

• Recall the principle of Maximum Likelihood Estimation (MLE):

Probabilistic text classification using BOW features



$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} P(\mathbf{X}, \mathbf{y}; \boldsymbol{\theta})$$

$$= \operatorname{argmax}_{\boldsymbol{\theta}} \left(\prod_{i=1}^M P(\mathbf{x}_i, y_i ; \boldsymbol{\theta}) \right)^{\frac{1}{M}}$$

$$= \operatorname{argmax}_{\boldsymbol{\theta}} \frac{1}{M} \sum_{i=1}^M \log P(\mathbf{x}_i, y_i ; \boldsymbol{\theta})$$

$$= \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^M \log P(\mathbf{x}_i, y_i ; \boldsymbol{\theta})$$

$\theta =$ parameters to estimate

X **E** **R** **M** **X** **N**

$$y \in \{1, \dots, K\}^M$$

Number of observations

$K = \text{number of categories/classes}$

where

$N \equiv$ vocabulary size

Note: P is a generative model (not the joint distribution x, y , not $y|x$)

Note: y is a representation of label, not a non-vector, hence lower case

← BOW features

Probabilistic text classification using BOW features

- Recall the principle of Maximum Likelihood Estimation (MLE):

$$\hat{\theta} = \operatorname{argmax}_{\theta} P(\mathbf{X}, \mathbf{y}; \theta)$$

$$= \operatorname{argmax}_{\theta} \left(\prod_{i=1}^M P(\mathbf{x}_i, y_i; \theta) \right)^{\frac{1}{M}}$$

$$= \operatorname{argmax}_{\theta} \frac{1}{M} \sum_{i=1}^M \log P(\mathbf{x}_i, y_i; \theta)$$

$$= \operatorname{argmax}_{\theta} \sum_{i=1}^M \log P(\mathbf{x}_i, y_i; \theta)$$

where

θ = parameters to estimate

$\mathbf{X} \in \mathbb{R}^{M \times N}$ ← BOW features

$\mathbf{y} \in \{1, \dots, K\}^M$

M = number of observations

N = vocabulary size

K = number of categories/classes

Note: P is a generative model (note the joint distribution x, y , not $y|x$)

Note: y is a scalar representation of label, not a one-hot vector, hence lower case

Naive Bayes' classifier

- Assumption 1: The text and label in one document does not affect those of another.
- Assumption 2: Words in a sentence are independent, conditioned on the class label

- Taken from Eisenstein, 2019, Chp 2

Algorithm 1 Generative process for the Naïve Bayes classification model

for Instance $i \in \{1, 2, \dots, M\}$ **do**:

 Draw the label $y^{(i)} \sim \text{Categorical}(\boldsymbol{\mu})$;

 Draw the word counts $\mathbf{x}^{(i)} \mid y^{(i)} \sim \text{Multinomial}(\boldsymbol{\phi}_{y^{(i)}})$.

chain rule $P(\mathbf{x}, y) = P(\mathbf{x} \mid y)P(y)$

where $\boldsymbol{\mu} = [\mu_1, \dots, \mu_K]$ label probability

$\boldsymbol{\phi} = [\phi_1, \dots, \phi_N]$ word probability

$$P_{mult}(\mathbf{x} \mid y; \boldsymbol{\phi}) = B(\mathbf{x}) \prod_{j=1}^N \phi_j^{\mathbf{x}_j} \quad B(\mathbf{x}) = \frac{(\sum_{j=1}^N \mathbf{x}_j)!}{\prod_{j=1}^N \mathbf{x}_j!}$$