

```
ocinar1@remote04:~/Desktop/ML2$ python3 NaiveBayes.py train/spam train/ham test/spam  
test/ham
```

Accuracy for training:	98.0561555075594
Accuracy with stop words:	94.97907949790795
Accuracy without stop words:	94.35146443514644

Accuracy for training:	98.0561555075594
Accuracy with stop words:	94.97907949790795
Accuracy without stop words:	94.35146443514644

①

① Derivative maximum likelihood estimates

a) parameter : p Bernoulli (p) → sample of size n

$$\text{Binomial} \rightarrow f(x, p) = \begin{cases} p^x (1-p)^{n-x} & ; x=0,1,2 \\ 0 & ; \text{else} \end{cases}$$

$x_i \rightarrow \text{Bern}(p)$ and x_i independent then f :

$$\log \text{likelihood } L \rightarrow L(p) = \ln \prod_{i=1}^n f(x_i, p)$$

$$L(p) = \sum_{i=1}^n \ln [p^{x_i} (1-p)^{1-x_i}]$$

$$L(p) = \sum_{i=1}^n [x_i \ln p + (1-x_i) \ln (1-p)]$$

we are going to maximize this

$$\underset{p}{\text{arg max}} \ L(p) \Rightarrow \frac{\partial L(p)}{\partial p} = 0$$

$$\frac{\partial}{\partial p} = \left[\sum_{i=1}^n x_i \ln(p) + (1-x_i) \ln(1-p) \right] = 0$$

$$\sum_{i=1}^n x_i \cdot \frac{1}{p} + \sum_{i=1}^n (1-x_i) \cdot \frac{1}{1-p} \cdot (-1) = 0$$

$$\left(\sum_{i=1}^n x_i \right) \left(\frac{1}{p} + \frac{1}{1-p} \right) + \sum_{i=1}^n \frac{-1}{1-p} = 0$$

$$\frac{1}{p(1-p)} \sum_{i=1}^n x_i - \frac{n}{1-p} = 0$$

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n x_i$$

(2)

b) parameter p , Binomial(N, p) sample size n

$$\text{Binomial dist} \rightarrow f(x; N, p) = \binom{N}{x} p^x (1-p)^{N-x}$$

$$\hat{\prod}_{i=1}^n f(x_i) = \left(\prod_{i=1}^n \binom{N}{x_i} \right) p^{\sum x_i} (1-p)^{nN - \sum x_i}$$

$$L(p) = \ln \left(\prod_{i=1}^n f(x_i) \right) = \ln \left(\prod_{i=1}^n \binom{N}{x_i} \right) + \left(\sum_{i=1}^n x_i \right) \ln(p) + \\ + \left(nN - \sum_{i=1}^n x_i \right) \ln(1-p)$$

$$\Rightarrow \frac{dL(p)}{dp} = 0 \Rightarrow \left(\sum_{i=1}^n x_i \right) \left(\frac{1}{p} \right) + \left(nN - \sum_{i=1}^n x_i \right) \frac{1}{1-p} (-1) = 0$$

$$\sum_{i=1}^n x_i \left[\frac{1}{p} + \frac{1}{1-p} \right] - \frac{nN}{1-p} = 0 \Rightarrow \left(\sum_{i=1}^n x_i \right) \frac{1}{p(1-p)} - \frac{nN}{1-p} = 0$$

$$\Rightarrow \left(\sum_{i=1}^n x_i \right) \frac{1}{p} - nN = 0 \Rightarrow \hat{p} = \frac{1}{nN} \sum_{i=1}^n x_i$$

sample $\Rightarrow (3, 6, 2, 0, 0, 3) \quad N=10 \quad n=6$

$$\sum_{i=1}^n x_i = 3+6+2+0+0+3 = 14$$

$$\hat{p} = \frac{1}{6 \cdot 10} \cdot 14 = 0.2333$$

③

c) parameters a, b on a Uniform (a, b) samp. size n

$$\text{Unif dist} \rightarrow f(x_i, a, b) = \begin{cases} \frac{1}{b-a} & a \leq x_i \leq b \\ 0 & \text{else} \end{cases}$$

$$L(\theta) = \ln \prod_{i=1}^n f(x_i, a, b) \rightarrow \ln \left(\prod_{i=1}^n \frac{1}{b-a} \right) = \ln \frac{1}{(b-a)^n} = -n \ln(b-a)$$

$$\frac{\partial L(\theta)}{\partial a} = \frac{n}{b-a} \rightarrow \text{it is monotonically increasing}$$

$$\frac{\partial L(\theta)}{\partial b} = \frac{-n}{b-a} \rightarrow \text{it is monotonically decreasing}$$

$$\Rightarrow \hat{a} = \min(x_1, x_2, x_3, \dots, x_n) \quad \hat{b} = \max(x_1, x_2, x_3, \dots, x_n)$$

$$\text{d) } f(x_i, \mu) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

$$\Rightarrow \prod_{i=1}^n f(x_i, \mu) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\sum_{i=1}^n \frac{(x_i-\mu)^2}{2\sigma^2}}$$

$$\Rightarrow L(\mu) = \ln \left[\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\sum_{i=1}^n \frac{(x_i-\mu)^2}{2\sigma^2}} \right] = \sum_{i=1}^n \left[-\ln \sqrt{2\pi\sigma^2} - \frac{1}{2\sigma^2} (x_i-\mu)^2 \right]$$

$$\frac{\partial L(\mu)}{\partial \mu} = 0 \Rightarrow \sum_{i=1}^n 2(-1)(x_i-\mu) \cdot \frac{(-1)}{2\sigma^2} = 0 \Rightarrow \sum_{i=1}^n \frac{(x_i-\mu)}{\sigma^2} = 0$$

$$\sum_{i=1}^n (x_i - \mu) = 0 \Rightarrow \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

3

$$e) L(\theta) = \sum_{i=1}^n \left[-\ln(\sqrt{2\pi} \sigma) - \frac{1}{2\sigma^2} (x_i - \mu)^2 \right] \Rightarrow \frac{\partial L(\theta)}{\partial \sigma} = 0$$

$$\Rightarrow \sum_{i=1}^n \left[\frac{-1}{\sigma} + \frac{1}{\sigma^3} (x_i - \mu)^2 \right] = 0 \Rightarrow \frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2 = 0$$
$$\Rightarrow \sum_{i=1}^n (x_i - \mu)^2 = n\sigma^2 \Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

$$f) f(x_i, \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$L(\theta) = \sum_{i=1}^n \left[-\frac{1}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} (x_i - \mu)^2 \right]$$

$$\left. \begin{array}{l} \frac{\partial L(\theta)}{\partial \mu} = 0 \Rightarrow \sum_{i=1}^n \frac{1}{2\sigma^2} \cdot (-1)(x_i - \mu) = 0 \\ \frac{\partial L(\theta)}{\partial \sigma^2} = \sum_{i=1}^n \left[-\frac{1}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} (x_i - \mu)^2 \right] = 0 \end{array} \right\}$$

$$\Rightarrow \begin{cases} \sum_{i=1}^n (x_i - \mu) = 0 \\ \frac{-n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 = 0 \end{cases}$$

$$\Rightarrow \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

⑤

②

$$a) mLE \rightarrow P(H) = \theta \quad P(T) = 1 - \theta \Rightarrow \theta \in [0, 1]$$

$$P(D|\theta) = \theta^{n_H} (1-\theta)^{n_T}$$

$$\Rightarrow L(\theta) = \ln(P(D|\theta)) = \ln[\theta^{n_H} (1-\theta)^{n_T}] = n_H \ln \theta + n_T \ln(1-\theta)$$

$$\frac{\partial L(\theta)}{\partial \theta} = 0 = \frac{n_H}{\theta} - \frac{n_T}{1-\theta} = 0 \Rightarrow n_H(1-\theta) = n_T \theta$$

$$\hat{\theta}_{MLE} = \frac{n_H}{n_H + n_T}$$

$$MAP \rightarrow Bayes \rightarrow \frac{P(D|\theta)P(\theta)}{P(D)} = P(\theta|D)$$

$$\theta_{MAP} = \operatorname{argmax}_{\theta} P(\theta|D) = \operatorname{argmax}_{\theta} \frac{P(D|\theta)P(\theta)}{P(D)} = \operatorname{argmax}_{\theta} P(D|\theta)P(\theta)$$
$$= \operatorname{argmax}_{\theta} \prod_{i=1}^n P(x_i|\theta)P(\theta)$$

$$\operatorname{argmax}_{\theta} P(\theta|D) = \operatorname{argmax}_{\theta} \log(P(\theta|D)) = \operatorname{argmax}_{\theta} \sum_{i=1}^n \log P(x_i|\theta)P(\theta)$$

$$\Rightarrow L(\theta) = \sum_{i=1}^n [\ln \text{Bernoulli}(x_i|\theta) + \ln \text{Beta}(\theta|\alpha, \beta)]$$

$$\frac{\partial L(\theta)}{\partial \theta} = \sum_{i=1}^n \left[\frac{\partial}{\partial \theta} \ln \text{Bernoulli}(x_i|\theta) + \frac{\partial}{\partial \theta} \ln \text{Beta}(\theta|\alpha, \beta) \right]$$

$$\text{First} \rightarrow \sum_{i=1}^n \frac{\partial \text{Bernoulli}(x_i|\theta)}{\partial \theta} = \frac{1}{\theta} \sum_{i=1}^n x_i - \frac{1}{1-\theta} \sum_{i=1}^n (1-x_i)$$

$$\text{Second} \rightarrow \frac{\partial}{\partial \theta} \ln \text{Beta}(\theta|\alpha, \beta) = \frac{\partial}{\partial \theta} \ln \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \cdot \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$= \frac{\partial}{\partial \theta} \ln \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} + \frac{\partial}{\partial \theta} \ln \theta^{\alpha-1} (1-\theta)^{\beta-1} = \frac{\partial}{\partial \theta} \ln \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$= \frac{\partial}{\partial \theta} (\alpha-1) \frac{\partial}{\partial \theta} \ln \theta + (\beta-1) \frac{\partial}{\partial \theta} \ln (1-\theta) = \frac{\alpha-1}{\theta} - \frac{\beta-1}{1-\theta}$$

(b)

$$\hat{\theta}_{\text{MAP}} = \frac{\partial \ell(\theta)}{\partial \theta} = 0$$

$$= \frac{1}{\theta} \sum_{i=1}^n x_i - \frac{1}{1-\theta} \sum_{i=1}^n (1-x_i) + \frac{\alpha-1}{\theta} - \frac{\beta-1}{1-\theta} = 0$$

$$\Rightarrow \theta \left[\sum_{i=1}^n (1-x_i) + \beta - 1 \right] = (1-\theta) \left[\sum_{i=1}^n x_i + \alpha - 1 \right]$$

$$\Rightarrow \theta \left[\sum_{i=1}^n (1-x_i) + \sum_{i=1}^n x_i + \beta - 1 + \alpha - 1 \right] = \sum_{i=1}^n x_i + \alpha - 1$$

$$\Rightarrow \theta \left[\sum_{i=1}^n 1 + \beta + \alpha - 2 \right] = \sum_{i=1}^n x_i + \alpha - 1$$

$$\Rightarrow \text{let head}=1 \text{ and tails}=0 \Rightarrow \sum_{i=1}^n x_i = n_H$$

$$\theta \left[\sum_{i=1}^n 1 + \beta + \alpha - 2 \right] = \sum_{i=1}^n x_i + \alpha - 1$$

$$\theta \left[n + \beta + \alpha - 2 \right] = n_H + \alpha - 1$$

$$\Rightarrow \hat{\theta}_{\text{MAP}} = \frac{n_H + \alpha - 1}{n + \beta + \alpha - 2}$$

b) thumbtack

n	n_H	α	β	$\hat{\theta}_{\text{MLE}}$	$\hat{\theta}_{\text{MAP}}$
100	30	100	100	0.3	0.32
100	30	1000	1000	0.3	0.31
100	30	10000	10000	0.3	0.30

coin \rightarrow

n	n_H	α	β	$\hat{\theta}_{\text{MLE}}$	$\hat{\theta}_{\text{MAP}}$
100	60	1	1	0.6	0.6
100	60	40	60	0.6	0.5
100	60	30	70	0.6	0.45

$$\hat{\theta}_{\text{MLE}} = \frac{n_H}{n} \text{ ad}$$

$$\hat{\theta}_{\text{MAP}} = \frac{n_H + \alpha - 1}{n + \beta + \alpha - 2}$$

$$\text{In coin } \rightarrow \text{if } \beta = 1 \text{ and } \alpha = 1 \Rightarrow \hat{\theta}_{\text{MLE}} = \hat{\theta}_{\text{MAP}}$$

* so \rightarrow MLE is a special case of MAP when the prior is uniform,

(2)

c) False \Rightarrow MLE approaches MAP in just one case
which is $\alpha = \beta = 1$

d) True \Rightarrow example \Rightarrow Cth $\Rightarrow n_+ = 60; n = 100 \alpha = \beta = 100.000$

$$\hat{\theta}_{MLE} = 0.6 \text{ and } \hat{\theta}_{MAP} = 0.6$$

$$\text{threshold} = 1 \quad n_+ = 40; n = 100 \quad \alpha = \beta = 100.000$$

$$\hat{\theta}_{MLE} = 0.4 \text{ and } \hat{\theta}_{MAP} = 0.5$$

when we use large prior for both of them, they will have the same MAP