

Lec 8

Tuesday, September 24, 2019 11:10

Recap: Naive Bayes

Naive assumption: individual features are independent of one another given label value

or: just the conditional densities factor:

$$P(X=x|Y=j) = f_j(x) = \prod_{k=1}^p f_{jk}(x_k)$$

\uparrow conditional density of $X|Y=j$
 \uparrow conditional density of $X_k|Y=j$

Then:

$$P(Y=j|X) \stackrel{\text{Bayes Law}}{=} \frac{\pi_j f_j(x)}{\sum_k \pi_k f_k(x)}$$

π_j - fraction of label j in pop

$$\stackrel{\text{Naive assumption}}{=} \frac{\pi_j \cdot \prod_{k=1}^p f_{jk}(x_k)}{\sum_k \pi_k \cdot \prod_{l=1}^p f_{kl}(x_l)}$$

Naive Bayes Classifier:

- estimate each π_j , f_{jk}
- plug in & max over j

Bag of words features

BoW is an approach to featurizing text data.

\uparrow
make into quantitative features

Given text: S. ... S.

Given texts s_1, s_2, \dots

e.g. $s_1 = \text{"My dog likes your dog"}$

We treat each as a bag of its words (ignore order but (potentially) remember counts)

$\text{BoW}(s_1) = \{ \text{"your"}, \text{"likes"}, \text{"My"}, \text{"dog"} \times 2 \}$

Then we make a dictionary of all words

$\mathcal{D} = \bigcup_{i=1}^n \text{BoW}(s_i)$ ← set of all words that appear in any text in dataset

and feature each text by

$X_{ij} = \# \text{ times that word } j \text{ appears in text } i$ $V_j \in \mathcal{D}$

(alt ver: 0/1 whether appears)

Intent: X_i capture the overall nature of the text

Important considerations: (you will explore HW2)

- de capitalization "Dog" → "dog"
- lemmatization "likes" → "like"
- prune contentless words, e.g. "a", "the", etc
- normalize or no normalization
- Bag-of-grams $n \geq 2$

Naïve Bayes w/ BoW

BoW generates very high-dim features

→ perfect for Naïve assumption,
which will help deal w/ long texts

Let's focus on $X_{ij} = 0$ or 1 whether
 S_i contains word j

$$\begin{aligned} \text{Then } \hat{P}(X_j=1 | Y=k) &= \hat{P}_{jk} \\ &= \frac{\sum_{i=1}^n \mathbb{I}[Y_i=k, X_{ij}=1] + \alpha}{\sum_{i=1}^n \mathbb{I}[Y_i=k] + \alpha} \\ &= \alpha\text{-smoothed frac of} \\ &\quad k\text{-texts that} \\ &\quad \text{have word } j \end{aligned}$$

Consider binary case (e.g. $Y \in \{\text{spam, not spam}\}$)

$$\begin{aligned} \text{logit}(\hat{P}(Y=1 | X=x)) &= \text{logit}(\hat{\pi}_1) \\ &+ \sum_{j=1}^p \left(x_j \log\left(\frac{\hat{p}_{j1}}{\hat{p}_{j0}}\right) + (1-x_j) \log\left(\frac{1-\hat{p}_{j1}}{1-\hat{p}_{j0}}\right) \right) \end{aligned}$$

Density estimation

Consider data $Y_1, \dots, Y_n \in \mathbb{R}$

drawn from a distribution

w/ CDF (cumulative dist fn) F ($F(y) = P(Y \leq y)$)

& PDF (prob density fn) f ($f(y) = F'(y)$)

We want to understand this unknown dist
from the data.

Estimating the CDF F is easy

empirical CDF: $\hat{F}_n(y) = \frac{1}{n} \sum_i \mathbb{I}[y_i \leq y]$

What about f ?

In particular if $y \notin \{y_1, \dots, y_n\}$
 (haven't seen y) then it's particularly
 difficult to say how likely y is

Histogram density estimation

Choose cutoffs $y_1 < y_2 < \dots < y_{m+1}$

Get bins $[y_1, y_2), [y_2, y_3), \dots, [y_m, y_{m+1})$

(we need to make sure that $y_1 \leq y_i \forall i$
 $y_{m+1} > y_i \forall i$)

Count the data in each bin

$$n_j = \sum_{i=1}^n \mathbb{I}[y_i \in [y_{j-1}, y_j)] = \# \text{ data pts in } j\text{th bin}$$

A histogram is just a bar chart
 w/ these counts

And a histogram density estimate

$$\hat{f}_n^{\text{hist}}(y) = \frac{n_j}{n} \frac{1}{y_j - y_{j-1}} \quad \text{where } j \text{ is s.t. } y \in [y_{j-1}, y_j)$$

But histograms are not smooth
 — doesn't really look like f

Kernel Density Estimate (KDE)

$$\hat{f}_n^{\text{KDE}}(y) = \frac{1}{n} \sum_{i=1}^n K_{\lambda}(y - Y_i)$$

$$\text{Where } K_{\lambda}(u) = \frac{1}{\lambda} \phi(u/\lambda)$$

$$\begin{array}{c} \uparrow \\ \text{std normal pdf} \\ \phi(u) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \|u\|^2} \end{array}$$