

1. (a) We need  $d$  times additions so  $O(d)$  for  $a+v$   
We need  $d$  times multiplication, so  $O(d)$  for  $aTv$

(b)  $A$  has  $n \cdot d$  elements  $A+B$  needs  $O(n \cdot d)$  addition  
also needs  $O(n \cdot d)$  space to store

(c)  $A \begin{pmatrix} d \\ \vdots \\ d \end{pmatrix} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$  one multiplication needs  $d$ , there are  $n$   
so  $Av$  is  $O(nd)$

$A^T B$  is  $O(nd^2)$

(d)  $A^T B v$  method 1 :  $(A^T B)$  takes  $O(nd^2)$   $(A^T B) \cdot v$  takes  $O(d^2)$

$O(nd^2)$  dominant

method 2:  $B \cdot v$   $A^T (B \cdot v)$

takes  $O(nd)$  takes  $O(d \cdot n)$

so  $O(nd)$  after drops constant

the dominate terms

2. 1.  $X^T X \rightarrow X^T X + \lambda I \rightarrow (X^T X + \lambda I)^{-1}$

$$O(nl^2) + O(l^2) + O(l^3) = O(l^3 + nl^2) \quad (1)$$

2.  $XX^T \rightarrow XX^T + \lambda I \rightarrow (XX^T + \lambda I)^{-1}$

$$O(n^2 l) + O(n^2) + O(n^3) = O(n^3 + n^2 l) \quad (2)$$

3. if  $n \gg l$  (2)'s complexity becomes  $O(n^3)$ , (1) is preferable

if  $n \ll l$  (1)'s complexity becomes  $O(l^3)$ , (2) is preferable

4. 1. if  $l \gg n$ , like the gene dataset, then we can use the kernel trick to reduce the computation cost

2. Sometimes it's easier to deal with the data in transformed feature space, such as predictors become linear in feature space, while not in the original space.

$$5. (a) K(x, z) = (x_1 x_2, \frac{x_1^2}{\sqrt{2}}, \frac{x_2^2}{\sqrt{2}})^T (z_1 z_2, \frac{z_1^2}{\sqrt{2}}, \frac{z_2^2}{\sqrt{2}})$$

$$= (\frac{1}{\sqrt{2}} x_1 z_1 + \frac{1}{\sqrt{2}} x_2 z_2)^2$$

$$= (\frac{1}{\sqrt{2}} x^T z)^2$$

$$(b) K(x, z) = (1, x, \frac{x^2}{2!}, \dots)^T (1, z, \frac{z^2}{2!}, \dots)$$

$$= 1^2 + xz + \frac{(xz)^2}{2!} + \frac{(xz)^3}{3!} + \dots$$

$$= e^{xz}$$

$$(c) K(x, z) = (1 + x^T z)^2 + x^T z$$

$$= 1 + 3x_1 z_1 + 3x_2 z_2 + 2x_1 z_1 x_2 z_2 + x_1^2 z_1^2 + x_2^2 z_2^2$$

$$= (x_1^2, x_2^2, \sqrt{2}x_1 x_2, \sqrt{3}x_1, \sqrt{3}x_2, 1)^T (z_1^2, z_2^2, \sqrt{2}z_1 z_2, \sqrt{3}z_1, \sqrt{3}z_2, 1)$$

$$= \langle \phi(x), \phi(z) \rangle$$

$$(d) e^{-(x-z)^2} \text{ can not be written as } \phi(x)\phi(z)$$

as  $e^{-(x-z)^2} = 1 + (-(x-z)^2) + \dots$  where  $-(x-z)^2$  cannot be splitted as a product of  $f(x)$  and  $g(z)$

6.

$$\hat{\theta} = (\Phi^T \Phi + \lambda I_d)^{-1} \Phi^T y$$

or

$$\hat{\theta} = \Phi^T (\Phi \Phi^T + \lambda I_n)^{-1} y$$

7. (1) suppose a new point is  $x_{\text{new}}$ , we'd like to predict  $\hat{y}_{\text{new}}$

with kernelized  $\hat{y}_{\text{new}} = \langle \phi(x_{\text{new}}), \hat{\theta} \rangle = \phi(x_{\text{new}})^T (\Phi \Phi^T + \lambda I_n)^{-1} y$

$$= (K(x_1, x_{\text{new}}), \dots, K(x_n, x_{\text{new}})) (K + \lambda I)^{-1} y$$

$$K = \begin{pmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots \\ K(x_2, x_1) & & \\ \vdots & & \\ & & K(x_n, x_n) \end{pmatrix}$$

$n \times n$

computing  $K$  takes  $n \times n$  times the cost of  $K(x, z)$

$$K(x, z) = (1 + x^T z)^p \text{ where } x^T z \text{ takes } O(L)$$

assume  $p = 2^m$  ( $p = 2^{m+1}$  doesn't affect the cost much)

$$\text{Then } (1 + x^T z)^p = ((1 + x^T z)^2)^{2^{m-1}} = (1 + x^T z)^{\underbrace{2 \times 2 \times 2 \dots \times 2}_m}$$

Squaring a number can be seen as  $O(1)$

We need  $m = \log p$  times squaring

So the overall cost for computing  $K(x, z)$  is  $O(L + \log p)$

inverting a  $n \times n$  matrix is  $O(n^3)$

The two dominant costs are  $O(n^3 + n^2(L + \log p))$

(2)  $\hat{y}_{\text{new}} = \langle \phi(x_{\text{new}}), \hat{\theta} \rangle = \phi(x_{\text{new}})^T (\Phi^T \Phi + \lambda I_d)^{-1} \Phi^T y$  Non-kernelized

$\Phi^T \Phi$  is  $d \times d$

$d \left( \begin{array}{c} n \\ \hline \end{array} \right) n \left( \begin{array}{c} d \\ \hline \end{array} \right)$  the total computation for  $\Phi^T \Phi$  is  $O(d^2 n)$   
 $\Phi^T$   $\Phi$  while inverting  $d \times d$  matrix takes  $O(d^3)$

These two costs dominate, so the overall is  $O(d^3 + d^2 n)$

Compare: if  $n \gg d$ , we prefer the non-kernelized, which is faster than the kernelized

if  $n \ll d$ , we prefer the kernelized  $(\Phi^T \Phi + \lambda I_d)^{-1} \Phi^T y$ , which has the cost lower than non-kernelized

4.2

$$(a) \hat{y} = \underset{K}{\operatorname{argmax}} P(Y=K) P(X|Y=K)$$

 $\pi(u_k) = P(Y=k)$ ,  $P(X|Y=k)$  is multivariate normal

$$\pi(u_2) \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(X-u_2)^T \Sigma^{-1}(X-u_2)\right\} > \pi(u_1) \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(X-u_1)^T \Sigma^{-1}(X-u_1)\right\}$$

↓ take log on both sides

$$-\frac{1}{2}(X-u_2)^T \Sigma^{-1}(X-u_2) + \ln(\pi(u_2)) + \text{constants} > -\frac{1}{2}(X-u_1)^T \Sigma^{-1}(X-u_1) + \ln(\pi(u_1)) + \text{constants}$$

$$X^T \Sigma^{-1} u_2 - \frac{1}{2} u_2^T \Sigma^{-1} u_2 + \ln(\pi(u_2)) > X^T \Sigma^{-1} u_1 - \frac{1}{2} u_1^T \Sigma^{-1} u_1 + \ln(\pi(u_1))$$

$$X^T \Sigma^{-1} (u_2 - u_1) > \frac{1}{2} u_1^T \Sigma^{-1} u_2 - \frac{1}{2} u_1^T \Sigma^{-1} u_1 + \underbrace{\ln\left(\frac{N_1}{N}\right) - \ln\left(\frac{N_2}{N}\right)}_{\ln\left(\frac{N_2}{N_1}\right)}$$

if left likelihood > right hand side likelihood, we classify to class 2  $\ln\left(\frac{N_2}{N_1}\right)$

if left < right, we classify to class 1

$$(c) X^T X \beta = X^T Y$$

$$\begin{pmatrix} 1^T & 1^T \\ X^{(2)T} & X^{(1)T} \end{pmatrix} \begin{pmatrix} 1 & X^{(2)} \\ 1 & X^{(1)} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta \end{pmatrix} = \begin{pmatrix} 1^T & 1^T \\ X^{(2)T} & X^{(1)T} \end{pmatrix} \begin{pmatrix} N/N_2 \\ -N/N_1 \end{pmatrix}$$

$$\begin{pmatrix} N & N\hat{u}_1^T + N_2\hat{u}_2^T \\ N\hat{u}_1 + N_2\hat{u}_2 & X^{(2)T}X^{(2)} + X^{(1)T}X^{(1)} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta \end{pmatrix} = \begin{pmatrix} 0 \\ N(\hat{u}_2 - \hat{u}_1) \end{pmatrix} \quad N\beta_0 + (N\hat{u}_1^T + N_2\hat{u}_2^T)\beta = 0 \quad (1)$$

$$X^{(2)T}X^{(2)} + X^{(1)T}X^{(1)} = \bar{X}^{(2)T}\bar{X}^{(2)} + \bar{X}^{(1)T}\bar{X}^{(1)} + N_1\hat{u}_1\hat{u}_1^T + N_2\hat{u}_2\hat{u}_2^T = (N-2)\hat{\Sigma} + N_1\hat{u}_1\hat{u}_1^T + N_2\hat{u}_2\hat{u}_2^T$$

$$(N_1\hat{u}_1 + N_2\hat{u}_2)\beta_0 + [(N-2)\hat{\Sigma} + N_1\hat{u}_1\hat{u}_1^T + N_2\hat{u}_2\hat{u}_2^T]\beta = N(\hat{u}_2 - \hat{u}_1) \quad \text{plug in (1)}$$

$$[(N-2)\hat{\Sigma} + \frac{N_1 N_2}{N}(\hat{u}_2 - \hat{u}_1)(\hat{u}_2 - \hat{u}_1)^T]\beta = N(\hat{u}_2 - \hat{u}_1)$$

(c) from part (b) we have

$$(N-2)\Sigma\beta + \frac{N_1 N_2}{N} (\hat{u}_2 - \hat{u}_1) (\hat{u}_2 - \hat{u}_1)^T \beta = N(\hat{u}_2 - \hat{u}_1)$$

Since  $(\hat{u}_2 - \hat{u}_1)^T \beta$  is scalar

$\Sigma\beta$  is a linear combination of  $\hat{u}_2 - \hat{u}_1$ , and  $\Sigma\beta$  is a multiple of  $\hat{u}_2 - \hat{u}_1$   
hence it's in the direction of  $\hat{u}_2 - \hat{u}_1$

$$\Sigma\beta \propto \hat{u}_2 - \hat{u}_1$$

$$\beta \propto \Sigma^{-1}(\hat{u}_2 - \hat{u}_1)$$

(d) Since in part (b) and (c), we assume labels are different,  
so if two classes are distinct,  $\beta \propto \Sigma^{-1}(\hat{u}_2 - \hat{u}_1)$  still holds.

$$(e) \quad \beta = \lambda \Sigma^{-1}(\hat{u}_2 - \hat{u}_1) \quad \beta_0 = -\frac{\lambda}{N} (N_1 \hat{u}_1^T + N_2 \hat{u}_2^T) \Sigma^{-1}(\hat{u}_2 - \hat{u}_1)$$

$$\beta_0 + X^T \beta = X^T \lambda \Sigma^{-1}(\hat{u}_2 - \hat{u}_1) \geq \frac{\lambda}{N} (N_1 \hat{u}_1^T + N_2 \hat{u}_2^T) \Sigma^{-1}(\hat{u}_2 - \hat{u}_1) \quad (2)$$

We see when  $N_1 = N_2$  (2) is equivalent to LDA in part (a).  $\ln(\frac{N_2}{N_1}) = 0$   
When  $N_1 \neq N_2$ , they are different.