

9/12/2023: Naive Bayes

Classification Algorithm

Discriminative

- Logistic Regression
- Softmax Regression

Directly model $p(y|x)$

e.g. in logistic Regression

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

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

Generative

- Naive Bayes

Model joint distribution $p(x, y)$

$$p(x, y) = p(y) p(x|y)$$

prior distribution over labels

For a given label, what does a plausible x look like?

What to do at test time?

Given x , want to predict y

Bayes Rule:

$$p(y|x) = \frac{p(y) p(x|y)}{p(x)}$$

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

Naive Bayes for text classification

Training Data

input x is a document

y	x
+1	great acting and score
-1	terrible directing
+1	great execution
-1	terrible
+1	amazing

Test time:

New review e.g. "great directing"

Naive Bayes: Make a key simplifying assumption so that estimating $p(x|y)$ is not too difficult

$$|V| = 8$$

Parameters: $- p(y) : +1 = 3/5 = .6$
 $-1 = 2/5 = .4$

$- p(x_i|+1)$: Distribution over 8 words

$- p(x_i|-1)$: Different distribution over 8 words

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

(where x_j is j -th word of x)

① Parameters of model

- $P(y)$: If we have C classes, need parameter vector $\pi \in \mathbb{R}^C$ where $P(y=k) = \pi_k$
- $P(x|y)$: Because of NB assumption, we just have to model $p(x_j|y)$

Imagine C different dice
Each dice has $|V|$ sides

\uparrow vocabulary is set of possible words

\hookrightarrow use $\tau_{wk} = p(x_j=w|y=k)$ for some $w \in V$

② Learning parameters: Apply MLE to joint distribution

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

$$= \underbrace{\left(\sum_{i=1}^n \log P(y^{(i)}; \pi, \tau) \right)}_{\text{only depends on } \pi} + \underbrace{\sum_{i=1}^n \log P(x^{(i)}|y^{(i)}; \tau)}_{\text{only depends on } \tau}$$

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

maximized when

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

For each label k , we have $|V|$ -sided dice, estimate

$$P(x_j=w|y=k) = \frac{\text{count}(w, y=k)}{\sum_{w' \in V} \text{count}(w', y=k)}$$

is τ_{wk}

DO NOT USE THIS FORMULA \rightarrow

Training Data	
y	x
+1	great acting and score
-1	terrible directing
+1	great execution
-1	terrible
+1	amazing

$P(x_j = \text{"great"} | +1) = 2/7$
 acting = $1/7$
 score = $1/7$
 No Smoothing
 terrible = $0/7$ because it generates 0 probabilities

With Smoothing: $P(x_j = \text{"great"} | +1) = \frac{2+1}{7+8} = \frac{3}{15} = \frac{1}{5}$
 of $\lambda = 1$

Suppose: Test example $x = \text{"great directing"} \text{"terrible"} = \frac{1}{15}$

- compute $P(x, y=1) = P(y=1)P(x|y=1) = 3/5 \times 2/2 \times 0/7 = 0$
- compute $P(x, y=-1) = 2/5 \times 0/3 \times 1/3 = 0$
- by Bayes Rule, get $P(y|x) = \frac{0}{0}$

Fix: Laplace Smoothing

Imagine that every (word, label) pair was seen an additional λ times

"pseudocounts"

λ hyperparameter, > 0

Actual Formula:
$$T_{wk} = \frac{\text{Count}(w, y=k) + \lambda}{\left(\sum_{w' \in V} \text{Count}(w', y=k) \right) + |V| \cdot \lambda}$$

total # of imaginary counts that were added

Avoiding Underflow

$x = \text{"great directing but the music was ..."}'$

$P(y)P(x|y) = P(y)P(\text{"great"}|y) \cdot P(\text{"directing"}|y) \cdot \dots$

Multiplying many small #'s on a computer leads to "underflow" i.e. answer is 0

Trick: Work in log space

ie: Don't compute $P(x|y)$

$$\text{Instead compute } \log P(x|y) = \sum_{j=1}^d \log P(x_j|y)$$

How to predict?

- For each label k , compute $\log P(y=k) + \log P(x|y=k)$
- k where this is largest also has largest $P(x, y) \Leftrightarrow$ largest $P(y=k|x)$

Naive Bayes for general feature vectors

- So far, x was a document with d words
- Now: x is a list/vector of d features

Black/white

Image:

28



28

You have $28^2 = 784$

{0,1} features

Song

Feature 1

genre \rightarrow pop, classical, country...

Feature 2

artist

⋮

Key diff: Each feature means something different

Same: Naive Bayes Assumption $P(x|y) = \prod_{i=1}^d P(x_i|y)$

Different: $P(x, y)$ different from $P(x_2|y)$

e.g. Songs:

Feature 1

$$P(\text{genre} = \text{pop} | y = 1)$$

$$P(\text{genre} = \text{rock} | y = 1)$$

$$P(\text{genre} = \text{pop} | y = 0)$$

$$P(\text{genre} = \text{rock} | y = 0)$$

⋮

Each one is parameter specific to feature 1