

1/31/2023: Naive Bayes

Classification Algorithms

Discriminative

- Logistic Regression
- Softmax Regression

Directly model $p(y|x)$

e.g. logistic regression

$$p(y=1|x;w) = \sigma(w^T x)$$

Don't try to model $p(x)$

Generative

- Naive Bayes

Jointly model $p(x,y)$

$$p(x,y) = \boxed{P(y)} \boxed{P(x|y)}$$

↑
prior distribution
over labels

↑
Given a label,
what does a
plausible x
look like?

$$P(y|x) = \frac{P(y) P(x|y)}{\boxed{P(x)}}$$

↑
Normalizing constant

$$= \sum_{k=1}^C p(y=k) P(x|y=k)$$

Naive Bayes: assume $x \in \mathbb{R}^d$

The Naive Bayes Assumption is: $P(x|y) = \prod_{j=1}^d P(x_j|y)$

i.e. all the x_j 's are conditionally independent given y

(we don't assume they are "independent")

$$p(y=0) = p(y=1) = 0.5$$

Suppose $x \in \mathbb{R}^2$
 $x_1, x_2 \in \{0, 1\}$

$$P(x_1=1 | y=0) = 0.9$$

$$P(x_1=1 | y=1) = 0.2$$

$$P(x_2=1 | y=0) = 0.8$$

$$P(x_2=1 | y=1) = 0.05$$

$$P(x | y=0)$$

$$P(x | y=1)$$

Note: x_1 and x_2 not independent (non-conditionally)


If $x_1=1$, $y=0$ is likely $\Rightarrow x_2=1$ is more likely

Common case: $x_j \in \{0, 1\} \forall j$

For this, we use Multivariate Bernoulli NB

many of these \rightarrow distribution over $\{0, 1\}^d$

Example 1: Black & white images

28×28  $\rightarrow 28^2 = 784$ -dim vector, each entry is $\{0, 1\}$

Example 2: Text classification

input = document \rightarrow $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$ aardvark

does the word occur in document?

\vdots

$\begin{bmatrix} 0 \end{bmatrix}$ zebra

Vocabulary V of size $|V|$

"bag of words"

Parameters of Multivariate Bernoulli NB model

• $P(y)$: Distribution over C classes \Rightarrow

$\pi \in \mathbb{R}^C$ where $p(y=k) = \pi_k$
parameter #1

• $P(x_j | y=k) \forall j \in \{1, \dots, d\}$, \rightarrow Each one is a Bernoulli

$K \in \{1, \dots, C\}$

So we have parameter $\tau \in \mathbb{R}^{d \times C}$
where $p(x_j=1 | y=k) = \tau_{jk}$

parameter #2

How to choose π & τ ? Apply MLE

$LL(\pi, \tau) = \sum_{i=1}^n \log P(x^{(i)}, y^{(i)}; \pi, \tau)$

log-likelihood

General form for generative classifier

$= \sum_{i=1}^n \log P(y^{(i)}; \pi) + \log P(x^{(i)} | y^{(i)}; \tau)$

$\sum_{i=1}^n \log P(y^{(i)}; \pi)$

$= \sum_{i=1}^n \sum_{k=1}^C \mathbb{1}[y^{(i)}=k] \log P(y=k; \pi)$
 $= \pi_k$

$= \sum_{i=1}^n \sum_{k=1}^C \mathbb{1}[y^{(i)}=k] \log \pi_k$

Let $\text{count}(y=k)$ means $\sum_{i=1}^n \mathbb{1}[y^{(i)}=k]$

$= \sum_{k=1}^C \text{count}(y=k) \log \pi_k$

If $C=2$:

$= \text{count}(y=1) \log \pi_1 + (n - \text{count}(y=1)) \cdot \log(1 - \pi_1)$

From HW0: maximized when $\pi_1 = \frac{\text{count}(y=1)}{n}$

Even when $C > 2$, MLE estimate for π is

$$\pi_k = \frac{\text{Count}(y=k)}{n}$$

what about τ ?

maximize $\sum_{i=1}^n \log P(x^{(i)} | y^{(i)}; \tau)$

By similar derivation, to maximize this, set

$$\tau_{jk} = \frac{\text{Count}(x_j=1, y=k)}{\text{Count}(y=k)}$$

\uparrow
 $= P(x_j=1 | y=k; \tau)$ Don't actually use!!

y	x ₁	x ₂
1	0	1
2	1	1
3	1	0
2	0	0
1	0	0
1	0	1
1	1	0
2	1	0

Estimate $\tau_{11} = P(x_1=1 | y=1) = \frac{1}{3}$

$\tau_{21} = P(x_2=1 | y=1) = \frac{2}{3}$

what happens when some counts are zero?

Text classification

- "giraffe" never occurred when $y=1$
- "choir" never occurred when $y=2$

What if a document has "giraffe" and "choir"?

$P(x | y=1) = 0$ b/c giraffe is "impossible"

$P(x | y=2) = 0$ b/c choir is "impossible"

Assuming zero probability for possible events is bad

Solution: Laplace Smoothing ("pseudocounts")

↓
Pretend we've seen every (feature, label) pair

Hyperparameter λ times outside of the dataset
 $\lambda = 1$ reasonable

Better formula for T_{jk} :

$$T_{jk} = \frac{\text{Count}(y=k, x_j=1) + \lambda}{\text{Count}(y=k) + 2 \cdot \lambda}$$

↑
once for $(y=k, x_j=1)$
once for $(y=k, x_j=0)$

If no training data then

$$T_{jk} = \frac{1}{2}$$

With enough training data, ignore λ 's

Another variant: Multinomial NB
(for text classification)

Input $x^{(i)}$ is document = list of words w/ length d_i

By Naive Bayes Assumption:

$$P(x^{(i)} | y^{(i)}) = \prod_{j=1}^{d_i} P(x_j^{(i)} | y^{(i)})$$

Note: preserves frequency info

multinomial distribution over V (vocabulary)

Additional assumption:

$P(x_j | y)$ is same for all j

Distribution of 1st word $| y$ = Distribution of 27th word $| y$

$$\text{doc} = \left[\text{dog dog dog} \dots \right] \\ P(\text{dog} | y)^3$$

Parameters

- $P(y)$ - same as before: π_k where $P(y=k) = \pi_k$
- $P(x_j | y=k)$ = Distribution over vocabulary V for each k

$$P_{uk} = P(x_j = u | y = k)$$

$u \in V$ $k \in \{1, \dots, C\}$

To estimate best P_{uk} , we count

y	x
1	I saw a...
0	Help.
1	... 100 words

$$P_{uk} = \frac{\text{count}(x_j = u, y = k) + \lambda}{\sum_{i=1}^n \mathbb{1}[y^{(i)} = k] d_i + |V| \cdot \lambda}$$

words in i -th doc

Add λ for every possible word in dictionary

whole denominator = # words that go with $y=k$