

Naïve Bayes classifier

Easy

For categorical
variables

Positive or negative movie review?

- ⇒ 0 + ⇒ 1

- ▶ This movie is disappointing. —
- ▶ This movie does not disappoint. +
- ▶ I would love to have that two hours of my life back. —
- ▶ This is one of my favorite if not favorite films. +
- ▶ I have seen so many bad low budget films lately, but I love this one. +

Deterministic model

$$f(\text{I love, love this movie.}) = 0 \text{ or } 1$$

Probabilistic model

$$\begin{aligned} f(\text{I love, love this movie.}) \\ = \mathbb{P}(y = 1 | \text{I love, love this movie.}) \end{aligned}$$

Naïve Bayes

- ▶ Simple (naïve) classifier based on **Bayes rule**:

For a document $d = \{w_1, w_2, \dots, w_n\}$ and a class $y = 0$ or 1

$$\mathbb{P}(y = 1 | w_1, w_2, \dots, w_n) = \frac{\mathbb{P}(w_1, w_2, \dots, w_n | y = 1) \mathbb{P}(y = 1)}{\mathbb{P}(w_1, w_2, \dots, w_n)}$$

$$\mathbb{P}(y = 0 | w_1, w_2, \dots, w_n) = \frac{\mathbb{P}(w_1, w_2, \dots, w_n | y = 0) \mathbb{P}(y = 0)}{\mathbb{P}(w_1, w_2, \dots, w_n)}$$

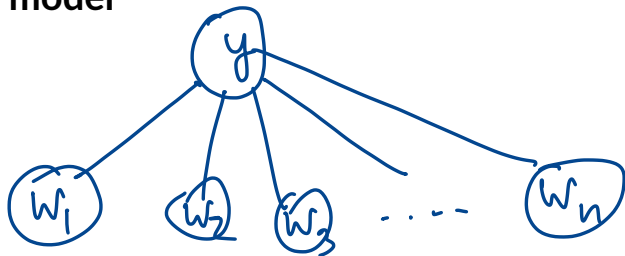
- ▶ Compute $\mathbb{P}(w_1, w_2, \dots, w_n | y)$ and $\mathbb{P}(y)$ from the (labeled) data.
- ▶ What about $\mathbb{P}(w_1, w_2, \dots, w_n)$?

Conditional independence

We assume that $\mathbb{P}(w_i|y)$ are independent given the class y .

$$\mathbb{P}(w_1, w_2, \dots, w_n|y) = \mathbb{P}(w_1|y)\mathbb{P}(w_2|y) \dots \mathbb{P}(w_n|y).$$

Graphical model



Maximum a posteriori (MAP)

$$\mathbb{P}(w_1, w_2, \dots, w_n | y) \mathbb{P}(y) = \mathbb{P}(y) \mathbb{P}(w_1 | y) \mathbb{P}(w_2 | y) \dots \mathbb{P}(w_n | y)$$

where

$$\mathbb{P}(w_i | y) = \frac{\text{count}(w_i, y)}{\sum_{w \in V} \text{count}(w, y)}$$

and

$$\mathbb{P}(y) = \frac{\text{countdoc}(Y = y)}{\text{count}(\text{Documents})}$$

Example

- ▶ This movie is disappointing. —
- ▶ This movie does not disappoint. +
- ▶ I would love to have that two hours of my life back. —
- ▶ This is one of my favorite if not favorite films. +
- ▶ I have seen so many bad low budget films lately, but I love this one. +

$$\mathbb{P}(\text{Positive}) = \frac{3}{5}$$

$$p(\text{negative}) = \frac{2}{5}$$

Example

- ▶ This movie is disappointing. —
- ▶ This movie does not disappoint. f
- ▶ I would love to have that two hours of my life back. —
- ▶ This is one of my favorite if not favorite films. f
- ▶ I have seen so many bad low budget films lately, but I love this one.

$$\mathbb{P}(\text{favorite}|\text{Positive}) = \frac{2}{30}$$

+

Example

- ▶ This movie is disappointing. —
- ▶ This movie does not disappoint. †
- ▶ I would love to have that two hours of my life back. —
- ▶ This is one of my favorite if not favorite films. †
- ▶ I have seen so many bad low budget films lately, but I love this one. †

$$\mathbb{P}(\text{disappoint} | \text{Positive}) = \frac{1}{30}$$

Example

$\mathbb{P}(y = 1 | \text{I love, love this movie.})$

$$= P(I | y=1) P(\text{love} | y=1) P(\text{this} | y=1) P(\text{movie} | y=1)$$

$$= \frac{2}{30} \times \frac{1}{30} \times \frac{3}{30} \times \frac{1}{30} \times \frac{3}{5} P(y=1)$$

$$= \frac{18}{9050000}$$

Example

$\Rightarrow P(y=0 | \dots) > P(y=1 | \dots)$
 \Rightarrow classify $y=0$

$\mathbb{P}(y = 0 | \text{I love, love this movie.})$

$$= P(I | y=0) P(\text{love} | y=0) P(\text{this} | y=0) P(\text{movie} | y=0) \\ \times P(y=0)$$

$$= \frac{1}{16} \times \frac{1}{16} \times \frac{1}{16} \times \frac{1}{16} \times \frac{2}{5}$$

$$= \frac{2}{16^4 \times 5} > \frac{18}{49,500,000} = P(y=1 | \text{I love but this movie})$$

- ▶ What if there is a word **slept** in the test set, but not in the training set?

$$\mathbb{P}(\text{slept}|y) = \frac{\text{count}(\text{slept}, y)}{\sum_{w \in V} \text{count}(w, y)} = 0.$$

- ▶ There is no best y in this case.

$$\mathbb{P}(y|\text{slept}, \dots) = \mathbb{P}(y)\mathbb{P}(\text{slept}|y) \dots = 0$$

Laplace smoothing

Fix $\alpha > 0$.

$$\begin{aligned}\mathbb{P}(w_i|y) &= \frac{\text{count}(w_i, y) + \alpha}{\sum_{w \in V} (\text{count}(w, y) + \alpha)} \\ &= \frac{\text{count}(w_i, y) + \alpha}{\sum_{w \in V} \text{count}(w, y) + \alpha |V|}\end{aligned}$$

of words in training set

In the previous example, if we choose $\alpha = 1$,

$$\begin{aligned}\mathbb{P}(\text{slept}|y) &= \frac{1}{\sum_{w \in V} \text{count}(w, y) + 1 |V|} \\ &= \frac{1}{\text{\# of words in class } y + |V|}\end{aligned}$$

Learning Naïve Bayes

- ▶ From the training corpus, extract the **Vocabulary**.
- ▶ For each class y , calculate $\mathbb{P}(y)$
 - ▶ Count number of documents in class y .
 - ▶ $\mathbb{P}(y) = \frac{\text{countdoc}(Y=y)}{\text{count}(\text{Documents})}$
- ▶ For each word w_i and class y
 - ▶ Merge all documents in class y
 - ▶ $n_i \leftarrow \#$ of occurrence of each word in class y
 - ▶ $\mathbb{P}(w_i|y) = \frac{n_i + \alpha}{\sum_j n_j + \alpha |\text{Vocab}|}$