

## Linear classification

### Q1:

a)

The posterior distribution is a sigmoid of a linear function, or equivalently the Bernoulli distribution.

b)

According to the Bayes' Rule,

$$\begin{aligned} p(y = 1|x) &= \frac{p(x|y = 1)p(y = 1)}{p(x|y = 0)p(y = 0) + p(x|y = 1)p(y = 1)} \\ p(y = 0|x) &= \frac{p(x|y = 0)p(y = 0)}{p(x|y = 0)p(y = 0) + p(x|y = 1)p(y = 1)} \end{aligned} \quad (1)$$

As

$$p(y = 0) = p(y = 1) = \frac{1}{2}$$

then

$$\begin{aligned} p(y = 1|x) &= \frac{p(x|y = 1)}{p(x|y = 0) + p(x|y = 1)} \\ p(y = 0|x) &= \frac{p(x|y = 0)}{p(x|y = 0) + p(x|y = 1)} \end{aligned} \quad (2)$$

Let

$$\begin{aligned} p(y = 1|x) - p(y = 0|x) &\geq 0 \\ \text{namely} \\ \frac{p(x|y = 1) - p(x|y = 0)}{p(x|y = 0) + p(x|y = 1)} &\geq 0 \end{aligned} \quad (3)$$

As

$$p(x|y = 0) + p(x|y = 1) \geq 0$$

Then the expression becomes as

$$\begin{aligned} p(x|y = 1) - p(x|y = 0) &\geq 0 \\ \lambda_1 e^{-\lambda_1 x} - \lambda_0 e^{-\lambda_0 x} &\geq 0 \end{aligned} \quad (4)$$

As  $\lambda_i > 0$  and  $\lambda_0 \neq \lambda_1$

$$(\lambda_0 - \lambda_1)x \geq \ln\left(\frac{\lambda_1}{\lambda_0}\right) \quad (5)$$

x will be classified as class 1 when

$$\begin{cases} x \geq \frac{\ln\lambda_1 - \ln\lambda_0}{\lambda_0 - \lambda_1}, & \text{if } \lambda_0 > \lambda_1 \\ x \leq \frac{\ln\lambda_1 - \ln\lambda_0}{\lambda_0 - \lambda_1}, & \text{if } \lambda_0 < \lambda_1 \end{cases}$$

**Q2:**

As we discussed in the class, in the extreme situation, if the dataset is linear separable, the sigmoid function will tend to have the shape of step function. This means  $\mathbf{w}^T \mathbf{x} \rightarrow \infty$ , as

$$\lim_{\mathbf{w}^T \mathbf{x} \rightarrow \infty} \sigma(\mathbf{w}^T \mathbf{x}) = \lim_{\mathbf{w}^T \mathbf{x} \rightarrow \infty} \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = 1 \quad (6)$$

Similar to the linear regression problem, we can use a regularization term to penalize large weights.

**Q3:**

The sigmoid in 2-class has the form:

$$\begin{aligned} p(y = 0|\mathbf{x}) &= \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \\ p(y = 1|\mathbf{x}) &= \frac{e^{-\mathbf{w}^T \mathbf{x}}}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \end{aligned} \quad (7)$$

The softmax function in 2-class:

$$\begin{aligned} p(y = c|\mathbf{x}) &= \frac{e^{\mathbf{w}_c^T \mathbf{x}}}{\sum_{c'} e^{\mathbf{w}_{c'}^T \mathbf{x}}} \\ p(y = 0|\mathbf{x}) &= \frac{e^{\mathbf{w}_0^T \mathbf{x}}}{e^{\mathbf{w}_0^T \mathbf{x}} + e^{\mathbf{w}_1^T \mathbf{x}}} \\ &= \frac{1}{1 + e^{(\mathbf{w}_1 - \mathbf{w}_0)^T \mathbf{x}}} \\ p(y = 1|\mathbf{x}) &= \frac{e^{\mathbf{w}_1^T \mathbf{x}}}{e^{\mathbf{w}_0^T \mathbf{x}} + e^{\mathbf{w}_1^T \mathbf{x}}} \\ &= \frac{e^{(\mathbf{w}_1 - \mathbf{w}_0)^T \mathbf{x}}}{1 + e^{(\mathbf{w}_1 - \mathbf{w}_0)^T \mathbf{x}}} \end{aligned} \quad (8)$$

Assume  $\mathbf{w}_1 - \mathbf{w}_0 = -\mathbf{w}$ , then the expressions are equivalent.

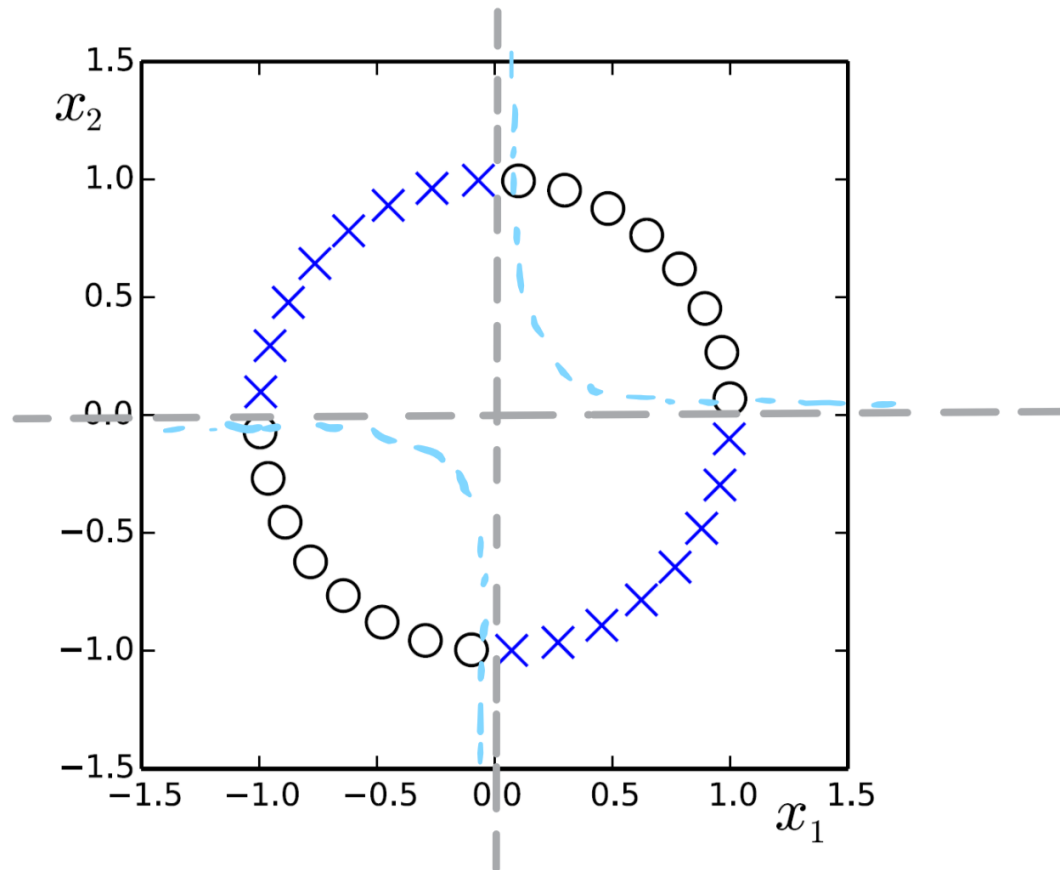
**Q4:**

Figure 1:  $\phi(x_1, x_2) = x_1x_2$  separates the data

The basis equation can separate the data as in the figure.

$$\phi(x_1, x_2) = x_1x_2$$