

Naive Bayes

Recap: MLE: model $P(x, y)$ with P_θ where $\theta \in \Theta$

$$\hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} P_\theta(D)$$

model θ that maximizes likelihood of data

MAP:
$$\hat{\theta}_{MAP} = \underset{\theta \in \Theta}{\operatorname{argmax}} P(\theta | D) = P(D | \theta) P(\theta)$$

\downarrow likelihood \downarrow Prior

Eg	Y	X = Favorite dish	X = {Soup, Mac N cheese, Tacos}
	Adult	Soup	
	child	Mac N Cheese	Estimate
	child	Mac N Cheese	$P(Y = \text{"child"} X = \text{"Mac or cheese"})?$
	Adult	Tacos	
	Adult	Soup	
	Child	Tacos	
	Adult	Soup	
	Adult	Mac N Cheese	
	child	Mac N Cheese	

Estimate $P(Y = y | X = x) ?$

What is the issue with this?

Eg	Y	$X(1)$ = Favorite dish	$X(2)$ = # words known	$X(3)$ = Favorite Movie	$X(4)$ = Hours of sleep
	Adult	Soup	20000	Godfather	8
	child	Mac N Cheese	200	Frozen	11
	child	Mac N Cheese	400	Frozen	12
	Adult	Tacos	17000	La la land	6
	Adult	Soup	15000	Godfather	5
	Child	Tacos	1000	Eternals	10
	Adult	Soup	21000	Avengers	10
	Adult	Mac N Cheese	11000	Avengers	8
	child	Mac N Cheese	700	Avengers	11

$\hat{P}(Y = \text{"Adult"} | X = (\text{"soup"}, 20000, \text{"Avengers"}, 8)) ?$

Naive Bayes Model

$$\text{Assumption: } P(X = \vec{x} | Y = y) = \prod_{\alpha=1}^d P(X^{(\alpha)} = x^{(\alpha)} | Y = y)$$

"Given it's a child, favorite dish, # words known, hours of sleep. are all independent"

why is this useful?

$$\begin{aligned} P(Y = y | X = x) &= \frac{P(X = x | Y = y) P(Y = y)}{P(X = x)} \\ &= \frac{\prod_{\alpha=1}^d P(X^{(\alpha)} = x^{(\alpha)} | Y = y) P(Y = y)}{P(X = x)} \\ &= \frac{\prod_{\alpha=1}^d P(X^{(\alpha)} = x^{(\alpha)} | Y = y) P(Y = y)}{\sum_{c \in \mathcal{Y}} \prod_{\alpha=1}^d P(X^{(\alpha)} = x^{(\alpha)} | Y = c) P(Y = c)} \end{aligned}$$

$P(Y = y)$ and $\forall \alpha, P(X^{(\alpha)} = x^{(\alpha)} | Y = y)$
are easy to estimate

Eg. estimate

$$\hat{P}(Y = \text{"Adult"} | X = (\text{"soup"}, 20000, \text{"Avengers"}, 8)) ?$$

$$\begin{aligned} h(x) &= \operatorname{argmax}_{y \in \mathcal{Y}} \hat{P}(Y = y | X = x) \\ &= \operatorname{argmax}_{y \in \mathcal{Y}} \prod_{\alpha} P(X^{(\alpha)} = x^{(\alpha)} | Y = y) P(Y = y) \\ &= \operatorname{argmax}_{y \in \mathcal{Y}} \sum_{\alpha} \log(P(X^{(\alpha)} = x^{(\alpha)} | Y = y)) \\ &\quad + \log P(Y = y) \end{aligned}$$

when $X_j(\mathcal{d})$ are counts

Eg $X(\mathcal{d}) = \cdot$ means α^{th} word in the dictionary occurs j times in the document

x is an m word document: $X(\mathcal{d}) \in \{0, 1, \dots, m\}$ $\sum_{\alpha=1}^d X(\mathcal{d}) = m$

Multinomial Distribution

$$P(X=x | m, Y=y) = \frac{m!}{X(1)! X(2)! \dots X(d)!} \prod_{i=1}^d (\theta_{\alpha,y})^{X(\mathcal{d})}$$

MLE estimate:
$$\hat{\theta}_{\alpha,y} = \frac{\sum_{i=1}^n \mathbb{1}\{y_i=y\} X_i(\mathcal{d})}{\sum_{i=1}^n \mathbb{1}\{y_i=y\} m_i}$$

$m_i = \#$ words in document i

$$h(x) = \operatorname{argmax}_{y \in Y} P(Y=y) \prod_{\alpha=1}^d \hat{\theta}_{\alpha,y}^{X(\mathcal{d})}$$

$X(\mathcal{d})$'s are continuous variables.: Gaussian distribution conditioned on Y

$$p(X(\mathcal{d})=x | Y=y) = \frac{1}{\sqrt{2\pi}\sigma_y} \exp\left(-\frac{(x-\mu_y)^2}{2\sigma_y^2}\right)$$

Parameter estimation:

$$\hat{\mu}_{y(\omega)} = \frac{\sum_{i=1}^n \mathbb{1}\{y_i=y\} X_i(\mathcal{d})}{\sum_{i=1}^n \mathbb{1}\{y_i=y\}} \quad \hat{\sigma}_y^2 = \frac{\sum_{i=1}^n \mathbb{1}\{y_i=y\} (X_i(\mathcal{d}) - \hat{\mu}_{y(\omega)})^2}{\sum_{i=1}^n \mathbb{1}\{y_i=y\}}$$

1. For both multinomial case and Gaussian case (with variance between class per feature fixed) classification boundary is linear.

2. For Gaussian case

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

logistic link function.

