

Unique name: lyifeng UMID:44603814

EECS 545 HW 1

Study group: Kunyi Lu, Chenrui Shu

Problem 1:

(a) The code is attached. The error rate on test data is : 0.001250.

(b) For a two-class Naive Bayes classifier, y equal c (1) or \bar{c} (0) and

x has dimension D with every element equal 0 or 1, we can have:

$$\mathbf{w}_c = (\log \frac{\theta_{c11}}{1 - \theta_{c11}}, \log \frac{\theta_{c21}}{1 - \theta_{c21}}, \dots, \log \frac{\theta_{cD1}}{1 - \theta_{cD1}}, \log \pi_c + \sum_{d=1}^D \log(1 - \theta_{cd1}))^T$$

$$\mathbf{w}_{\bar{c}} = (\log \frac{\theta_{\bar{c}11}}{1 - \theta_{\bar{c}11}}, \log \frac{\theta_{\bar{c}21}}{1 - \theta_{\bar{c}21}}, \dots, \log \frac{\theta_{\bar{c}D1}}{1 - \theta_{\bar{c}D1}}, \log \pi_{\bar{c}} + \sum_{d=1}^D \log(1 - \theta_{\bar{c}d1}))^T$$

$$\mathbf{w} = \mathbf{w}_c - \mathbf{w}_{\bar{c}}$$

Then we can get:

$$\mathbf{w}^T * \varphi(x) > 0 \Leftrightarrow (\mathbf{w}_c^T - \mathbf{w}_{\bar{c}}^T) * \varphi(x) > 0 \Leftrightarrow \mathbf{w}_c^T * \varphi(x) > \mathbf{w}_{\bar{c}}^T * \varphi(x)$$

Since $(1 - \theta_{cd1}) = \theta_{cd0}$, $(1 - \theta_{\bar{c}d1}) = \theta_{\bar{c}d0}$

$$\mathbf{w}_c^T * \varphi(x) = \log \pi_c * \prod_{d=1}^D (\frac{\theta_{cd1}}{1 - \theta_{cd1}})^{I(x_d=1)} * \prod_{d=1}^D (1 - \theta_{cd1}) =$$

$$\log \pi_c * \prod_{d=1}^D (\theta_{cd1})^{I(x_d=1)} * \prod_{d=1}^D (\theta_{cd0})^{I(x_d=0)} = \log P(y = c | \mathbf{x} = (x_1, x_2, \dots, x_D)^T)$$

$$\mathbf{w}_{\bar{c}}^T * \varphi(x) = \log \pi_{\bar{c}} * \prod_{d=1}^D (\frac{\theta_{\bar{c}d1}}{1 - \theta_{\bar{c}d1}})^{I(x_d=1)} * \prod_{d=1}^D (1 - \theta_{\bar{c}d1})$$

$$= \log \pi_{\bar{c}} * \prod_{d=1}^D (\theta_{\bar{c}d1})^{I(x_d=1)} * \prod_{d=1}^D (\theta_{\bar{c}d0})^{I(x_d=0)} = \log P(y = \bar{c} | \mathbf{x} = (x_1, x_2, \dots, x_D)^T)$$

So we can point out that

$$\begin{aligned} \mathbf{w}_c^T * \phi(\mathbf{x}) > \mathbf{w}_{\bar{c}}^T * \phi(\mathbf{x}) &\Leftrightarrow P(\mathbf{y} = c \mid \mathbf{x} = (x_1, x_2, \dots, x_d)^T) > P(\mathbf{y} = \bar{c} \mid \mathbf{x} = (x_1, x_2, \dots, x_d)^T) \\ &\Leftrightarrow P(\mathbf{y} = 1 \mid \mathbf{x} = (x_1, x_2, \dots, x_d)^T) > P(\mathbf{y} = 0 \mid \mathbf{x} = (x_1, x_2, \dots, x_d)^T) \end{aligned}$$

Which means the linear classifier listed in the question has the same function with two-class Naive Bayes binary classifier hence the two-class Naive Bayes classifier can be represented as a linear classifier. Therefore Naive Bayes binary classifier is a linear classifier.

Problem 2:

Pr. 2. Since k_1, k_2 are positive-definite kernel. For any $X = (x_1 \dots x_n)^T$ $X^T k_1 X > 0$ $X^T k_2 X > 0$
 [A matrix G is positive-definite $\Leftrightarrow X^T G X > 0$ for any X]

(i) $\because X^T k_1 X > 0$ for any X $\therefore a X^T k_1 X > 0$ when $a > 0$ $\therefore X^T (a k_1) X > 0$
 $\therefore K(x, x') = a k_1(x, x')$ is a valid kernel

(ii) $\because X^T k_1 X > 0$ $X^T k_2 X > 0$ for any X . $\therefore X^T k_1 X + X^T k_2 X = X^T (k_1 + k_2) X > 0$
 $\therefore K(x, x') = k_1(x, x') + k_2(x, x')$ is a valid kernel

(iii) Do a decomposition that $k_2 = R^T R$ where $R = (r_{ij})$ is a real-valued matrix
 Then for any Y . $Y^T K Y = \sum_{i,j} y_i k(x_i, x_j) y_j$

$$= \sum_{i,j} y_i k_1(x_i, x_j) k_2(x_i, x_j) y_j$$

$$= \sum_{i,j} y_i k_1(x_i, x_j) \left(\sum_k r_{ki} r_{kj} \right) y_j$$

$$= \sum_k \sum_{i,j} (r_{ki} y_i) k_1(x_i, x_j) (r_{kj} y_j)$$

$$= \sum_k h_k^T K h_k > 0 \quad [h_{ki} = r_{ki} y_i]$$

$\therefore K(x, x') = k_1(x, x') + k_2(x, x')$ is a valid kernel

(iv) $\because K = f(x_1) f(x_2)$ \therefore For any Y , $Y^T K Y = \sum_{i,j} y_i f(x_i) f(x_j) y_j = [y_1^2 f(x_1)^2 + \dots + y_n^2 f(x_n)^2 + 2 y_1 f(x_1) + \dots + 2 y_n f(x_n)]$
 $= \left(\sum_{i=1}^n y_i f(x_i) \right)^2$
 > 0

$\therefore K = f(x) f(x')$ is a valid kernel

(v) from (iii) we know $k = k_1 k_2$ is a valid kernel

$\therefore k = k_1^d = k_1 \cdot k_1 \cdot k_1 \dots k_1$, every time we multiply two valid kernel

\therefore By apply (iii), $k = k_1^d$ is a valid kernel

(vi) $k' = p(k)$ where k is a arbitrary kernel.

① when k is a valid kernel, Then $k' = \sum a_k^d$, k' is the sum of polynomial term with product of valid kernels and positive coefficient. By applying (i) and (v), $k' = p(k)$ is a valid kernel

② If k is not a valid kernel. Then $k' = p(k)$ is not a valid kernel.

(b) The Gaussian kernel $k(x, x') = \exp(-\frac{\|x - x'\|^2}{2\sigma^2}) = \exp(-\frac{(x - x')^2}{2\sigma^2}) = \exp(-\frac{x^2 - 2xx' + x'^2}{2\sigma^2})$

As $e^x = 1 + x + \frac{x^2}{2} + \dots + \frac{x^n}{n!}$ (By power series and Taylor's theorem)

$$\therefore k(x, x') = \exp(-\frac{x^2 - 2xx' + x'^2}{2\sigma^2}) \left(1 + \frac{2xx'}{2\sigma^2} + \frac{(2xx')^2}{2!} + \dots + \frac{(2xx')^n}{n!} \right)$$

$$= \exp(-\frac{x^2 - x'^2}{2\sigma^2}) \left(1 + \sqrt{\frac{1}{\sigma^2}} x \sqrt{\frac{1}{\sigma^2}} x' + \dots + \sqrt{\frac{1}{\sigma^{2n}}} x^n \sqrt{\frac{1}{\sigma^{2n}}} x'^n \right)$$

$$= \phi(x)^T \phi(x')$$

$$\text{where } \phi(x) = \exp(-\frac{x^2}{2\sigma^2}) \left[1, \sqrt{\frac{1}{\sigma^2}} x, \sqrt{\frac{1}{\sigma^{2 \cdot 2!}}} x^2, \dots, \sqrt{\frac{1}{\sigma^{2 \cdot n!}}} x^n \right]^T$$

Problem 3:

(a) The prediction for $y(x)$ for Kernel Perceptron when given a new sample x' is:

$$\bar{y} = \text{sgn}\left(\sum_i \alpha_i y_i K(x_i, x')\right)$$

And the pseudocode is shown below:

Initialize α to an all-zeros vector of length n , the number of training samples.

For some fixed number of iterations, or until some stopping criterion is met:

For each training example x_j with ground truth label $y_j \in \{-1, 1\}$:

Let $\bar{y} = \text{sgn}\left(\sum_i \alpha_i y_i K(x_i, x_j)\right)$

If $\hat{y} \neq y_i$, perform an update by incrementing the mistake counter:

$$a_j \leftarrow a_j + 1$$

(b) Two plots are shown in Fig 1 and Fig 2.

decision boundary learned by Gaussian Kernel Perceptron with $\sigma = 0.1$

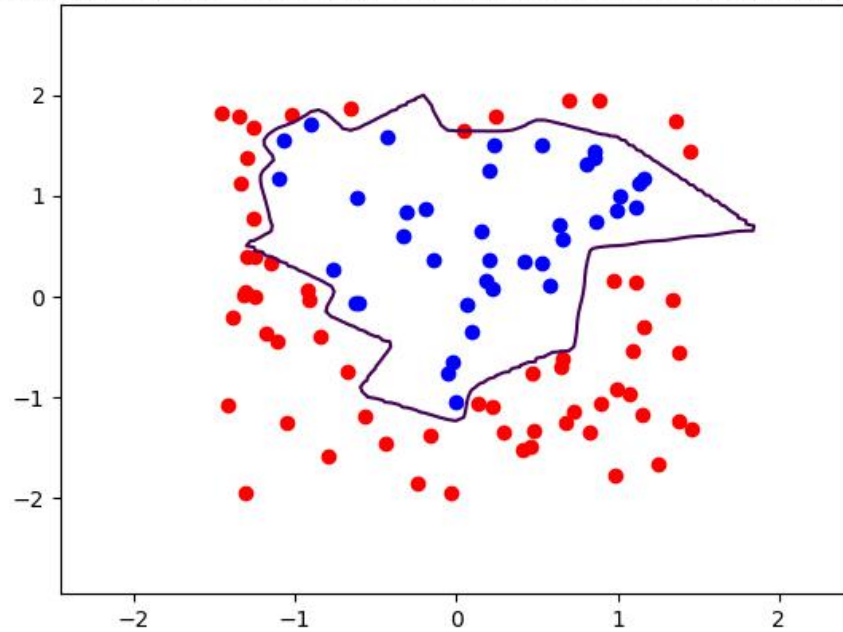


Fig 1 decision boundary with $\sigma = 0.1$

decision boundary learned by Gaussian Kernel Perceptron with $\sigma = 1.0$

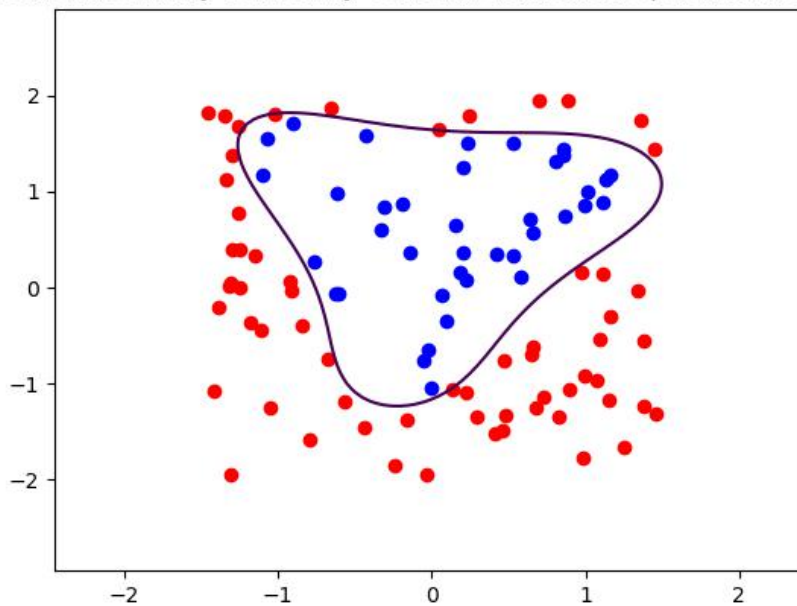


Fig 2 decision boundary with $\sigma = 1.0$

Problem 4:

Two plots are shown in Fig 3 and Fig 4.

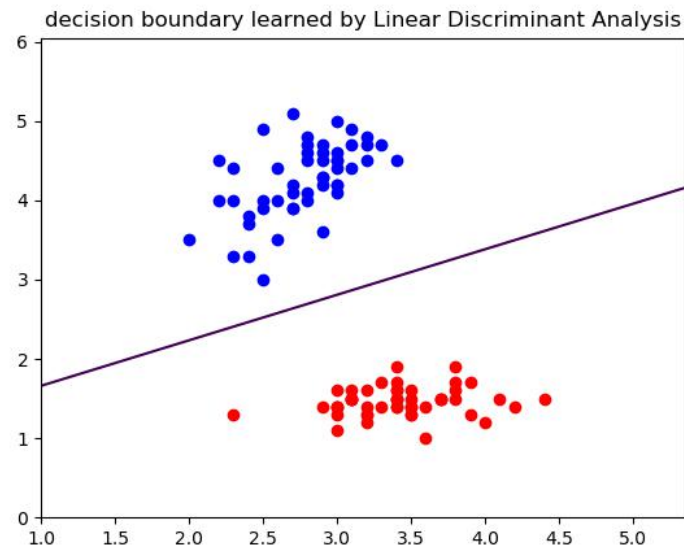


Fig 3 decision boundary by Linear Discriminant Analysis

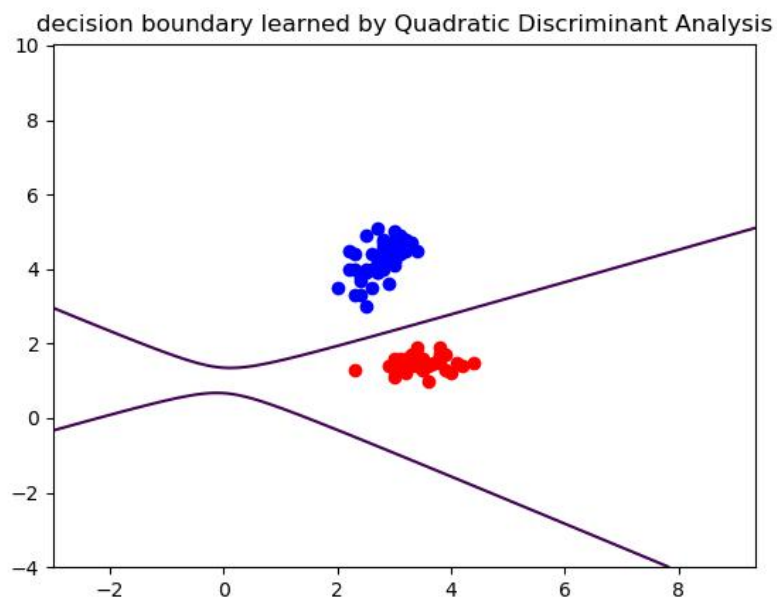


Fig 4 decision boundary by Quadratic Discriminant Analysis