

Q1 Given  $P_{\mu}[y_i = c] = \pi_c$

$$x_i | y_i \sim N(\mu_c, \Sigma_c).$$

Gaussian Naive Bayes classifier is similar to GDA. The difference is in the construction of  $P(x_i | y_i)$ . NB makes a (naive but strong) assumption that features are independent so that

$$\begin{aligned} P(x | y = k) &= P(x_1 | y = k) \cdot P(x_2 | y = k) \cdot P(x_3 | y = k) \cdots P(x_n | y = k) \\ &= \prod_{j=1}^n P(x_j | y = k) \end{aligned}$$

For multivariate normal distribution,  $P(x_j | y = k) \sim N(\mu_j^k, \sigma_j^2)$

We define probability for any class given  $x$ .

$$P(y | x) = P(y) \prod_{j=1}^n P(x_j | y)$$

a) maximum likelihood estimation for the model parameters.

$$\text{L} = - \sum_{i=1}^n \log P(y^{(i)}) \prod_{j=1}^n P(x_j^{(i)} | y^{(i)}).$$

$$= - \left( \sum_{i=1}^n \log P(y^{(i)} = c) + \sum_{i=1}^n \log \left( \prod_{j=1}^d P(x_j^{(i)} | y^{(i)}) \right) \right) \quad \text{--- (1)}$$

If we constrain  ~~$\Sigma_c$  to be diagonal~~ then  $\Sigma_c$  (where  $c = 1, \dots, k$ ) to be diagonal then we can rewrite  $P(x, y)$  as a product of  $P(x_j | y_i)$  for  $j = 1, \dots, d$  features.

$$p(\mathbf{x}/y^{(i)}) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma_{y^{(i)}})}} \exp \left( -\frac{1}{2} (\mathbf{x}_j - \mu_{j y^{(i)}})^T \Sigma_{y^{(i)}}^{-1} (\mathbf{x}_j - \mu_{j y^{(i)}}) \right)$$

$$= \prod_{j=1}^d \frac{1}{\sqrt{(2\pi) \Sigma_{j y^{(i)}}}} \exp \left( -\frac{1}{2 \Sigma_{j y^{(i)}}} \|\mathbf{x}_j - \mu_{j y^{(i)}}\|_2^2 \right)$$

$$= \prod_{j=1}^d p(\mathbf{x}_j | t)$$

Using this fact we can modify our equation (1) as.

$$= - \left( \sum_{i=1}^n \log \pi_{y^{(i)}} + \sum_{i=1}^n \sum_{j=1}^d \log \frac{1}{\sqrt{2\pi \Sigma_{j y^{(i)}}}} \exp \left( -\frac{(\mathbf{x}_j - \mu_{j y^{(i)}})^2}{2 \Sigma_{j y^{(i)}}} \right) \right)$$

$$= - \left( \sum_{i=1}^n \log \pi_{y^{(i)}} + \sum_{i=1}^n \sum_{j=1}^d \log \frac{1}{\sqrt{2\pi \Sigma_{j y^{(i)}}}} + \sum_{i=1}^n \sum_{j=1}^d \log \exp \left( -\frac{(\mathbf{x}_j - \mu_{j y^{(i)}})^2}{2 \Sigma_{j y^{(i)}}} \right) \right)$$

$$= - \left( \sum_{i=1}^n \log \pi_{y^{(i)}} + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^d \log (2\pi \Sigma_{j y^{(i)}}) - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^d \log \exp \left( \frac{(\mathbf{x}_j - \mu_{j y^{(i)}})^2}{2 \Sigma_{j y^{(i)}}} \right) \right)$$

differentiate with respect to  $\mu_{j y^{(i)}}$

$$\mu_{j y^{(i)}} = \frac{\sum_{t=1}^n 1[y^{(t)} = y^{(i)}] \mathbf{x}_j^{(t)}}{\sum_{t=1}^n 1[y^{(t)} = y^{(i)}]}$$

taking derivative with respect to  $u_j(y^{(i)}=c)$ .

$$u_{jc} = \frac{\sum_{i=1}^n 1[y^{(i)}=c] x_j^{(i)}}{\sum_{i=1}^n 1[y^{(i)}=c]}$$

$\swarrow$   $j$ th feature.       $\searrow$   $c$ th class.  
 $\rightarrow$   $i$ th training example.       $\rightarrow$  1 if  $y^{(i)}=c$  or 0.

taking derivative with respect to  $\nabla_{j,c}^2$   $i$ th training example.

$$\nabla_{j,c}^2 = \frac{\sum_{i=1}^n 1[y^{(i)}=c] (x_j^{(i)} - u_{jc})^2}{\sum_{i=1}^n 1[y^{(i)}=c]}$$

$\swarrow$   $j$ th feature       $\searrow$   $c$ th class.

taking derivative with respect to  $\pi_c$ .

$$\pi_c = \frac{\sum_{i=1}^n 1(y^{(i)}=c)}{n}$$

(b.) Now with estimated model parameters  $\pi_c, u_c, \nabla_c$ , given any new data point we can predict  $\hat{y}_0$

in following way.

$$\hat{y}_0 = \underset{y_c, c=1, \dots, K}{\operatorname{argmax}} P(y_c) P(x_0/y_c).$$

$$P(x/y_c) = \prod_{i=1}^d P(x_i/y_c) \quad \because \text{we know } \pi_c, \mu_{ic}, \sigma_{ic}^2.$$

$$= \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_{ic}^2}} \exp \left[ \frac{-(x_i - \mu_{ic})^2}{2\sigma_{ic}^2} \right]$$

$$\underset{y_c}{\operatorname{argmax}} \pi_c \cdot \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_{ic}^2}} \exp \left[ \frac{-(x_i - \mu_{ic})^2}{2\sigma_{ic}^2} \right] \rightarrow \text{new value of } x_0 \text{ for which prediction has to be found.}$$

plug the estimated values in the equation to get predicted  $\hat{y}_0$



Q2

Given:- Lasso regularized K-class logistic regression problem.

$$\underset{\theta_1, \dots, \theta_K}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \left( \log \sum_{c=1}^K \exp(\theta_c^T \phi_i) - \theta_{y_i}^T \phi_i \right) + \frac{\lambda}{2} \sum_{c=1}^K \|\theta_c\|_1$$

This is of the form  $f(x) = g(x) + h(x)$  so using proximal stochastic gradient descent algorithm to solve this will be our approach of the algorithm.

Pseudo code for solving this can be as follows:-

→ Choose or initialize  $\theta^{(0)}$  and repeat following steps for  $k=1, 2, 3, \dots$  max-iteration.

→ For a random sample  $i$  uniformly drawn from  $\{1, \dots, n\}$ .  
The stochastic step.

→ Calculate the gradient  $\nabla g$  for  $i$ .

if  $i = y_i$ ; input class is equal to  $y_i$

$$\nabla g(\theta^{(k-1)}) = \frac{\phi_i e^{\theta_i^T \phi_i}}{\sum_{c=1}^K e^{\theta_c^T \phi_i}} - \phi_i \quad \left\{ \begin{array}{l} \text{using } f(z) = \log \sum_{c=1}^K e^{z_c} \\ \nabla f(z) = \frac{1}{\sum_{c=1}^K e^{z_c}} \exp(z) \end{array} \right.$$

$$\text{else } \nabla g(\theta^{(k-1)}) = \frac{\phi_i e^{\theta_i^T \phi_i}}{\sum_{c=1}^K e^{\theta_c^T \phi_i}}$$

→ update  $\theta^{(k)}$

$$\theta^{(k)} = \theta^{(k-1)} - t_k \nabla g(\theta^{(k-1)})$$

where  $t_k$  is the  
diminized step size  
of our algorithm.

→ feed the new  $\theta^{(k)}$  into proximal operation

where  $\text{prox}_{\frac{\lambda}{2}}(\tilde{\theta}) = \arg\min_{\theta} \frac{1}{2} \|\theta - \tilde{\theta}\|^2 + \lambda h(\theta)$  → lasso regularization in our case.

prox. operator tries to find a point that makes  $h(\theta)$  small but is approximately close to  $\tilde{\theta}$ .

For lasso regularization, proximal operator can be explicitly expressed as soft thresholding.

$$\text{prox}_{\frac{\lambda}{2}}(\tilde{\theta})_k = \begin{cases} \tilde{\theta}^{(k)} - \frac{\lambda}{2} & \tilde{\theta}^{(k)} > \frac{\lambda}{2} \\ 0 & |\tilde{\theta}^{(k)}| \leq \frac{\lambda}{2} \\ \tilde{\theta}^{(k)} + \frac{\lambda}{2} & \tilde{\theta}^{(k)} < -\frac{\lambda}{2} \end{cases}$$

→ repeat for given number of iteration.

→ End of pseudo code.

Q3

Given

$$p_{\pi}(y_i = c) = \pi_c \quad p_{\pi}(x_i / y_i = c) \sim \text{Multi}(p_c, L_i)$$

$$p(x_i / y_i = c) = \frac{L_i!}{\prod_{j=1}^d x_{ij}!} \prod_{j=1}^d p_{c,j}^{x_{ij}}$$

(a) For maximum likelihood formulation for estimating  $p_1, p_2, \dots, p_K$  and  $\pi$  we find the log likelihood of probability mass function.

$$\Rightarrow \sum_{i=1}^n \log p(x_i, y_i; \pi, p_c) \quad \text{where } c = 1, 2, \dots, K.$$

$$\Rightarrow \sum_{i=1}^n \log p(y_i = c) + \sum_{i=1}^n \log p(x_i / y_i; p_c).$$

$$\Rightarrow \sum_{i=1}^n \log p(y_i = c) + \sum_{i=1}^n \sum_{j=1}^d x_{ij} \log p_{cj} + \underbrace{\text{constant}}_{\text{parameter-free term.}}$$

$$\Rightarrow \sum_{i=1}^n \sum_{c=1}^K \mathbb{I}\{y_i = c\} \log \pi_c + \sum_{i=1}^n \sum_{c=1}^K \sum_{j=1}^d x_{ij} \mathbb{I}\{y_i = c\} \log p_{cj} + \text{constant}$$

where this is either ① or ② ①

(b) The E-step on expectation step requires that we compute the expected complete log likelihood under multinomial distribution.

$$E p(\cdot | x_i, \pi, p_c) \left[ \sum_{i=1}^n \log p(x_i, y_i; \pi, p_c) \right],$$

$$= \sum_{i=1}^n \sum_{c=1}^K \phi_{i,c} \log \pi_c + \sum_{i=1}^n \sum_{c=1}^K \sum_{j=1}^d x_{ij} \phi_{i,c} \log p_{cj}$$

where we define  $\phi_{i,c} = p(y_i = c | x_i; \pi, p_c)$  as in the GMM case. studies in class, it suffices to compute  $\phi_{i,c}$  to complete this step.

By Bayes' Rule.

$$p(y_i = c | x_i; \pi, p_c) \propto p(y_i = c; \pi) p(x_i / y_i; p_c)$$

$$P(y_i = c | x_i; \pi, p_c) = \frac{\pi_c f(x_i, p_c)}{\sum_{\ell=1}^K \pi_\ell f(x_i, p_\ell)}$$

where  $f(x_i, p_c) \propto \prod_{j=1}^d p_{cj}^{x_{ij}}$  using multinomial distribution.

The M-step maximize the expected log likelihood we have a closed form solution for the parameter updates.

$$\pi_c = \frac{1}{n} \sum_{i=1}^n \phi_{i,c} \quad \text{taking derivative of equ ① function w.r.t to } \pi_c$$

$$p_{cj} = \frac{\sum_{i=1}^n x_{ij} \phi_{i,c}}{\sum_{i=1}^n \phi_{i,c}}$$

(c)

Top 10 key words for cluster: 0

['ooo\n', 'colour\n', 'moslems\n', 'italian\n', 'subjects\n', 'atom\n', 'explaining\n', 'league\n', 'crucifixion\n', 'honor\n']

Top 10 key words for cluster: 1

['koran\n', 'innermost\n', 'geoffrey\n', 'wesleyan\n', 'nrp\n', 'league\n', 'coincidence\n', 'disobeying\n', 'creatures\n', 'exist\n']

Top 10 key words for cluster: 2

['koran\n', 'innermost\n', 'disobeying\n', 'nrp\n', 'creatures\n', 'exist\n', 'coincidence\n', 'necessity\n', 'shits\n', 'fri\n']

Top 10 key words for cluster: 3

['koran\n', 'innermost\n', 'disobeying\n', 'nrp\n', 'geoffrey\n', 'creatures\n', 'coincidence\n', 'exist\n', 'wesleyan\n', 'shits\n']

Top 10 key words for cluster: 4



['ooo\n', 'colour\n', 'moslems\n', 'subjects\n', 'league\n',  
'italian\n', 'geoffrey\n', 'atom\n', 'crucifixion\n', 'explaining\n']  
Top 10 key words for cluster: 5  
['moslems\n', 'explaining\n', 'italian\n', 'colour\n', 'ooo\n',  
'gaul\n', 'atom\n', 'climbed\n', 'honor\n', 'wherein\n']  
Top 10 key words for cluster: 6  
['koran\n', 'geoffrey\n', 'league\n', 'wesleyan\n', 'blood\n',  
'husbands\n', 'ooo\n', 'subjects\n', 'crucifixion\n', 'innermost\n']  
Top 10 key words for cluster: 7  
['disobeying\n', 'innermost\n', 'koran\n', 'creatures\n', 'exist\n',  
'nrp\n', 'coincidence\n', 'necessity\n', 'shits\n', 'fri\n']  
Top 10 key words for cluster: 8  
['disobeying\n', 'innermost\n', 'koran\n', 'creatures\n', 'nrp\n',  
'exist\n', 'coincidence\n', 'necessity\n', 'shits\n', 'fri\n']  
Top 10 key words for cluster: 9  
['disobeying\n', 'koran\n', 'innermost\n', 'creatures\n', 'nrp\n',  
'exist\n', 'coincidence\n', 'necessity\n', 'shits\n', 'fri\n']  
Top 10 key words for cluster: 10  
['koran\n', 'geoffrey\n', 'league\n', 'wesleyan\n', 'innermost\n',  
'blood\n', 'husbands\n', 'ooo\n', 'subjects\n', 'omnipotent\n']  
Top 10 key words for cluster: 11  
['ooo\n', 'moslems\n', 'colour\n', 'italian\n', 'explaining\n',  
'atom\n', 'subjects\n', 'gaul\n', 'honor\n', 'climbed\n']  
Top 10 key words for cluster: 12  
['disobeying\n', 'innermost\n', 'koran\n', 'creatures\n', 'nrp\n',  
'exist\n', 'coincidence\n', 'necessity\n', 'shits\n', 'fri\n']  
Top 10 key words for cluster: 13  
['koran\n', 'innermost\n', 'disobeying\n', 'nrp\n', 'creatures\n',  
'coincidence\n', 'exist\n', 'shits\n', 'necessity\n', 'fri\n']  
Top 10 key words for cluster: 14  
['moslems\n', 'explaining\n', 'colour\n', 'italian\n', 'ooo\n',  
'gaul\n', 'atom\n', 'climbed\n', 'honor\n', 'subjects\n']  
Top 10 key words for cluster: 15  
['moslems\n', 'ooo\n', 'colour\n', 'italian\n', 'explaining\n',  
'atom\n', 'gaul\n', 'subjects\n', 'climbed\n', 'honor\n']  
Top 10 key words for cluster: 16  
['ooo\n', 'league\n', 'colour\n', 'moslems\n', 'subjects\n',  
'geoffrey\n', 'italian\n', 'crucifixion\n', 'blood\n', 'atom\n']  
Top 10 key words for cluster: 17  
['explaining\n', 'moslems\n', 'italian\n', 'colour\n', 'ooo\n',  
'gaul\n', 'atom\n', 'climbed\n', 'honor\n', 'fondly\n']  
Top 10 key words for cluster: 18  
['ooo\n', 'moslems\n', 'colour\n', 'italian\n', 'explaining\n',  
'subjects\n', 'atom\n', 'gaul\n', 'honor\n', 'climbed\n']  
Top 10 key words for cluster: 19  
['koran\n', 'innermost\n', 'disobeying\n', 'nrp\n', 'creatures\n',  
'exist\n', 'coincidence\n', 'necessity\n', 'shits\n', 'fri\n']

Q4 ~~The given~~

(a) Given

$$\text{minimize } \sum_{i=1}^n \sum_{j=1}^n (\tilde{x}_i^T \tilde{x}_j - \tilde{y}_i^T \tilde{y}_j)^2$$

$$\Rightarrow \text{minimize } \|\tilde{X}^T \tilde{X} - \tilde{Y}^T \tilde{Y}\| \text{ --- (1)}$$

Using Eckart-Young theorem.

$$\|A - B\| \geq \|A - A_k\| \text{ where } \text{rank}(B) \leq k.$$

and  $A_k$  is the best approximation of  $A$  to rank  $k$ ,  
 $k$  as taking largest  $k$  eigenvalues in SVD.

Using this theorem we can state that.

$$\|\tilde{X}^T \tilde{X} - \tilde{Y}^T \tilde{Y}\| \geq \|\tilde{X}^T \tilde{X} - (\tilde{X}^T \tilde{X})_k\|.$$

At optimality.

$$\tilde{Y}^T \tilde{Y} = (\tilde{X}^T \tilde{X})_k \text{ --- (2)}$$

Now

Take SVD and largest  $k$  eigenvalues.

$$\Rightarrow \tilde{X}^T \tilde{X} = U \Lambda U^T$$

$$\Rightarrow (\tilde{X}^T \tilde{X})_k = U_k \Lambda_k U_k^T \text{ --- (3)}$$

from (2) and (3)

$$\tilde{Y}^T \tilde{Y} = U_k \Lambda_k U_k^T$$

$$Y = \Lambda^{1/2} U^T \text{ as } Y \text{ is positive definite}$$

Since we only keep  $K$  eigenvalues, the corresponding columns in  $Y$  using  $Y = \Lambda^{1/2} U^T$  are the  $i^{\text{th}}$  column of  $Y$  which is  $y_i$  which is same as PCA.

(b)

$$\text{let } K = X^T X$$

$$\text{with } \Lambda = \text{diag}(\Lambda_{ii}) \in \mathbb{R}^n$$

$$D = d_{ij}^2 \text{ given in question becomes}$$

$$= K \cdot \mathbf{1}^T + \mathbf{1} \cdot K^T - 2K \quad \text{--- (1)}$$

Also given in question.

$$\tilde{x}_i = x_i - \frac{1}{n} \sum_{j=1}^n x_j \} \rightarrow \text{nothing but mean.}$$

$$\tilde{X} = X - \frac{1}{n} X \cdot \mathbf{1} \cdot \mathbf{1}^T \quad \text{--- (2)}$$

$$\text{Given } B = -\frac{1}{2} \left( \mathbf{1} - \frac{1}{n} \mathbf{1} \cdot \mathbf{1}^T \right) D \left( \mathbf{1} - \frac{1}{n} \mathbf{1} \cdot \mathbf{1}^T \right)^T$$

$$\text{we have to prove } B = \tilde{X}^T \tilde{X}$$

using (1)

$$B = -\frac{1}{2} \left( \mathbf{1} - \frac{1}{n} \mathbf{1} \cdot \mathbf{1}^T \right) (K \cdot \mathbf{1}^T + \mathbf{1} \cdot K^T - 2K) \left( \mathbf{1} - \frac{1}{n} \mathbf{1} \cdot \mathbf{1}^T \right)^T$$

expanding this we have.

$$B = HKH^T - \frac{1}{2} H \cdot \mathbf{1} \cdot \mathbf{1}^T \cdot H^T - \frac{1}{2} H \cdot \mathbf{1} \cdot \mathbf{1}^T \cdot H^T$$

where  $H = \left( I - \frac{1}{n} \cdot \mathbf{1} \cdot \mathbf{1}^T \right)$ .

$$\begin{aligned} H \cdot \mathbf{1} \cdot \mathbf{1}^T \cdot H^T &= H \cdot \mathbf{1} \cdot \mathbf{1}^T \cdot H^T = K \cdot \mathbf{1} \cdot \left( I - \frac{1}{n} \cdot \mathbf{1} \cdot \mathbf{1}^T \right)^T \\ &= K \cdot \mathbf{1} - K \left( \frac{\mathbf{1}^T \cdot \mathbf{1}}{n} \right) \cdot \mathbf{1} \\ &= 0. \end{aligned}$$

$$B = HKH^T - \frac{1}{2} H \cdot \mathbf{1} \cdot \mathbf{1}^T \cdot H^T - \frac{1}{2} H \cdot \mathbf{1} \cdot \mathbf{1}^T \cdot H^T$$

$$B = HKH^T = \left( I - \frac{1}{n} \cdot \mathbf{1} \cdot \mathbf{1}^T \right) K \left( I - \frac{1}{n} \cdot \mathbf{1} \cdot \mathbf{1}^T \right)^T$$

$$= K - \frac{1}{n} \cdot \mathbf{1} \cdot \mathbf{1}^T \cdot K - \frac{1}{n} \cdot K \cdot \mathbf{1} \cdot \mathbf{1}^T + \frac{1}{n^2} \cdot \mathbf{1} \cdot \mathbf{1}^T \cdot K \cdot \mathbf{1} \cdot \mathbf{1}^T$$

$$= K = X^T X$$

Hence proved.