

Naive Bayes II

Recall: Naive Bayes assumption:

features X_1, X_2, \dots, X_d independent given label y .

$$P(X|y) = \prod_{\alpha=1}^d P(X_\alpha|y)$$

Categorical Naive Bayes

assumption: $X_\alpha \in \{f_1, f_2, f_3, \dots, f_{K_\alpha}\} = X_\alpha$

eg. $X = X_1 \times X_2 \times X_3 \times \dots \times X_d$

$$|X_\alpha| = K_\alpha.$$

model: ~~$P(X_\alpha|Y)$~~ .

$$P(X_\alpha = j|Y) = \theta_{j,Y,\alpha}$$

where $\theta_{j,Y,\alpha} \geq 0$ and $\sum_{j \in X_\alpha} \theta_{j,Y,\alpha} = 1$.

categorical distribution.

estimating θ :

$$\hat{\theta}_{j,Y,\alpha} = \frac{\text{\# of examples with label } y \text{ and } \alpha\text{th feature } j}{\text{\# of examples with label } y}$$

$I(\text{true}) = 1$
 $I(\text{false}) = 0$.

$$\hat{\theta}_{j,c,\alpha} = \frac{\sum_{i=1}^n I(y_i=c) I(x_{i,\alpha}=j) + l}{\sum_{i=1}^n I(y_i=c) + l \cdot K_\alpha}$$

Prediction:

$$h(x) = \arg \max_{c \in \mathcal{Y}} \hat{P}(y=c|x)$$

$$= \arg \max_{c \in \mathcal{Y}} \hat{P}(x|y=c) \cdot \hat{P}(y=c)$$

$$= \arg \max_{c \in \mathcal{Y}} \frac{1}{\prod_{\alpha} c} \cdot \frac{m!}{x_1! x_2! \dots x_d!} \prod_{\alpha=1}^d (\theta_{\alpha c})^{x_{\alpha}}$$

$$= \arg \max_{c \in \mathcal{Y}} \frac{1}{\prod_{\alpha} c} \cdot \prod_{\alpha=1}^d (\theta_{\alpha c})^{x_{\alpha}}$$

$$= \arg \max_{c \in \mathcal{Y}} \underbrace{\log\left(\frac{1}{\prod_{\alpha} c}\right)}_{b_c} + \sum_{\alpha=1}^d x_{\alpha} \underbrace{\log(\theta_{\alpha c})}_{w_{\alpha c}}$$

$$= \arg \max_{c \in \mathcal{Y}} b_c + \sum_{\alpha=1}^d w_{\alpha c} \cdot x_{\alpha}$$

$$= \arg \max_{c \in \mathcal{Y}} (Wx + b)_c.$$

$W \in \mathbb{R}^{C \times d}$

$$\hat{\theta}_{j,c,\alpha} = \frac{|\{i \mid y_i = c \text{ and } x_{i,\alpha} = j\}| + l}{|\{i \mid y_i = c\}| + l \cdot K_\alpha}$$

$$= \frac{|\{(x,y) \in \mathcal{D} \mid y = c \text{ and } x_\alpha = j\}| + l}{|\{(x,y) \in \mathcal{D} \mid y = c\}| + l \cdot K_\alpha}$$

Prediction

$$h(x) = \arg \max_{c \in \mathcal{Y}} \hat{P}(y=c|x)$$

$$= \arg \max_{c \in \mathcal{Y}} \frac{\hat{P}(x|y=c) \cdot \hat{P}(y=c)}{\hat{P}(x)} \quad (\text{Bayes Thm})$$

$$= \arg \max_{c \in \mathcal{Y}} \hat{P}(x|y=c) \hat{P}(y=c)$$

$$= \arg \max_{c \in \mathcal{Y}} \hat{P}(y=c) \cdot \prod_{\alpha=1}^d \hat{P}(x_\alpha|y=c) \quad (\text{Naive Bayes assumption})$$

$$h(x) = \arg \max_{c \in \mathcal{Y}} \hat{\pi}_c \prod_{\alpha=1}^d \hat{\theta}_{x_\alpha, c, \alpha}$$

$$\hat{\pi}_c = \frac{|\{(x,y) \in \mathcal{D} \mid y=c\}|}{n} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(y_i=c)$$

Q: For ~~multi~~ categorical naive Bayes, how many parameters are there? ($\hat{\pi}$ and $\hat{\theta}$)
 in terms of $n, d, K_\alpha, |\mathcal{Y}|=C$ i.e. what is the dimension of the parameter vector?

"I am spam. Spam I am. Do you like perceptrons?"

I → 2

am → 2

spam → 2

do → 1

you → 1

like → 1

perception → 1.

$$x = (\dots, 2, \dots, 1, \dots, 1, \dots)$$

sparse!

$d = \#$ of english words

$m = 7$

Multinomial Naive Bayes

features $x_\alpha \in \mathbb{N}$, sum up to $\sum_{\alpha=1}^d x_\alpha = m$. length

Model: x are multinomially distributed

$$P(x|y) = \frac{m!}{x_1! x_2! \dots x_d!} \prod_{\alpha=1}^d (\theta_{\alpha y})^{x_\alpha}$$

e.g. a 6-sided die: $P(x) = \frac{m!}{x_1! x_2! x_3! x_4! x_5! x_6!} \cdot \prod_{\alpha=1}^6 \left(\frac{1}{6}\right)^{x_\alpha}$

Estimate parameters:

$$\hat{\theta}_{\alpha c} = \frac{\sum_{i=1}^n I(y_i=c) \cdot x_{i\alpha} + l}{\sum_{i=1}^n I(y_i=c) \cdot m_i + l \cdot d}$$

length of i th training example $m_i = \sum_{\alpha=1}^d x_{i\alpha}$

of params: $C + dC$

Gaussian Naive Bayes!

Features: ~~x_α~~ $x_\alpha \in \mathbb{R}$, $x \in \mathbb{R}^d$

Model: $x_\alpha \sim \mathcal{N}(\mu_{\alpha c}, \sigma_{\alpha c}^2)$

$$P(x_\alpha | y=c) = \frac{1}{\sqrt{2\pi\sigma_{\alpha c}^2}} \exp\left(-\frac{(x_\alpha - \mu_{\alpha c})^2}{2\sigma_{\alpha c}^2}\right)$$

Estimating Parameters:

$$\hat{\mu}_{\alpha c} = \frac{\sum_{i=1}^n I(y_i=c) x_{i\alpha}}{\sum_{i=1}^n I(y_i=c)}$$

$$\hat{\sigma}_{\alpha c}^2 = \frac{\sum_{i=1}^n I(y_i=c) (x_{i\alpha} - \hat{\mu}_{\alpha c})^2}{\sum_{i=1}^n I(y_i=c) - 1}$$

or! $\mathbb{P} x_\alpha \sim \mathcal{N}(\mu_{\alpha c}, \sigma_{\alpha c}^2)$

here: $\hat{\sigma}_{\alpha c}^2 = \frac{1}{n} \sum_{i=1}^n (x_{i\alpha} - \hat{\mu}_{\alpha c})^2$

Parameters are $\hat{\pi}_1, \hat{\pi}_2, \dots, \hat{\pi}_c, \dots, \hat{\mu}_{\alpha y}, \dots, \hat{\sigma}_{\alpha}^2$

$$\hat{P}(x_{\alpha}|y) = \frac{1}{\sqrt{2\pi} \cdot \hat{\sigma}_{\alpha}} \cdot \exp\left(-\frac{(x_{\alpha} - \hat{\mu}_{\alpha y})^2}{2\hat{\sigma}_{\alpha}^2}\right)$$

What is $\hat{P}(y|x)$? (As a function of π, μ, σ , and x).

$$\hat{P}(y|x) = \frac{\hat{P}(x|y) \cdot \hat{P}(y)}{\hat{P}(x)} \quad (\text{Bayes rule}) \quad \text{Naive Bayes (Assumption + Gaussian)}$$

$$= \frac{\hat{P}(y)}{\hat{P}(x)} \cdot \prod_{\alpha=1}^d \hat{P}(x_{\alpha}|y) = \frac{\hat{\pi}_y}{\hat{P}(x)} \cdot \prod_{\alpha=1}^d \frac{1}{\sqrt{2\pi} \cdot \hat{\sigma}_{\alpha}} \cdot \exp\left(-\frac{(x_{\alpha} - \hat{\mu}_{\alpha y})^2}{2\hat{\sigma}_{\alpha}^2}\right)$$

$$\text{let } \frac{1}{Z} = \frac{1}{\hat{P}(x)} \cdot \prod_{\alpha=1}^d \frac{1}{\sqrt{2\pi} \cdot \hat{\sigma}_{\alpha}} \Rightarrow = \frac{\hat{\pi}_y}{Z} \prod_{\alpha=1}^d \exp\left(-\frac{(x_{\alpha} - \hat{\mu}_{\alpha y})^2}{2\hat{\sigma}_{\alpha}^2}\right)$$

$$= \frac{\hat{\pi}_y}{Z} \cdot \exp\left(\sum_{\alpha=1}^d -\frac{1}{2\hat{\sigma}_{\alpha}^2} (x_{\alpha} - \hat{\mu}_{\alpha y})^2\right)$$

$$= \frac{\hat{\pi}_y}{Z} \cdot \exp\left(\sum_{\alpha=1}^d -\frac{x_{\alpha}^2}{2\hat{\sigma}_{\alpha}^2}\right) \exp\left(\sum_{\alpha=1}^d \frac{2\hat{\mu}_{\alpha y}x_{\alpha} - \hat{\mu}_{\alpha y}^2}{2\hat{\sigma}_{\alpha}^2}\right)$$

$$\frac{\hat{\pi}_y}{Z'} \quad \text{i.e. } \frac{1}{Z'} = \frac{1}{Z} \exp\left(\sum_{\alpha=1}^d -\frac{x_{\alpha}^2}{2\hat{\sigma}_{\alpha}^2}\right)$$

$$= \frac{\hat{\pi}_y}{Z'} \exp\left(\sum_{\alpha=1}^d \frac{\hat{\mu}_{\alpha y}x_{\alpha}}{\hat{\sigma}_{\alpha}^2}\right) \cdot \exp\left(\sum_{\alpha=1}^d -\frac{\hat{\mu}_{\alpha y}^2}{2\hat{\sigma}_{\alpha}^2}\right)$$

$$\text{let } \beta_y = \hat{\pi}_y \cdot \exp\left(\sum_{\alpha=1}^d -\frac{\hat{\mu}_{\alpha y}^2}{2\hat{\sigma}_{\alpha}^2}\right)$$

$$= \frac{\beta_y}{Z'} \cdot \exp\left(\sum_{\alpha=1}^d \frac{\hat{\mu}_{\alpha y}}{\hat{\sigma}_{\alpha}^2} \cdot x_{\alpha}\right)$$

$$\hat{P}(y|x) = \frac{\beta_y}{Z'} \exp\left(\sum_{\alpha=1}^d \frac{\hat{\mu}_{y\alpha}}{\hat{\sigma}_\alpha^2} \cdot x_\alpha\right)$$

let $W \in \mathbb{R}^{C \times d}$
 s.t. $W_{y\alpha} = \frac{\hat{\mu}_{y\alpha}}{\hat{\sigma}_\alpha^2}$

$$= \frac{\beta_y}{Z'} \cdot \exp\left((Wx)_y\right)$$

$C = \# \text{ of classes}$
 $= |Y|$

$$= \frac{1}{Z'} \exp\left((Wx+b)_y\right)$$

let $b \in \mathbb{R}^C$
 s.t. $b_y = \log \beta_y$

Know: $\sum_{y=1}^C \hat{P}(y|x) = 1 = \sum_{y=1}^C \frac{1}{Z'} \exp\left((Wx+b)_y\right)$

$$\Rightarrow Z' = \sum_{y=1}^C \exp\left((Wx+b)_y\right)$$

$$\hat{P}(y|x) = \frac{\exp\left((Wx+b)_y\right)}{\sum_{j=1}^C \exp\left((Wx+b)_j\right)}$$

$(\text{Softmax}(u))_i = \frac{\exp(u_i)}{\sum_{j=1}^C \exp(u_j)}$
 $u \in \mathbb{R}^d$

$$\hat{P}(y|x) = \text{softmax}(Wx+b)_y$$

$h(x) = \underset{c \in Y}{\text{arg max}} \text{softmax}\left((Wx+b)_c\right)$

If we have just 2 classes:

$= \underset{c \in Y}{\text{arg max}} (Wx+b)_c$

$$\hat{P}(y=1|x) = \frac{\exp\left((Wx+b)_1\right)}{\exp\left((Wx+b)_1\right) + \exp\left((Wx+b)_2\right)} = \frac{1}{1 + \exp\left((Wx+b)_2 - (Wx+b)_1\right)}$$

now let $w \leftarrow w_2 - w_1$
 $w \in \mathbb{R}^d$

let $b \leftarrow b_2 - b_1$
 $b \in \mathbb{R}$

$$= \frac{1}{1 + \exp(w^T x + b)} = \sigma(w^T x + b)$$

$\sigma(z) = \frac{1}{1 + e^z}$

predict $y=1$, when $w^T x + b \leq 0$.

Discriminative modeling: model just $P(y|x)$.

Logistic Regression. $y \in \{-1, 1\}$. Params $w \in \mathbb{R}^d$

Assume:

$$P(y|x) = \frac{1}{1 + \exp(-y(w^T x + b))}.$$

if $x \leftarrow \begin{bmatrix} x \\ 1 \end{bmatrix}$ $w \leftarrow \begin{bmatrix} w \\ b \end{bmatrix}$ $w^T x \leftarrow w^T x + b$

so: $P(y|x) = \frac{1}{1 + \exp(-y(w^T x))}$

$$P(y|x) = \frac{1}{1 + \exp(-y w^T x)}.$$

$h(x) = \text{sign}(w^T x)$. \Leftarrow a linear model.

MAP MLE estimate for w :

$$\max_{w \in \mathbb{R}^d} \text{arg max } P(D; w)$$

$$\hat{w} = \arg \max_{w \in \mathbb{R}^d} P(D; w)$$

$$= \arg \max_{w \in \mathbb{R}^d} P((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n); w)$$

$$= \arg \max_{w \in \mathbb{R}^d} \prod_{i=1}^n P((x_i, y_i); w)$$

$$= \arg \max_{w \in \mathbb{R}^d} \prod_{i=1}^n \underbrace{P(y_i | x_i; w)}_{\text{no distribution}} / \underbrace{P(x_i; w)}_{\text{not modeling}}$$

$$\hat{w}_{\text{MLC}} = \arg \max_{w \in \mathbb{R}^d} P(y | x; w) \quad (\text{discriminative MLE})$$

$$= \arg \max_{w \in \mathbb{R}^d} \prod_{i=1}^n P(y_i | x_i; w)$$

$$= \arg \max_{w \in \mathbb{R}^d} \prod_{i=1}^n \frac{1}{1 + \exp(-y_i x_i^T w)} \quad \log \frac{1}{x} = -\log x$$

$$= \arg \max_{w \in \mathbb{R}^d} \sum_{i=1}^n -\log(1 + \exp(-y_i x_i^T w))$$

$$\hat{w}_{\text{MLE}} = \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w))$$

logistic regression loss.