

**Problem 1.**

For noise free observations:

$$\text{cov}(f_x) = k_{xx} - k_x^T k^{-1} k_x$$

As  $x^*$  is a training point. Let's assume  $x^* = x_i \quad i \in \{1, \dots, n\}$

$$\text{Then } k_x = k e_i, \quad k_{xx} = e_i^T k e_i$$

where  $e_i$  is a column vector with the  $i$ -th element be 1 and others be 0.

$$\begin{aligned}\therefore \text{cov}(f_x) &= e_i^T k e_i - (k e_i)^T k^{-1} (k e_i) \\ &= e_i^T k e_i - e_i^T k^T k^{-1} k e_i \\ &= e_i^T k e_i - e_i^T k^T e_i\end{aligned}$$

$\because K$  is the covariance matrix of  $X$ ,  $K$  is symmetric.  
 $\hookrightarrow$  (gram matrix if we use kernel)

$$\therefore k^T = K$$

$$\begin{aligned}\therefore \text{cov}(f_x) &= e_i^T k e_i - e_i^T k e_i \\ &= 0\end{aligned}$$

**Problem 2.**

$$a) p(\vec{y} | x, \vec{\theta}) = p(\vec{f} + \vec{\varepsilon} | x)$$

$$= N(0, K_\theta(x, x) + \sigma^2 I)$$

$$\begin{aligned}\therefore \log p(\vec{y} | x, \vec{\theta}) &= \log \left( \frac{1}{\sqrt{(2\pi)^n |K_\theta(x, x) + \sigma^2 I|}} e^{-\frac{1}{2} \vec{y}^T (K_\theta(x, x) + \sigma^2 I)^{-1} \vec{y}} \right) \\ &= -\frac{1}{2} \log [(2\pi)^n \cdot |K_\theta(x, x) + \sigma^2 I|] - \frac{1}{2} \vec{y}^T (K_\theta(x, x) + \sigma^2 I)^{-1} \vec{y} \\ &= -\frac{n}{2} \log 2\pi - \frac{1}{2} \log (|K_\theta(x, x) + \sigma^2 I|) - \frac{1}{2} \vec{y}^T (K_\theta(x, x) + \sigma^2 I)^{-1} \vec{y}\end{aligned}$$

b)  $\log(\vec{y} | x, \vec{\theta}) = -\frac{n}{2} \log 2\pi - \frac{1}{2} \log(|k_y|) - \frac{1}{2} \vec{y}^T k_y^{-1} \vec{y}$

$$= -\frac{1}{2} \log 2\pi - \frac{1}{2} \log(|LL^T|) - \frac{1}{2} \vec{y}^T (LL^T)^{-1} \vec{y}$$

$$= -\frac{1}{2} \log 2\pi - \frac{1}{2} \log |LL^T| - \frac{1}{2} \vec{y}^T L^{-1} L^T \vec{y}$$

$$\therefore \vec{\alpha} = L^T \backslash (L \backslash \vec{y})$$

$$\therefore L^T \vec{\alpha} = L \backslash \vec{y}$$

$$\therefore LL^T \vec{\alpha} = \vec{y}$$

$$\therefore \log(\vec{y} | x, \vec{\theta}) = -\frac{1}{2} \log 2\pi - \frac{1}{2} \log |LL^T| - \frac{1}{2} \vec{\alpha}^T LL^T (L^T)^{-1} L^{-1} LL^T \vec{\alpha}$$

$$= -\frac{1}{2} \log 2\pi - \frac{1}{2} \log |LL^T| - \frac{1}{2} \vec{\alpha}^T LL^T \vec{\alpha}$$

$$= -\frac{1}{2} \log 2\pi - \sum_{i=1}^n \log L_{ii} - \frac{1}{2} \vec{\alpha}^T LL^T \vec{\alpha}$$

c)

$$\frac{\partial \log(\vec{y} | x, \vec{\theta})}{\partial \vec{\theta}} = \frac{\partial \log(\vec{y} | x, \vec{\theta})}{\partial k_y} \cdot \frac{\partial k_y}{\partial \vec{\theta}}$$

$$= -\frac{n}{2} \log 2\pi - \frac{1}{2} \log(|k_y|) - \frac{1}{2} \vec{y}^T k_y^{-1} \vec{y}$$

$$= -\frac{1}{2} \frac{(k_y | k_y^{-1})}{|k_y|} \cdot \frac{\partial k_y}{\partial \vec{\theta}} - \frac{1}{2} \vec{y}^T k_y^{-1} \frac{\partial k_y}{\partial \vec{\theta}} k_y^{-1} \vec{y}$$

$$= -\frac{1}{2} k_y^{-1} \frac{\partial k_y}{\partial \vec{\theta}} - \frac{1}{2} \vec{y}^T k_y^{-1} \frac{\partial k_y}{\partial \vec{\theta}} k_y^{-1} \vec{y}$$

$$= -\frac{1}{2} (LL^T)^{-1} \frac{\partial k_y}{\partial \vec{\theta}} - \frac{1}{2} \vec{y}^T (LL^T)^{-1} \frac{\partial k_y}{\partial \vec{\theta}} (LL^T)^{-1} \vec{y}$$

$$= -\frac{1}{2} (L^{-1})^T L^{-1} \frac{\partial k_y}{\partial \vec{\theta}} - \frac{1}{2} \vec{\alpha}^T \frac{\partial k_y}{\partial \vec{\theta}} \vec{\alpha}$$

## Problem 3.

(i)  $\Leftrightarrow$  (ii)

prove (i)  $\Rightarrow$  (ii):

Assume in  $i$ -th iteration:

we got Assignment A.

Then we will update the center with:

$$\mu_\alpha = \frac{\sum_{i=1}^n [z_i]_\alpha x_i}{\sum_{i=1}^n [z_i]_\alpha}$$

AKA the mean vector of datas that assigned to the same cluster.

Then in  $(i+1)$ -th iteration:

if we still got Assignment A

Then  $\mu_\alpha$  will also be the same:  $\mu_\alpha =$

$$\frac{\sum_{i=1}^n [z_i]_\alpha x_i}{\sum_{i=1}^n [z_i]_\alpha}$$

$\therefore$  (i)  $\Rightarrow$  (ii)

prove (ii)  $\Rightarrow$  (i)

Assume in  $i$ -th iteration:

we got Assignment A

Then we update the center to be  $\mu_1, \dots, \mu_k$ .

Then in  $(i+1)$ -th iteration:

we got Assignment B (assign according to the closest center)

Then we update the center, but the center stay the same.

Then in the  $(i+1)$ th iteration.

As we still assign data according to the closest center to it, we will end up with the same assignment  $B$ .

$$\therefore (j)_i \Rightarrow (i)_i$$

In conclusion :

$$(i)_i \Leftrightarrow (i+1)_i$$

b) The k-means will always converge:

As there are only  $k^n$  different assignments on  $n$  data points.

As we want to minimize the "variance" inside each cluster. There must be a infimum bound of all this assignments. As we update the cluster each time to achieve lower "variance" inside each cluster.

And the infimum is finite. Then we must will converge.

c). It is possible that k-means generate empty cluster. for example when we initial the init-center into this:



However we should not allow it.

Because if k-means generate empty cluster, we can always fetch a data point from any other not-empty cluster and make it as this empty cluster center. Then we will achieve lower in cluster variance and achieve better result.

## Problem 4.

(Reference: 12 lecture note. Kernel k means)

a) Randomly initialize  $z_1, \dots, z_n$ 

repeat

for  $i = 1, \dots, n$  dofor  $\alpha = 1, \dots, k$  do

$$\text{d}i\alpha \leftarrow k(x_i, x_i) - \frac{2}{\sum_{j=1}^n [z_j]_\alpha} \sum_{j=1}^n [z_j]_\alpha K(x_i, x_j) + \frac{1}{(\sum_{j=1}^n [z_j]_\alpha)^2} \sum_{j=1}^n \sum_{l=1}^n [z_j]_\alpha [z_l]_\alpha K(x_j, x_l)$$

end for

end for

update  $z_i$  using  $\alpha = \operatorname{argmin}_\alpha \text{d}i\alpha$ if  $z_i$  do not change for all  $i$ , then

STOP

end if

until convergence.

- b) Because the cluster center "lies" in the transformed feature space. For some kernel like RBF, it has infinite dimensions. Thus is impossible to visualize or use them as prototypes for learning.

## Problem 6

a) As  $\|\vec{x} - \vec{u}\vec{u}^T\vec{x}\|^2$  (st.  $\vec{u}^T\vec{u}=1$ ) is convex and differentiable

$\therefore$  local minimizer is the global minimizer. (convex)

and global minimizer will satisfy:  $\frac{\partial \|\vec{x} - \vec{u}\vec{u}^T\vec{x}\|^2 - \lambda(\vec{u}^T\vec{u}-1)}{\partial \vec{u}} = \vec{0}$

$$\begin{aligned} & \|\vec{x} - \vec{u}\vec{u}^T\vec{x}\|^2 \\ &= (\vec{x} - \vec{u}\vec{u}^T\vec{x})^T (\vec{x} - \vec{u}\vec{u}^T\vec{x}) \\ &= \vec{x}^T \vec{x} - 2\vec{x}^T \vec{u}\vec{u}^T \vec{x} + \vec{x}^T \vec{u}\vec{u}^T \vec{u}\vec{u}^T \vec{x} \quad (\vec{u} \text{ is orthogonal}) \\ &= \vec{x}^T \vec{x} - \vec{u}^T \vec{x} \vec{x}^T \vec{u} \end{aligned}$$

$$\frac{\partial \|\vec{x} - \vec{u}\vec{u}^T\vec{x}\|^2 - \lambda(\vec{u}^T\vec{u}-1)}{\partial \vec{u}} = -2\vec{x}^T \vec{u} - 2\lambda \vec{u}$$

In PCA,  $\vec{u} = \vec{u}_e$  when  $\vec{u}_e$  is the eigen vector matrix of  $C$ .

$$C = \vec{x} \vec{x}^T$$

$$\therefore -2\vec{x}^T \vec{u} - 2\lambda \vec{u} = -2C \vec{u} - 2\lambda \vec{u}$$

$$\therefore C \vec{u} = \lambda \vec{u}$$

$$\therefore \frac{\partial \|\vec{x} - \vec{u}\vec{u}^T\vec{x}\|^2 - \lambda(\vec{u}^T\vec{u}-1)}{\partial \vec{u}} = -2(C \vec{u} - \lambda \vec{u}) = \vec{0}$$

$\therefore$  PCA solution is optimal.

b)

PCA feature transformation  $\tilde{x} = \vec{u}^T x$  where  $\vec{u}$  is the eigen-vector matrix of  $xx^T$ .

$$\text{Normalization: } x = \frac{x - \mu(x)}{\text{std}(x)}$$

$$\text{Scaling: } x = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

After PCA, each feature are uncorrelated with each other. This is not true for either normalization nor scaling. PCA can be used as dimensional reduction method, while normalization and scaling can not.

After normalization:  $\mu=0$  and  $\sigma = 1$ . This might not be true for scaling. After normalization, the relative mahalanobis distance between two points do not change. While after scaling, the relative euclidean distance between two points do not change.

Scaling works better when the data do not obey gaussian. After scaling the data lies in  $[0, 1]$  which is not true for normalization

- C)
- ① if we use gradient method to find optimal weights for both logistic regression and ridge regression, It will converge faster.
  - ② with feature scaling / feature normalization, we can avoid that large-scale feature matter more than those small-scale features. (Because the regularization might limit the weights in front of small-scale features, so they can not contribute to the prediction as much as those large-scale features.)

## Problem 5.

a)  $\pi_j \cdot \frac{1}{N(2\pi)^d |\Sigma_j|} \exp(-\frac{1}{2} (\vec{x} - \vec{\mu}_j)^T \Sigma_j^{-1} (\vec{x} - \vec{\mu}_j))$

b)  $L(x|\theta) = \prod_{i=1}^n \left( \sum_{j=1}^k \pi_j P(x_i|\theta_j) \right)$

if we put log over it:  $\Rightarrow \sum_{i=1}^n \log \sum_{j=1}^k \pi_j P(x_i|\theta_j)$

c) We could plot GMM against  $k$ . And use the  $k$  after which GMM likelihood stop increasing rapidly.

d) The log-likelihood is:

$$\begin{aligned} \log(P(x_i|\theta_j)) &= \log \prod_{m=1}^d (\theta_{jm})^{x_{im}} (1-\theta_{jm})^{(1-x_{im})} \\ &= \sum_{m=1}^d (x_{im} \log \theta_{jm} + (1-x_{im}) \log (1-\theta_{jm})) \end{aligned}$$

maximize  $B = \sum_{i=1}^n \sum_{j=1}^k Z_{ij} \log \left( \frac{\pi_j P(x_i|\theta_j)}{Z_{ij}} \right)$

maximize  $\sum_{i=1}^n \sum_{j=1}^k Z_{ij} \left( \log \pi_j + \sum_{m=1}^d (x_{im} \log \theta_{jm} + (1-x_{im}) \log (1-\theta_{jm})) - \log Z_{ij} \right)$

$$\begin{aligned} \frac{\partial}{\partial \theta_{jm}} \left( \sum_{i=1}^n \sum_{j=1}^k Z_{ij} \left( \log \pi_j + \sum_{m=1}^d (x_{im} \log \theta_{jm} + (1-x_{im}) \log (1-\theta_{jm})) - \log Z_{ij} \right) \right) \\ = \sum_{i=1}^n Z_{ij} \left( \frac{x_{im}}{\theta_{jm}} - \frac{1-x_{im}}{1-\theta_{jm}} \right) \end{aligned}$$

to maximize this: set the derivative to be 0:

$$\sum_{i=1}^n Z_{ij} \left( \frac{x_{im}}{\theta_{jm}} - \frac{1-x_{im}}{1-\theta_{jm}} \right) = 0$$

$$\sum_{i=1}^n Z_{ij} \left( \frac{x_{im} - \theta_{jm}}{\theta_{jm}(1-\theta_{jm})} \right) = 0$$

CSE517 HW4

Ziyang Jiao  
Flora Sun

475589  
475178

$$\therefore \sum_{i=1}^n z_{ij} (x_{im} - \theta_{jm}) = 0 \quad (0 < \theta_{jm} < 1)$$

$$\therefore \theta_{jm} = \frac{\sum_{i=1}^n z_{ij} x_{im}}{\sum_{i=1}^m z_{ij}}$$