

Lecture 4: Naive Bayes (the Frequentist approach and the Bayesian approach)

Julia Hockenmaier

juliahmr@illinois.edu
3324 Siebel Center

<http://www.cs.uiuc.edu/class/sp10/cs598jhm>

The model

Today's class

The task: text classification (sentiment analysis)

Assign (sentiment) label $L_i \in \{+, -\}$ to a document $\mathbf{W}_i = (w_{i1} \dots w_{iN})$.

$\mathbf{W}_1 =$ "This is an amazing product: great battery life, amazing features and it's cheap."

$\mathbf{W}_2 =$ "How awful. It's buggy, saps power and is way too expensive."

The data:

A set \mathcal{D} of N documents with (or without) labels

The model:

Naive Bayes

Comparing different estimation techniques:

- Supervised MLE
- Unsupervised MLE with EM
- Unsupervised Bayesian Estimation with Gibbs sampling
- Supervised Bayesian Estimation

CS598JHM: **Advanced NLP**

2

A Naive Bayes model

The task:

Assign (sentiment) label $L_i \in \{+, -\}$ to a document \mathbf{W}_i .

$\mathbf{W}_1 =$ "This is an amazing product: great battery life, amazing features and it's cheap."

$\mathbf{W}_2 =$ "How awful. It's buggy, saps power and is way too expensive."

The model:

- Use Bayes' Rule:

$$L_i = \operatorname{argmax}_L P(L | \mathbf{W}_i) = \operatorname{argmax}_L P(\mathbf{W}_i | L)P(L)$$

- Assume \mathbf{W}_i is a "bag of words":

$\mathbf{W}_1 = \{an: 1, and: 1, amazing: 2, battery: 1, cheap: 1, features: 1, great: 1, \dots\}$

$\mathbf{W}_2 = \{awful: 1, and: 1, buggy: 1, expensive: 1, \dots\}$

- $P(\mathbf{W}_i | L)$ is a multinomial distribution: $\mathbf{W}_i \sim \text{Multinomial}(\theta_L)$

We have a vocabulary of V words. Thus: $\theta_L = (\theta_1, \dots, \theta_V)$

- $P(L)$ is a Bernoulli distribution: $L \sim \text{Bernoulli}(\pi)$

CS598JHM: **Advanced NLP**

4

Using this model

The model:

$P(\mathbf{W}_i | L)$ is a multinomial distribution: $\mathbf{W}_i \sim \text{Multinomial}(\theta_L)$

$P(L)$ is a Bernoulli distribution: $L \sim \text{Bernoulli}(\pi)$

The “frequentist” approach (MLE):

Estimate π, θ_+, θ_- , then:

$$P(L_i = + | \mathbf{W}_i) \propto P(\mathbf{W}_i | \theta_+) \pi$$

The Bayesian approach:

Choose priors for $\pi \sim \text{Beta}(\alpha, \beta)$,

$\theta_+ \sim \text{Dirichlet}(\gamma), \theta_- \sim \text{Dirichlet}(\gamma)$ then compute the following expectation:

$$P(L_i = + | \mathbf{W}_i) \propto \iint P(\mathbf{W}_i | \theta_+) \theta_+ P(\theta_+; \gamma) P(\pi; \alpha, \beta) d\theta_+ d\pi$$

The frequentist approach

Supervised MLE

The data is labeled:

We have a set \mathbf{D} of D documents $\mathbf{W}_1 \dots \mathbf{W}_d$ with N words

Each document \mathbf{W}_i has N^i words

D^+ documents (subset \mathbf{D}^+) have a positive label and N^+ words

D^- documents (subset \mathbf{D}^-) have a negative label and N^- words

Each word $w_1 \dots w_i \dots w_V$ appears $N^+(w_i)$ times in \mathbf{D}^+ , $N^-(w_i)$ times in \mathbf{D}^-

Each word $w_1 \dots w_i \dots w_V$ appears $N(w_i)$ times in D^i

MLE: relative frequency estimation

- Labels: $L \sim \text{Bernoulli}(\pi)$ with $\pi = D^+ / d$

- Words: $\mathbf{W}_i | + \sim \text{Multinomial}(\theta^+)$ with $\theta_i^+ = N^+(w_i) / N^+$

- Words: $\mathbf{W}_i | - \sim \text{Multinomial}(\theta^-)$ with $\theta_i^- = N^-(w_i) / N^-$

Inference with MLE

The inference task:

Given a new document \mathbf{W}_{i+1} , what is its label L_{i+1} ?

Word w_j occurs $N_{i+1}(w_j)$ times in \mathbf{W}_{i+1} .

$$\begin{aligned} P(L = + | \mathbf{W}_{i+1}) &\propto P(+) P(\mathbf{W}_{i+1} | +) \\ &= \pi \prod_{j=1}^V \theta_{+j}^{N_{i+1}(w_j)} \end{aligned}$$

Unsupervised MLE

The data is *unlabeled*:

We have a set \mathbf{D} of D documents $W_1 \dots W_d$ with N words

Each document W_i has N^i words

Each word $w_1 \dots w_i \dots w_V$ appears $N^i(w_i)$ times in W_j

EM algorithm: “expected rel. freq. estimation”

Initialization: pick initial $\pi^{(0)}$, $\theta^{+(0)}$, $\theta^{-(0)}$

Iterate:

- Labels: $L \sim \text{Bernoulli}(\pi)$ with $\pi^{(l)} = \langle N_{+} \rangle_{(t-1)} / \langle N \rangle_{(t-1)}$
- Words: $W_i | + \sim \text{Multinomial}(\theta^+)$ with $\theta_i^{+(l)} = \langle N^+(w_i) \rangle_{(t-1)} / \langle W^+ \rangle_{(t-1)}$
- Words: $W_i | - \sim \text{Multinomial}(\theta^-)$ with $\theta_i^{-(l)} = \langle N^-(w_i) \rangle_{(t-1)} / \langle W^- \rangle_{(t-1)}$

CS598JHM: Advanced NLP

9

The Bayesian approach

The Bayesian approach

We need to compute an integral

$$P(L_i = + | W_i) \propto \int P(W_i | \theta_+) \theta_+ P(\theta_+; \gamma) P(\pi; \alpha, \beta) d\theta_+ d\pi$$

Case 1: we have labeled data

Case 2: we do not have labeled data

CS598JHM: Advanced NLP

11

Bayesian: supervised

The data is labeled:

We have a set \mathbf{D} of D documents $W_1 \dots W_d$ with N words

Each document W_i has N^i words

D^+ documents (subset \mathbf{D}^+) have a positive label and N^+ words

D^- documents (subset \mathbf{D}^-) have a negative label and N^- words

Each word $w_1 \dots w_i \dots w_V$ appears $N^+(w_i)$ times in \mathbf{D}^+ , $N^-(w_i)$ times in \mathbf{D}^-

Bayesian estimation

- $P(+) = (D^+ + \alpha) / (D + \alpha + \beta)$
- $P(w_i | +) = (N^+(w_i) + \gamma_i) / (N^+(w_i) + \gamma_0)$
- $P(W_i | +) = \prod P(w_j | +)^{N^i(w_j)}$

CS598JHM: Advanced NLP

12

Bayesian: unsupervised

We need to approximate the integral/expectation:
 $P(L_i = + | \mathbf{W}_i) \propto \int P(\mathbf{W}_i | \theta_+) P(\theta_+; \gamma) P(\pi; \alpha, \beta) d\theta_+ d\pi$

We can approximate the expectation of $f(x)$ by sampling a finite number of points x_1, \dots, x_N according to $p(x)$, evaluating $f(x_i)$ for each of them, and computing the average.

How can we sample according to $p(x)$?

For us $p(x) = p(D, \mathbf{L}, \pi, \theta_+, \theta_-; \alpha, \beta, \gamma)$

Gibbs sampling

We will visit states according to transition probabilities $P(y|x)$

That is, we will go from state $\mathbf{x} = (x_1, \dots, x_k)$
 to state $\mathbf{y} = (y_1, \dots, y_k)$

For $i = 1 \dots k$:
 pick y_i by sampling from $P(Y_i | y_1, \dots, y_{i-1}, x_{i+1}, \dots, x_k)$

$$P(Y_i = y_i | y_1, \dots, y_{i-1}, x_{i+1}, \dots, x_k) = \frac{P(y_1, \dots, y_{i-1}, y_i, x_{i+1}, \dots, x_k)}{P(y_1, \dots, y_{i-1}, x_{i+1}, \dots, x_k)}$$

Markov Chain Monte Carlo

If we had discrete distribution $p(\mathbf{x}) = p(x_1, \dots, x_k)$, $p(\mathbf{x})$ has only a finite number of outcomes.

Markov Chain Monte Carlo methods construct a Markov chain whose states are the outcomes of $p(\mathbf{x})$.

The probability of visiting state x_j is $p(x_j)$

We sample from $p(\mathbf{x})$ by visiting a sequence of states from this Markov chain.

Gibbs sampling

For us $p(x) = p(D, \mathbf{L}, \pi, \theta_+, \theta_-; \alpha, \beta, \gamma)$

π, θ_+, θ_- are real-valued, but they will disappear:

$$P(L_j = + | \mathbf{L}^{(-j)}; \alpha, \beta) = \frac{\alpha + N_+^{(-j)}}{\alpha + \beta + N - 1}$$

$$P(w_k = y | D_+^{(-j)}; \gamma) = \frac{N_{D_x^{(-j)}}(y) + \gamma_y}{\gamma_0 + N_{D_x^{(-j)}}$$

The Gibbs sampler

Initialize:

Define priors α, β, γ .

Assign initial labels $\mathbf{L}^{(0)}$ to documents

Iterate:

For each iteration $t = 1 \dots T$:

For every document \mathbf{W}_i (with current label $x = L_i^{(t-1)}$)

(Temporarily) remove its word counts $N_i(w_j)$ from its class x :

$$N_{x_i}^{(t-1)}(w_j) = N_x^{(t-1)}(w_j) - N_i^{(t-1)}(w_j)$$

(Temporarily) remove \mathbf{W}_i from the documents in its class x :

$$D_{x_i}^{(t-1)} = D_x^{(t-1)} - 1$$

Assign a new label $x' = L_i^{(t)}$ to \mathbf{W}_i with

$$P(L_i | \mathbf{W}_i, L_0^{(0)} \dots L_{i-1}^{(0)}, L_{i+1}^{(0)} \dots L_D^{(0)}, \alpha, \beta, \gamma)$$

Add \mathbf{W}_i to the documents in class x'

Add its word counts $N_i(w_j)$ to word counts for class x'

Final estimate:

Use (some of the) snapshots $\mathbf{L}^{(1)} \dots \mathbf{L}^{(T)}$ to estimate $P(+)$, $P(w_i | +)$, $P(w_i | -)$

Why we don't need to estimate π

$$\begin{aligned}
 P(L_j = + | \mathbf{L}^{(-j)}; \alpha, \beta) &= \int P(L_j = + | \pi) P(\pi | \mathbf{L}^{(-j)}; \alpha, \beta) d\pi \\
 &= \int \pi P(\pi | \mathbf{L}^{(-j)}; \alpha, \beta) d\pi \\
 &= \int \pi \frac{\Gamma(\alpha + \beta + N - 1)}{\Gamma(\alpha + N_+^{(-j)}) \Gamma(\beta + N_-^{(-j)})} \pi^{\alpha + N_+^{(-j)} - 1} (1 - \pi)^{\beta + N_-^{(-j)} - 1} d\pi \\
 &= \frac{\Gamma(\alpha + \beta + N - 1)}{\Gamma(\alpha + N_+^{(-j)}) \Gamma(\beta + N_-^{(-j)})} \int \pi^{\alpha + N_+^{(-j)} - 1} (1 - \pi)^{\beta + N_-^{(-j)} - 1} d\pi \\
 &= \frac{\Gamma(\alpha + \beta + N - 1)}{\Gamma(\alpha + N_+^{(-j)}) \Gamma(\beta + N_-^{(-j)})} \frac{\Gamma(\alpha + N_+^{(-j)} + 1) \Gamma(\beta + N_-^{(-j)})}{\Gamma(\alpha + \beta + N)} \\
 &= \frac{\Gamma(\alpha + \beta + N - 1)}{\Gamma(\alpha + N_+^{(-j)}) \Gamma(\beta + N_-^{(-j)})} \frac{\Gamma(\alpha + N_+^{(-j)} + 1)}{\Gamma(\alpha + \beta + N)} \\
 &= \frac{\Gamma(\alpha + \beta + N - 1)}{\Gamma(\alpha + N_+^{(-j)})} \frac{\Gamma(\alpha + N_+^{(-j)} + 1)}{(\alpha + \beta + N - 1) \Gamma(\alpha + \beta + N - 1)} \\
 &= \frac{\Gamma(\alpha + N_+^{(-j)} + 1)}{\Gamma(\alpha + N_+^{(-j)}) (\alpha + \beta + N - 1)} \\
 &= \frac{(\alpha + N_+^{(-j)}) \Gamma(\alpha + N_+^{(-j)})}{\Gamma(\alpha + N_+^{(-j)}) (\alpha + \beta + N - 1)} \\
 &= \frac{\alpha + N_+^{(-j)}}{\alpha + \beta + N - 1}
 \end{aligned}$$