

(1)

CSE575: SML Assignment 1

92 / 15 / 2019

(ii) Given,

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

here $\hat{\sigma}^2$ $\hat{\sigma}^2$ - sample \rightarrow variance $\hat{\mu}$ - sample mean.

The actual or population standard deviation is always greater than sample standard deviation.

 μ - population mean \bar{x}_{MLE}^2 = $\hat{\mu}_{MLE}$ or Sample mean $\hat{\sigma}^2$ = population std variance (unbiased) $\hat{\sigma}^2$ = $\hat{\sigma}_{MLE}^2$ or sample variance [biased]

$$\hat{\mu}_{MLE} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

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

$$E[\hat{\sigma}^2] = E\left[\frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_{MLE})^2\right]$$

$$E[\sigma^2 - \hat{\sigma}^2] = E\left[\sum ((x_i - \mu)^2 - (x_i - \hat{\mu})^2)\right]$$

(2)

$$= \frac{1}{n} E \left[\sum_{i=1}^n \left((\hat{x}_i^2 - 2\hat{x}_i u + u^2) - (\hat{x}_i^2 - 2\hat{x}_i \hat{u} + \hat{u}^2) \right) \right]$$

$$= \frac{1}{n} E \left[\sum_{i=1}^n (u - \hat{u} - 2\hat{x}_i u + 2\hat{x}_i \hat{u}) \right]$$

$$= \frac{1}{n} E \left[n\hat{u}^2 - n\hat{u} - 2 \sum_{i=1}^n \hat{x}_i (u - \hat{u}) \right]$$

$$= E \left[\hat{u}^2 - \hat{u} - 2(u - \hat{u}) \sum_{i=1}^n \hat{x}_i \right]$$

$$= E \left[\hat{u}^2 - \hat{u} - 2(u - \hat{u}) \hat{u} \right]$$

$$= E \left[\hat{u}^2 - \hat{u} - 2u\hat{u} + 2\hat{u}^2 \right]$$

$$= E \left[\hat{u}^2 + \hat{u}^2 - 2u\hat{u} \right]$$

$$= E \left[(\hat{u} - u)^2 \right]$$

$$= \text{Var}(\hat{u})$$

$$= \frac{\sigma^2}{n}$$

$$\therefore \text{Var}(\hat{u}) = \frac{\sigma^2}{n}$$

$$\begin{aligned} \text{Var}(\hat{u}) &= \text{Var} \left(\frac{1}{n} \sum_{i=1}^n x_i \right) = \frac{1}{n^2} \sum \text{Var}(x_i) \\ &= \frac{1}{n^2} \times n\sigma^2 = \frac{\sigma^2}{n} \end{aligned}$$

(3)

$$\therefore E[\hat{\sigma}^2] = \sigma^2 - \frac{\sigma^2}{n}$$

$$= \sigma^2 \left[1 - \frac{1}{n} \right]$$

$$\therefore \hat{\sigma}^2 = \sigma^2 \frac{(n-1)}{n}$$

Thus we see that
 $\frac{n-1}{n} \neq \sigma^2$

which means variance is biased
 ie the unbiased variance is
 slightly less than than the actual
 variance

\therefore we divide with $(n-1)$ instead
 of n while calculating standard
 deviation.

This "correction" term is also
 called as Bessel's correction

② Given

lever distⁿ for mean $\mu \sim N(0, 1)$

$$U_{app} = ?$$

$$P(D|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(n-\mu)^2}{2\sigma^2}}$$

$$P(\mu|0, 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(\mu-0)^2}{2\cdot 1^2}}$$

Now,

$$P(\mu|D) = P(D|\mu) P(\mu) \quad (1)$$

(4)

plus) in the posterior which we have to minimize

$$P(D|u, \sigma) = \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x_i - u)^2}{2\sigma^2}}$$

$$P(u|\sigma, D) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(u - \bar{x})^2}{2\sigma^2}}$$

Now we take \log & set derivative w.r.t u to 0.

$$\therefore \ln P(u|D) = \ln [P(D|u, \sigma) P(u|\sigma, D)]$$

$$= \ln P(D|u, \sigma) + \ln P(u) \quad (2)$$

~~= const~~

Now,

$$\ln P(D|u, \sigma) = \ln \frac{1}{\sigma \sqrt{2\pi}} + \ln \left[\prod_{i=1}^n e^{-\frac{(x_i - u)^2}{2\sigma^2}} \right]$$

$$= C_1 + \sum_{i=1}^n \frac{(x_i - u)^2}{2\sigma^2} \quad (3)$$

$$\ln P(u) = \ln \frac{1}{2\sqrt{2\pi}} + \frac{(u - \bar{x})^2}{2\sigma^2} \quad (4)$$

gives λ not λ^2 since here λ is the variance

Taking derivative of ②

$$\frac{\partial}{\partial \mu} [\ln P(D|\mu, \sigma) + \ln P(\mu)] = 0$$

From ③ & ④

$$\frac{\partial}{\partial \mu} \left[-\frac{\sum (x_i - \mu)^2}{2\sigma^2} - \frac{(\mu - \theta)^2}{2\lambda} \right] = 0$$

$$\therefore \frac{\sum (x_i - \mu)}{\sigma^2} = \frac{(\mu - \theta)}{\lambda}$$

$$\therefore \frac{\sum x_i}{\sigma^2} - \frac{\sum \mu}{\sigma^2} = \frac{\mu}{\lambda} - \frac{\theta}{2\lambda}$$

$$\therefore \frac{\sum x_i}{\sigma^2} + \frac{\theta}{2\lambda} = \mu \left[\frac{1}{\lambda} + \frac{N}{\sigma^2} \right]$$

$$\therefore \mu = \frac{\frac{\sum x_i}{\sigma^2} + \frac{\theta}{2\lambda}}{\frac{1}{\lambda} + \frac{N}{\sigma^2}}$$

$$\therefore \hat{\mu}_{MAP} = \frac{\frac{\sum x_i}{\sigma^2} + \frac{\theta}{2\lambda}}{\frac{1}{\lambda} + \frac{N}{\sigma^2}}$$

(6)

Here, we can see,

the terms $\frac{\theta}{\sigma^2} \ln \frac{1}{\theta}$ same

from prior, if we ignore
them

$$\text{we get } \frac{\frac{\sum w_i}{\sigma^2}}{\frac{N}{\sigma^2}}$$

$$= \frac{1}{N} \sum w_i$$

which is MLE

we can see adding Prior to
MLE gives MAP.

If $N \rightarrow \infty$, above terms would
vanish (negligible) giving

$\text{MLE} = \text{MAP}$
i.e (Posterior = Likelihood)

(62)

$$P(k|z) = \frac{2^k e^{-2}}{L_k}$$

$\lambda_1, \lambda_2, \dots, \lambda_n \rightarrow$ Po samples.

$$P(\text{data} | \lambda) = \prod_{i=1}^n P(\lambda_i | \lambda)$$

$$= \prod_{i=1}^n \frac{\lambda^{\lambda_i}}{\lambda^{\sum \lambda_i}} e^{-\lambda}$$

Taking log.

$$\ln \prod \left(\frac{\lambda^{\lambda_i}}{\lambda^{\sum \lambda_i}} e^{-\lambda} \right)$$

$$= \sum_{i=1}^n \ln \left(\frac{\lambda^{\lambda_i}}{\lambda^{\sum \lambda_i}} e^{-\lambda} \right)$$

$$= \sum_{i=1}^n [\ln \lambda^{\lambda_i} - \lambda - \log (\sum \lambda_i)]$$

Solving derivative wrt $\lambda = 0$

$$\frac{d}{d\lambda} \sum_{i=1}^n [\ln \lambda^{\lambda_i} - \lambda - \log (\sum \lambda_i)] = 0$$

$$= \sum_{i=1}^n \frac{\lambda_i}{\lambda} - \sum_{i=1}^n 1 - 0 = 0$$

(8)

$$\therefore \sum_{i=1}^n \frac{k_i}{\lambda} - n = 0$$

$$\therefore \sum k_i = n\lambda$$

$$\therefore \lambda = \frac{\sum k_i}{n} \quad \boxed{i=1}^n$$

$$\therefore \hat{\lambda}_{MLE} = \frac{1}{n} \sum_{i=1}^n k_i$$

(2)

Poisson

$$P(k|\lambda) = \frac{\bar{e}^{-\lambda} \lambda^k}{k!}$$

$$E(X) = \sum k P(k|\lambda)$$

$$= \sum k \bar{e}^{-\lambda} \frac{\lambda^k}{k!}$$

$$= \sum k \bar{e}^{-\lambda} \frac{\lambda^k}{(k-1)!}$$

$$= \bar{e}^{-\lambda} \sum k \cdot 2 \cdot \frac{\lambda^{k-1}}{(k-1)!}$$

$$= \bar{e}^{-\lambda} \sum \frac{\lambda^{k-1}}{(k-1)!}$$

(9)

$$\therefore \bar{c}^2 \lambda \leq \frac{\lambda^{k+1}}{k+1}$$

Now, $\frac{\bar{c}^2 \lambda^{k+1}}{k+1} = \bar{c}^2$ for $k \geq 1$

$$\therefore \bar{c}^2 \lambda e^2$$

$$= \lambda$$

$$\therefore E(\lambda) = \lambda.$$

(3)

① Here, no. of features = 3
(size, color, shape)

without conditional independent assumption, we have

$$P(Y|X) = P(X_1, X_2, X_3 | Y) P(Y)$$

$$\text{So for } P(X_1, X_2, X_3 | Y)$$

we have to estimate $2(2^n) = 2^{n+1}$

Parameters, here $n=3$

$2^n \rightarrow$ Every X has 2 values,
This gives total combinations
of values.

$2 \times 2^n \rightarrow$ because 1 can have 2
values.

10

$P(Y|X) = 1$ parameter.
since Y can have 2 values.

With conditional independence, we have,

$$P(Y|X) = P(X_1, X_2, X_3|Y) = P(X_1|Y) P(X_2|Y) P(X_3|Y)$$

X has 2 values for each feature.

where we only need to estimate $P(X_1|Y)$ $P(X_2|Y)$ $P(X_3|Y)$

for both values of Y

which will be $3 \times 2 = 6$ parameters.

& $P(Y) \rightarrow 1$ Parameter

\therefore totally $- 6 + 1$

= 7 parameters.

where

$X_1 = \text{size}$

$P(X_1|Y=\text{yes})$

$X_2 = \text{shape}$

$P(X_1|Y=\text{no})$

$X_3 = \text{color}$

$P(X_2|Y=\text{yes})$

$Y = \text{good_apple}$

$P(X_2|Y=\text{no})$

$P(X_3|Y=\text{yes})$

$P(X_3|Y=\text{no})$

Here since only 2 values of X are there for every X , we can estimate other $P(Y)$

② For MLE, if one is known, we try to maximize

$P(X|Y)$ \rightarrow likelihood.

$$P(X_1, X_2, X_3|Y) = \prod_{i=1}^3 P(X_i|Y)$$

Prior:

$$P(Y=\text{yes}) = 4/10$$

$$P(Y=\text{no}) = 6/10$$

? from the table

MLE estimate

$$\hat{\pi}_k = \hat{p}(y=y_k)$$

where ~~y~~ $\in \{y_k \in \{\text{yes, no}\}$
 $0 \leq k \leq 1$

$$\hat{\pi}_k = \frac{\# D\{y=y_k\}}{|D|}$$

where D is the data

$$\hat{\pi}_0 = \frac{\# D\{y=y_0\}}{|D|} \quad \text{where } y_0 = \text{yes}$$

$y_1 = \text{no}$

$$= \frac{4}{10}$$

$$\hat{\pi}_1 = \frac{\# D\{y=y_1\}}{|D|}$$

$$= \frac{6}{10}$$

Now, we estimate cond probability

$$\hat{o}_{ijk} = \hat{p}(x_i = x_{ij} | y=y_k)$$

$$= \frac{\# D\{x_i = x_{ij} | y=y_k\}}{\# D\{y=y_k\}}$$

Now, $j = \text{Sample}$
 $j = \text{feature}$
 $k = \text{class}$

For $y = \text{Yes}$,

$$P(x_1 = \text{small} | y = \text{Yes}) = \frac{|x_1 = \text{small} \wedge y = \text{Yes}|}{|y = \text{Yes}|}$$

$$= \frac{1}{4}$$

$$P(x_1 = \text{large} | y = \text{Yes}) = \frac{3}{4}$$

$$P(x_2 = \text{green} | y = \text{Yes}) = 0$$

$$P(x_2 = \text{red} | y = \text{Yes}) = \frac{4}{4} = 1$$

$$P(x_3 = \text{irregular} | y = \text{Yes}) = \frac{1}{4}$$

$$P(\cancel{x_3 = \text{Regular}} | \cancel{y = \text{Yes}}) = ?$$

$$P(x_3 = \text{circle} | y = \text{Yes}) = \frac{3}{4}$$

For $y = \text{No}$,

$$P(x_1 = \text{small} | y = \text{No}) = \frac{3}{6} = \frac{1}{2}$$

$$P(x_1 = \text{large} | y = \text{No}) = \frac{3}{6} = \frac{1}{2}$$

$$P(X_2 = \text{green} | Y = \text{No}) = \frac{5}{6} = \frac{3}{5} = \frac{5}{6}$$

$$P(X_2 = \text{Red} | Y = \text{No}) = \frac{1}{6}$$

$$P(X_3 = \text{Irregular} | Y = \text{No}) = \frac{4}{6} = \frac{2}{3}$$

$$P(X_3 = \text{Circle} | Y = \text{No}) = \frac{2}{6} = \frac{1}{3}$$

~~Q = f + g - h~~

② Given,

$x = (\text{small, Red, Circle})$

$$P(Y = \text{No} | x)$$

We know, the features are independent

$$\therefore P(Y = \text{No} | x) \propto P(x|y) P(y)$$

$$\propto P(x_1, x_2, x_3 | y) P(y)$$

where x_i is feature

$$= P(x_1 | y) P(x_2 | y) P(x_3 | y) P(y)$$

Since they are independent

$$= P(X_1 = \text{Small} | Y = \text{No})$$

$$P(X_2 = \text{Red} | Y = \text{No}) P(X_3 = \text{Circle} | Y = \text{No}) P(\text{No})$$

14

$$= \frac{1}{2} \times \frac{1}{6} \times \frac{1}{3} \times \frac{6}{10} = \frac{1}{36}$$

$$= \frac{1}{36} \times \frac{1}{10}$$

$$= \frac{1}{60} = 0.0167$$

$$P(Y = No | x) \approx C \times 0.0167$$

where C is "normalization" constant

$$= \frac{1}{P(x)}$$

Now,

$$\begin{aligned} P(Y = Yes | x) &= P(X_1 = Small | Yes) \times \\ &\quad P(X_2 = Red | Yes) \times \\ &\quad P(X_3 = Green | Yes) \times \\ &\quad P(Yes) \end{aligned}$$

$$= \frac{1}{6} \times 1 \times \frac{3}{4} \times \frac{4}{10}$$

$$= \frac{3}{40} = 0.075$$

Now

$$P(Y = No | x) = 0.0167$$

$$P(Y = Yes | x) = 0.075$$

We see that $P(Y = Yes | x) > P(Y = No | x)$

\therefore we ~~say~~ predict Yes for this sample of x where $x = \{Small, red, green\}$

(15)

(5)

$$\textcircled{1} \quad P(x_i | y = k) = N(\mu_{ik}, \sigma_{ik}^2)$$

$k = 0, 1$
 $d = 1 \dots d$

So, here

~~for~~

it will have

There will $2 \times d$ different μ $\therefore 0 \leq i \leq 1$ (2 values) $0 \leq k \leq d$ (d values)There will be $2 \times d$ different σ since σ is diff for classes & features ie. σ_{ik}

$$\begin{matrix} 0 \leq i \leq 1 \\ 1 \leq k \leq d \end{matrix}$$

$$\therefore 2d + 2d = 4d \quad \text{for } \mu \text{ & } \sigma \text{ ie.}$$

Gaussian ie.
likelihood funct'

$$P(x|y)$$

$P(y)$ = Bernoulli with one
 (P_{prior}) parameter θ

Total Parameters =

Parameters of prior +
 Parameters of likelihood
 $- 1 + 4d$

$- 4d + 1$ parameters
 where d are the features.

(16)

② Yes, we can translate w into ANB parameters with assumption.
Using Baye's

$$P(Y_1 = 1 | x) = \frac{P(x | Y=1) P(Y=1)}{P(x | Y=1) P(Y=1) + P(x | Y=0) P(Y=0)}$$

$$= \frac{1}{1 + \frac{P(x | Y=0)}{P(x | Y=1)} \frac{P(Y=0)}{P(Y=1)}}$$

Taking exp & log. = (doing nothing)

$$= \frac{1}{1 + \exp(\ln\left(\frac{P(x | Y=0)}{P(x | Y=1)} \frac{P(Y=0)}{P(Y=1)}\right))} \quad (1)$$

Now

Since $P(Y)$ is Bernoulli

$$P(Y=0) = \theta$$

$$P(Y=1) = 1 - \theta$$

Also Likelihood is Gaussian

$$\begin{aligned} P(x | Y=0) &= \prod P(x_i | Y=0) \\ P(x | Y=1) &= \prod P(x_i | Y=1) \end{aligned}$$

from (1)

$$= \frac{1}{1 + \exp\left(\sum \ln P(x_i | Y=0) + \ln\left(\frac{1-\theta}{\theta}\right)\right)} \quad (2)$$

Now since $P(x_i | Y=Y_k)$ is Gaussian
 $\sim N(\mu_{ik}, \sigma_{ik})$

(17)

we have,

$$P(x_i = n | Y=k) = \frac{1}{\sigma_{ik} \sqrt{2\pi}} e^{-\frac{(x - \mu_{ik})^2}{2\sigma_{ik}^2}}$$

$$\therefore P(x_i = n | Y=0) = \frac{1}{\sigma_{i0} \sqrt{2\pi}} e^{-\frac{(n - \mu_{i0})^2}{2\sigma_{i0}^2}}$$

$$P(x_i = n | Y=1) = \frac{1}{\sigma_{i1} \sqrt{2\pi}} e^{-\frac{(n - \mu_{i1})^2}{2\sigma_{i1}^2}}$$

\therefore from (2)
we have

$$\begin{aligned} & \frac{d}{1 + \text{exp}} \\ \therefore \ln \left(\frac{P(x_i = n | Y=0)}{P(x_i = n | Y=1)} \right) &= \ln \left(\frac{\sigma_{i1}}{\sigma_{i0}} e^{-\frac{(n - \mu_{i0})^2}{2\sigma_{i0}^2} + \frac{(n - \mu_{i1})^2}{2\sigma_{i1}^2}} \right) \\ &= \ln \left(\frac{\sigma_{i1}}{\sigma_{i0}} \right) - \frac{(n - \mu_{i0})^2}{2\sigma_{i0}^2} + \frac{(n - \mu_{i1})^2}{2\sigma_{i1}^2} \end{aligned}$$

Plugging them in (2)
we get

$$\begin{aligned} & \frac{1}{1 + \text{exp}} \left[2 \ln \left(\frac{\sigma_{i1}}{\sigma_{i0}} \right) + \frac{(n - \mu_{i0})^2}{2\sigma_{i0}^2} - \frac{(n - \mu_{i1})^2}{2\sigma_{i1}^2} \right] \\ &+ \ln \left(\frac{1 + \text{exp}}{\sigma} \right) \end{aligned}$$

$$\therefore \ln \frac{P(x_i | y=0)}{P(x_i | y=1)} = \sum_i \ln \left(\frac{\sigma_{i0}}{\sigma_{i1}} \right) +$$

$$\sum_i \left(\frac{(x_i - \mu_{i1})^2}{2\sigma_{i1}^2} + \frac{(x_i - \mu_{i0})^2}{2\sigma_{i0}^2} \right)$$

⇒

Now,

$$\sum_i \left(\frac{(x_i - \mu_{i1})^2}{2\sigma_{i1}^2} - \frac{(x_i - \mu_{i0})^2}{2\sigma_{i0}^2} \right) \quad \text{③}$$

$$= \sum_i f(x_i)$$

Here, we assume that

$\sigma_{i1} = \sigma_{i0} = \sigma_i$ which means
 variance is different for diff features but similar for features of diff classes.

i.e. it doesn't vary with class.
 from ③

$$\therefore \sum_i \left(\frac{(x_i - \mu_{i1})^2}{2\sigma_i^2} - \frac{(x_i - \mu_{i0})^2}{2\sigma_i^2} \right)$$

$$= \sum_i \frac{x_i^2 - 2x_i \mu_{i1} + \mu_{i1}^2 - x_i^2 + 2x_i \mu_{i0} + \mu_{i0}^2}{2\sigma_i^2}$$

(19)

$$= \sum_i \left(\frac{(\rho - 2x_i^T \mu_{ij}) + 2x_i^T \mu_{iv}}{2\sigma_j^2} + \frac{\mu_{ij}^2 + \mu_{iv}^2}{2\sigma_j^2} \right)$$

$$= \sum_i \left(\frac{(\mu_{iv} - \mu_{ij})x_i}{\sigma_j^2} + \frac{\mu_{ij}^2 + \mu_{iv}^2}{2\sigma_j^2} \right)$$

(4)

Substituting (1) in (2)

we get

$$\frac{1}{1 + \exp \left[\ln \left(\frac{1-\theta}{\theta} \right) + \sum_i \left(\frac{\mu_{iv} - \mu_{ij}}{2\sigma_j^2} x_i + \frac{\mu_{ij}^2 + \mu_{iv}^2}{2\sigma_j^2} \right) \right]}$$

This is similar to logistic regression
function $\frac{1}{1 + e^{-z}}$

where $z = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d$
for d features

: Comparing values we can see
that

w_0 = constant with $\ln \theta$ term

$$\therefore w_0 = \ln \left(\frac{1-\theta}{\theta} \right) + \sum_i \left(\frac{\mu_{ij}^2 + \mu_{iv}^2}{2\sigma_j^2} \right)$$

$$L \sum_i^d w_i x_i \text{ from LK}$$

can be equated with

$$i = \frac{(m_{i0} - m_i)}{2\sigma_i} x_i$$

where

$$w_i = \frac{m_{i0} - m_i}{2\sigma_i}$$

∴ we can say that LR & GNB have same form with assumption that variance doesn't change among classes for given feature.

Now, we have to make this assumption because in order to convert AND form to LR form, GNB has to be a linear classifier since LR is a linear classifier

$$\text{ie. for } \frac{1}{1+e^{-z}} \quad z = w_0 x_0 + w_1 x_1 + w_d x_d -$$

we have all x_i^0 as first degree terms

∴ assuming $\sigma_k^2 = \sigma^2$ helps eliminate x^2 terms, making GNB linear.



① Pseudo code for Logistic Regression

Load data into memory array by
Split $x \in \mathbb{R}^4$ [features & labels] .

Divide the data into 3 parts
for repeat fell 3 times, one for each
Subset of data:

Take $\frac{2}{3}$ of data as train data
& remaining $\frac{1}{3}$ as test data

Initialise weight vector

Add a 0^{th} column to x

[so that we can do $x \cdot w$
ie $\sum_{i=0}^n w_i x_i$]

every

For iterations, usually 10,000:

do following:

$$z = x \cdot w$$

$h = \text{sigmoid}(z)$ i.e.

$$\left[h = \frac{1}{1 + e^{-z}} \right]$$

Calculate gradient
using formula

$$\text{gradient}_i = \sum_j x_j^i (y - h)$$

$$\text{gradient}_i = x_i \cdot (y - h)$$

Then update weights

$$\text{or } w_i = w_i + n \cdot \text{gradient}_i$$

$$w = w + n \cdot \text{gradient}$$

After all iterations, we'll get final weight vector which we will use to predict remaining $\frac{1}{3}$ (test) data.

// Predict & calculate loss.

Let x_{test} be test data

$z_k = x_{\text{test}} \cdot \text{final weight } (w)$

$h = \text{sigmoid } (z)$

This model give predicted value.

Find loss using fell function

$$-y \log h - (1-y) \log (1-h)$$

Average the loss for all different test data sets used (3 here)

* Pseudo code for Naive Bayes

Load data

Split train & test data

Separate

Split train values into the respective classes.

Now, for every class, for every feature calculate mean & standard deviation

for i in class

for j in features

$$m_{ij} = \frac{1}{m} \sum x_{ij}$$

$$\sigma_{ij}^2 = \frac{1}{m} \sum (x_{ij} - m_{ij})^2$$

where m is the no. of rows in training data

This is the training part,

Now, we have mean & std for all classes & ith features, we can use this to predict new values.

for feature in test row

$$p(x_i | y=k) = \frac{e^{-\frac{(x_i - m_{ik})^2}{2\sigma_{ik}^2}}}{\sqrt{2\pi}\sigma}$$

for i^{th} feature $\in k^{th}$ class

So

$$p(x_i | y=y_0) = \frac{1}{\sigma} \frac{\exp(-\frac{(x_i - m_{i0})^2}{2\sigma_{i0}^2})}{\sqrt{2\pi}}$$

$$p(x_i | y=y_1) = \frac{1}{\sigma} \frac{\exp(-\frac{(x_i - m_{i1})^2}{2\sigma_{i1}^2})}{\sqrt{2\pi}}$$

$$p(y|x) = \prod_{i=1}^4 p(x_i | y \neq 0)$$

$$i.e. \text{ we } p(y=1|x) = \prod_{i=1}^4 p(x_i | y=1)$$

if $P(Y=1|X) > P(Y=0|X)$
 we predict 1
 else we predict 0.

- ③ We can see from the O/P of Q5(c) code that if we use trained model which has μ_i & σ_i^2 (i^{th} feature & k^{th} class) we can generate n samples using these parameters ($n \approx 50$). After generating these samples, we saw that mean and standard deviation of these samples is similar to the mean & std of trained model which seems obvious. We generated the samples from the same model. The samples will be within 3 std of the model.

(Q4)

$$\begin{aligned} n_1 &= (1, 0, 0) & n_2 &= (0, 0, 1) & n_3 &= (0, 1, 0) \\ y_1 &= 1 & y_2 &= 0, 1 & y_3 &= 1 \\ &&&&&\rightarrow +ve samples \end{aligned}$$

$$\begin{aligned} n_4 &= (-1, 0, 0) & n_5 &= (0, -1, 0) & n_6 &= (0, 0, -1) \\ y_4 &= 0 & y_5 &= 0 & y_6 &= 0 \\ &&&&&\rightarrow -ve samples \end{aligned}$$

$$w_0 = (0, 0, 0, 0)$$

- 1st Case (i)

$$P(Y=1 | x^1, w^0) = \frac{e^0}{1+e^0} = 0.5$$

Here, since $w^0 = (0, 0, 0, 0)$
we get $w_0 + \sum w_i n_i = 0$

$$\begin{aligned} P(Y=1 | x^2, w^0) &= P(Y=1 | x^3, w^0) \\ P(Y=1 | x^4, w^0) &= P(Y=1 | x^5, w^0) = P(Y=1 | x^6, w^0) \\ &= \underline{0.5} \end{aligned} \quad (1)$$

Now,

$$\begin{aligned} w'_0 &= w_0 + \sum_j n_j (y_j - P(Y=1 | x^j, w_0)) \\ &= 0 + n [1 - 0.5 + \cancel{0.5} + 0 \\ &\quad \cancel{0.5} - \cancel{0.5} + \cancel{0.5}] \\ &= 0 + 0 \\ \therefore w'_0 &= 0 \end{aligned}$$

(26)

$$\omega_1' = \omega_1^0 + n \sum_j^n x_1^j (y_j^0 - p(y=1|x_1^j, \omega_1^0))$$

$$= 0 + n [1(1-0.5) + 0 + 0 + \\ -1(0-0.5)]$$

$$= 0 + n [0.5 + 0.5]$$

$$\therefore \omega_1' = n$$

$$\omega_2' = \omega_2^0 + n \sum_j^n x_2^j (y_j^0 - p(y=1|x_2^j, \omega_2^0))$$

$$= 0 + n [0 + 0 + 1(1-0.5) + \\ 0 + -1(0-0.5)]$$

$$= n [0.5 + 0.5]$$

$$= n$$

$$\omega_3' = \omega_3^0 + n \sum_j^n x_3^j (y_j^0 - p(y=1|x_3^j, \omega_3^0))$$

$$= 0 + n [(1(1-0.5)) + 0 + 0 + \\ 0 + 0 - 1(0-0.5)]$$

$$= 0 + n (0.5 + 0.5)$$

$$= n$$

$$\therefore \omega' = (0, n, n, n) \quad \dots \quad (27)$$

Now, we calculate p_{wb} .

$$p(y=1|x_i', \omega') = \frac{e^{w_0 + \sum_i w_i n_i}}{1 + e^{w_0 + \sum_i w_i n_i}}$$

(27)

$$e^{w_0 + \sum w_i x_i} = \frac{e^{\sum w_i x_i}}{1 + e^{\sum w_i x_i}} = \frac{e^{\sum w_i x_i}}{1 + e^{\sum w_i x_i}} = \frac{e^{\sum w_i x_i}}{1 + e^{\sum w_i x_i}}$$

$$\therefore P(Y=1 | x^1, w) = \frac{e^{\sum w_i x_i}}{1 + e^{\sum w_i x_i}}$$

Similarly for other positive samples we get

$$\begin{aligned} P(Y=1 | x^2, w) &= P(Y=1 | x^3, w) = \\ P(Y=1 | x^4, w) &= \frac{e^{\sum w_i x_i}}{1 + e^{\sum w_i x_i}} \quad \text{--- (3)} \end{aligned}$$

for negative samples

$$P(Y=0 | x^1, w) = \frac{e^{\sum w_i x_i}}{1 + e^{\sum w_i x_i}} = \frac{e^{-\sum w_i x_i}}{1 + e^{-\sum w_i x_i}} = \frac{e^{-\sum w_i x_i}}{1 + e^{-\sum w_i x_i}}$$

Similarly,

$$\begin{aligned} P(Y=0 | x^2, w) &= P(Y=0 | x^3, w) = P(Y=0 | x^4, w) \\ &= \frac{e^{-\sum w_i x_i}}{1 + e^{-\sum w_i x_i}} \quad \text{--- (4)} \end{aligned}$$

Now, we will calculate w^2 using Gradient Ascent

$$\begin{aligned} w_0^2 &= w_0 + n \sum_{j=1}^n (y_j - P(Y=1 | x^j, w)) \\ &= 0 + n \left[3 \left(1 - \frac{e^{-\sum w_i x_i}}{1 + e^{-\sum w_i x_i}} \right) + 10 + 0 + \right. \\ &\quad \left. 10 \phi \left(0 = \frac{e^{-\sum w_i x_i}}{1 + e^{-\sum w_i x_i}} \right) \right] \end{aligned}$$

Since $\phi'(0) = 0$

$$\therefore \phi \left(0 = \frac{e^{-\sum w_i x_i}}{1 + e^{-\sum w_i x_i}} \right)$$

$$= 0 + n \left[\frac{1}{1+e^n} + \frac{e^{-n}}{1+e^{-n}} \right]$$

$$= 0 + \frac{2n}{1+e^n} \quad \textcircled{5} \quad \textcircled{0}$$

$$\omega_1^2 = \omega_1 + n \sum_j x_j^j (p_j^j - p(4=1) x_j, \omega^*)$$

$$= n + n \left[1 \left(1 - \frac{e^n}{1+e^n} \right) - 1 \left(0 - \frac{e^{-n}}{1+e^{-n}} \right) \right]$$

$$n^2 + \frac{2n}{1+e^n} = n \left(1 + \frac{e^n}{1+e^n} \right) n \left(1 + \frac{2}{1+e^n} \right)$$

Similarly $\omega_1^3 = \frac{(3+e^n)}{1+e^n} n$

Similarly $\omega_2^2 = \omega_3^2 = \frac{(3+e^n)}{1+e^n} n$

$$\therefore \omega^2 = \left(\frac{2n}{1+e^n}, \frac{(3+e^n)}{1+e^n} n, \frac{(3+e^n)}{1+e^n} n, \frac{(3+e^n)}{1+e^n} n \right)$$

From $\textcircled{0}, \textcircled{1}, \textcircled{2}, \textcircled{3}$

we see

$$\omega_0^0 = (0, 0, 0, 0)$$

$$\omega_1^1 = (0, n, n, n)$$

$$\omega_2^2 = \left(\frac{2n}{1+e^n}, \frac{(3+e^n)}{1+e^n} n, \frac{(3+e^n)}{1+e^n} n, \frac{(3+e^n)}{1+e^n} n \right)$$

where we can say that ω_0^0 remain 0
 and $\omega_0^0 < \omega_1^1 < \omega_2^2$ $\forall i \text{ bet } 0 \text{ to } 3$
 which means ω is increasing constantly

$$x_1 = (1, 0, 0) \quad x_2 = (0, 1, 0) \quad x_3 = (0, 0, 1)$$

(2a)

Now

$$\text{for } w = (0, 0, 1, 0) \rightarrow 2^{\text{nd}} \text{ case}$$

$$w^0 = (0, 0, 1, 0)$$

$$\begin{aligned} P(y=1 | x, w^0) &= \frac{e^{w^0 \cdot x}}{1 + e^{w^0 \cdot x}} \\ &= \frac{e^{0 + [w_2 + 2]}}{1 + e^{0 + w_2 + 2}} \\ &= \frac{e^0}{1 + e^0} = 0.5 \end{aligned}$$

$$P(y=1 | x^2, w^0) = \frac{e^{w^0 \cdot x^2}}{1 + e^{w^0 \cdot x^2}} = \frac{e^0}{1 + e^0} = 0.5$$

$$P(y=1 | x^3, w^0) = \frac{e^0}{1 + e^0} = 0.5$$

$$P(y=1 | x^4, w^0) = \frac{e^0}{1 + e^0} = 0.5$$

$$P(y=1 | x^5, w^0) = \frac{e^{0-1}}{1 + e^{-1}} = 0.27$$

$$P(y=1 | x^6, w^0) = \frac{e^0}{1 + e^0} = 0.5$$

Now, we calculate w_j^1 using Gradient Descent

$$w_j^1 = w_j^0 + n \sum_i x_j^i (y_i^j - P(y=1 | x^j, w^0))$$

$$= 0 + n[0] \Rightarrow 0$$

$$w_j^1 = w_j^0 + n \sum_i x_j^i (y_i^j - P(y=1 | x^j, w^0))$$

$$\begin{aligned}
 &= 0 + n(1(1-0.5) - 1(0-0.5)) \\
 &= n(0.5 + 0.5) \\
 &= n \\
 \therefore w_1 &= n
 \end{aligned}$$

$$\begin{aligned}
 w_2' &= w_0^0 + n \left[1(1 - \frac{e}{1+e}) - 1 \left(0 - \frac{e^{-1}}{1+e^{-1}} \right) \right] \\
 &= 1 + n \left[\frac{1}{1+e} + \frac{1}{1+e} \right] \\
 &= 1 + \frac{2n}{1+e} \\
 \therefore w_2 &= 1 + \frac{2n}{1+e}
 \end{aligned}$$

$$\begin{aligned}
 w_3' &= w_0^0 + n \left[1(1-0.5) - 1(0-0.5) \right] \\
 &= 0 + n \\
 &= n
 \end{aligned}$$

$$\therefore w = (0, n, 1 + \frac{2n}{1+e}, n)$$

Now we calculate Probabilities.

$$P(Y=1 | x^1, w') = \frac{e^n}{1+e^n}$$

$$P(Y=1 | x^2, w') = \frac{e^{1+\frac{2n}{1+e}}}{1+e^{1+\frac{2n}{1+e}}}$$

$$P(Y=1 | x^3, w') = \frac{e^n}{1+e^n}$$

$$P(Y=1|X^4) = P(Y=1|X^6) = \frac{e^{-n}}{1+e^{-n}}$$

$$P(Y=1|X^5) = \frac{e^{-1-\frac{2n}{1+e}}}{1+e^{-1-\frac{2n}{1+e}}}$$

Now, we calculate \bar{w}^2 using G.A.

$$\bar{w}_0 = \bar{w}_0 + n[0] \\ = 0$$

$$\begin{aligned}\bar{w}_1 &= \bar{w}_0 + n \left[1 \left(1 - \frac{e^{-n}}{1+e^{-n}} \right) - 0 - \frac{e^{-n}}{1+e^{-n}} \right] \\ &= n + n \left[\frac{1}{1+e^{-n}} + \frac{1}{1+e^{-n}} \right] \\ &= n \frac{2n}{1+e^{-n}} - n \frac{3+e^{-n}}{1+e^{-n}} \\ &= n \left(\frac{3+e^{-n}}{1+e^{-n}} \right)\end{aligned}$$

$$\begin{aligned}\bar{w}_2 &= \bar{w}_1 + n \left[1 \left(1 - \frac{e^{-1-\frac{2n}{1+e}}}{1+e^{-1-\frac{2n}{1+e}}} \right) \right. \\ &\quad \left. - 0 - \frac{e^{-1-\frac{2n}{1+e}}}{1+e^{-1-\frac{2n}{1+e}}} \right] \\ &= n \frac{2n}{1+e^{-1-\frac{2n}{1+e}}} + n \left[\frac{1}{1+e^{-2}} + \frac{1}{1+e^{-2}} \right]\end{aligned}$$

$$= 1 + \frac{2n}{1+e} + \frac{2}{1+e^2} \quad z = 1 + \frac{2n}{1+e}$$

$$\therefore \bar{w}_2 = 1 + \frac{2n}{1+e} + \frac{2}{1+e^{1+\frac{2n}{1+e}}}$$

(32)

$$w_3^2 = w_1^2$$

$$= n + \frac{2n}{1+e^n} = \left(\frac{3+e^n}{1+e^n} \right) n$$

∴ we have

$$w^2 = \left(0, \left(\frac{3+e^n}{1+e^n} \right) n, 1 + \frac{2n}{1+e^n} + \frac{2n}{1+e^{\frac{2n}{1+e^n}}} \right)$$

$$w^0 = (0, 0, 1, 0)$$

$$w^1 = (0, n, 1 + \frac{2n}{1+e}, n)$$

$$w^2 = (0, \left(\frac{3+e^n}{1+e^n} \right) n, 1 + \frac{2n}{1+e} + \frac{2n}{1+e^{1+\frac{2n}{1+e}}})$$

$$, \left(\frac{3+e^n}{1+e^n} \right) n)$$

∴ we see that w^2 is same for all

values

$$w_1^1 = w_3^1 \text{ for all values } n$$

$$w_1^0 < w_1^1 < w_1^2 \text{ also}$$

$$w_2^0 < w_2^1 < w_2^2 \text{ also}$$

$$w_3^0 < w_3^1 < w_3^2$$

(33)

which means w's are constantly increasing to attain a final value

$$w^k = \begin{bmatrix} 0 \\ \infty \\ \infty \\ \infty \\ \infty \end{bmatrix}$$

This is similar to the first case as well.

i.e. weight increase till they attain a certain max value.

I think for initial iterations final weight vector will be different for ~~four~~ ~~with~~ different weight initializations but as we increase iterations they'll approach a max value where the difference will be negligible.

At I added, I observed that for 10^6 iterations, the difference in the 2nd value i.e. w_2 for the final vectors was nearly 0.01 only. (O/P attached)

We can also see that, initially the difference =
 $w_0^0 = (0, 0, 0, 0)$ & $w_0^1 = (0, 0, 1, 0)$

The difference was $(0, 0, 1, 0)$

with one iteration

$$\begin{aligned} \left| w_i^0 - w_i' \right| &= (0, n, n, n) - \left(0, n, 1 + \frac{2n}{1+e}, n \right) \\ (\text{1st Iterat}) &= (0, 0, \approx 0.99, 0) \end{aligned}$$

Now the difference is 0.99 which is less than 1

$$\begin{aligned} \left| w_2^0 - w_2' \right| &= \left[0, \frac{(3+e^n)}{1+e^n} n, \frac{(3+e^n)}{1+e^n} n, \frac{(3+e^n)}{1+e^n} n \right] \\ (\text{2nd Iterat}) &- \left[0, \frac{(3+e^n)}{1+e^n} n, 1 + \frac{2n}{1+e} + \frac{2n}{1+e(1+2n)}, \frac{(3+e^n)}{1+e^n} n \right] \\ &\quad \left(\frac{3+e^n}{1+e^n} n \right) \end{aligned}$$

Now the difference is

$$1 + \frac{2n}{1+e} + \frac{2n}{1+e(1+2n)} - \frac{(3+e^n)}{1+e^n} n$$

$$\approx 0.27 - 0.019$$

$$\approx 0.25$$

which is less than 0.99

for n ,

(3)

The difference will keep on decreasing with fractions until a point that it won't matter much, ie. both ratios will almost be similar