets come from a Russian bot

Political Speeches (Wilmott example)

6

- Using Bayes Rule

$$P(\text{Politician is left wing} \mid \text{used the word "Comrade"}) =$$

$$\frac{P(\text{used the word "Comrade"} \mid \text{left wing}) P(\text{left wing})}{P(\text{uses the word "comrade"})}$$

- So if a politician uses the word "comrade" in a speech, we can calculate the prob. of him/her being left-wing

- Apply to whole phrases and entire speeches not just single words

- Suppose a speech contains the phrase "Property is theft, Comrade"
and we want to know if they are left or right wing

$$P(\text{Left} \mid \text{"Property is theft, Comrade"})$$

and similar for

$$P(\text{Right} \mid \text{"Property is theft, Comrade"})$$

- Now comes the "Naïve" part

$$P(\text{"Property is that, Comrade"} | \text{Left}) =$$

$$P(\text{"property"} | \text{Left}) * P(\text{"that"} | \text{Left}) * P(\text{"Comrade"} | \text{Left})$$

- Drop the "is" because it's a stop word

In symbols

10

- The data, text we want to classify is called X as a vector consisting of x_d for $1 \leq d \leq D$ so that there are D words in the text

- We want to find

$$P(c_k | x)$$

for each of the K classes (political persuasions) c_k

- use MLE or MAP to estimate

Bayes Rule tells us that

$$P(C_k | x) = \frac{P(x | C_k) P(C_k)}{P(x)}$$

→ When the features are independent this simplifies to

$$P(C_k) \prod_{d=1}^D P(x_d | C_k)$$

→ This is what we compare for each class

→ We will get $P(x_d | C_k)$ from the data set

- Finally, because multiplying a bunch of possibly (very) small numbers can cause round-off error we take the log transform

$$\ln [P(C_k)] + \sum_{d=1}^D \ln [P(x_d | C_k)]$$

- $\ln [P(C_k)]$ is the prior
 - ~ We can ignore it (Frequentist)
 - Estimate it from data (Empirical Bayes)
 - or model it subjectively (Bayesian)

Wilmott's Example

13

- Churchill "Beaches" speech
- JFK inaugural address
- Benn: Mardon speech as MP
- Thatcher: "The Lady's not for turning" speech
- May: Syria speech
- Corbyn: Post-Brexit-Referendum speech
- Trump: state of the Union speech

* Predict MLK's "I have a dream" speech

* What is your prediction?