

Question 1

1. (1) True
- (2) False
- (3) True
- (4) True
- (5) True
- (6) True
- (7) True
- (8) False

2. Discriminative classifier learns a model for $p(t|x)$. An example is the logistic regression. Generative classifier learns a model for $p(x|t)$ with a prior assumption on the marginal density $P(x)$. One example is the Gaussian Discriminant Analysis.

3. (a) $\sigma_{ML} = \arg \max_{\sigma} \prod_{i=1}^N p(t_i | x_i, \sigma) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2} (t_i^{(1)} - \beta_0 - \beta^T x_i)^2\right)$ Since ϵ is i.i.d then t is i.i.d

$$= \arg \max_{\sigma} \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} (t_i - \beta_0 - \beta^T x_i)^2$$

(b) Define $L(\sigma^2) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} (t_i^{(1)} - \beta_0 - \beta^T x_i)^2$

$$\frac{\partial L(\sigma^2)}{\partial \sigma^2} = \sum_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma^2} \left(-\frac{1}{2} \frac{1}{(\sqrt{2\pi}\sigma^2)^3}\right) 2\pi - \frac{(t_i^{(1)} - \beta_0 - \beta^T x_i)^2}{2} \left(-\frac{1}{\sigma^4}\right)$$

$$= \sum_{i=1}^N \frac{1}{\sigma^2} - \frac{\pi}{2\pi\sigma^2} + \frac{(t_i^{(1)} - \beta_0 - \beta^T x_i)^2}{2} \frac{1}{\sigma^4}$$

$$= -\frac{N}{2\sigma^2} + \frac{1}{\sigma^4} \frac{\sum_{i=1}^N (t_i^{(1)} - \beta_0 - \beta^T x_i)^2}{2}$$

Set $\frac{\partial L}{\partial \sigma}$ to be 0

$$\frac{N}{2\sigma_{ML}^2} = \frac{1}{2\sigma_{ML}^4} \sum_{i=1}^N (t_i^{(1)} - \beta_0 - \beta^T x_i)^2$$

$$\sigma_{ML}^2 = \frac{1}{N} \sum_{i=1}^N (t_i^{(1)} - \beta_0 - \beta^T x_i)^2$$

4. Assume that we have K class with target value 1 to k

All data information is stored in data variable and target information is stored in target variable

```
X1min = np.min(data[:,0])
```

```
X1max = np.max(data[:,0])
```

```
X2min = np.min(data[:,1])
```

```
X2max = np.max(data[:,1])
```

```
X1 = np.linspace(X1min, X1max, 50)
```

```
X2 = np.linspace(X2min, X2max, 50)
```

```
xx1, xx2 = np.meshgrid(X1, X2)
```

```
Xprediction = np.hstack((xx1.reshape(-1,1), xx2.reshape(-1,1)))
```

```
from sklearn.linear_model import LinearRegression
```

```
prediction = np.zeros((len(xx1).flatten(), K-1))
```

```

for k in np.arange(k-1):
    indices_class_k = np.where(target == k)
    target_one_vs_rest = np.zeros((len(target), 1))
    target_one_vs_rest[indices_class_k] = 1
    reg = LinearRegression()
    reg.fit(data, target_one_vs_rest)
    prediction[:, k] = np.squeeze(reg.predict(data - prediction) > 1/2)

```

```

final_prediction = np.zeros((len(x).flatten(), 1))
for i in np.arange(len(final_prediction)):
    if np.argmax(prediction[i,:]) == 1:
        final_prediction[i] = np.argmax(prediction[i,:])
    elif np.argmax(prediction[i,:]) == 1:
        final_prediction[i] = k
    else:
        final_prediction[i] = k-1

```

Question 2

1. (1) False
- (2) True
- (3) True
- (4) False
- (5) False
- (6) False

2. (a) The normal equation is $X^T t = X^T X \beta$ where $X = \begin{bmatrix} \bar{x}_1^{(1)} & \bar{x}_2^{(1)} \\ \vdots & \vdots \\ \bar{x}_1^{(N)} & \bar{x}_2^{(N)} \end{bmatrix}$ $t = \begin{bmatrix} \bar{t}^{(1)} \\ \vdots \\ \bar{t}^{(N)} \end{bmatrix}$ $\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$

$$\begin{bmatrix} \bar{x}_1^{(1)} & \dots & \bar{x}_1^{(N)} \\ \bar{x}_2^{(1)} & \dots & \bar{x}_2^{(N)} \end{bmatrix} \begin{bmatrix} \bar{t}^{(1)} \\ \vdots \\ \bar{t}^{(N)} \end{bmatrix} = \begin{bmatrix} \bar{x}_1^{(1)} & \dots & \bar{x}_1^{(N)} \\ \bar{x}_2^{(1)} & \dots & \bar{x}_2^{(N)} \end{bmatrix} \begin{bmatrix} \bar{x}_1^{(1)} & \bar{x}_2^{(1)} \\ \vdots & \vdots \\ \bar{x}_1^{(N)} & \bar{x}_2^{(N)} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$$

$$\begin{bmatrix} \sum_{i=1}^N \bar{x}_1^{(i)} \bar{t}^{(i)} \\ \sum_{i=1}^N \bar{x}_2^{(i)} \bar{t}^{(i)} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N \bar{x}_1^{(i)2} & \sum_{i=1}^N \bar{x}_1^{(i)} \bar{x}_2^{(i)} \\ \sum_{i=1}^N \bar{x}_2^{(i)} \bar{x}_1^{(i)} & \sum_{i=1}^N \bar{x}_2^{(i)2} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \Rightarrow \frac{1}{N} \begin{bmatrix} \sum_{i=1}^N \bar{x}_1^{(i)} \bar{t}^{(i)} \\ \sum_{i=1}^N \bar{x}_2^{(i)} \bar{t}^{(i)} \end{bmatrix} = \frac{1}{N} \begin{bmatrix} \sum_{i=1}^N \bar{x}_1^{(i)2} & \sum_{i=1}^N \bar{x}_1^{(i)} \bar{x}_2^{(i)} \\ \sum_{i=1}^N \bar{x}_2^{(i)} \bar{x}_1^{(i)} & \sum_{i=1}^N \bar{x}_2^{(i)2} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 1 & \frac{1}{N} \sum_{i=1}^N \bar{x}_1^{(i)} \bar{x}_2^{(i)} \\ \frac{1}{N} \sum_{i=1}^N \bar{x}_1^{(i)} \bar{x}_2^{(i)} & 1 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N \bar{x}_1^{(i)} \bar{t}^{(i)} \\ \sum_{i=1}^N \bar{x}_2^{(i)} \bar{t}^{(i)} \end{bmatrix}$$

Since $\sum_{i=1}^N \bar{x}_R^{(i)2} = \sum \frac{(X_R^{(i)} - \frac{1}{N} \sum X_R^{(i)})^2}{\sigma_R^2} = \frac{\sum (X_R^{(i)} - \frac{1}{N} \sum X_R^{(i)})^2}{N \sum (X_R^{(i)} - \frac{1}{N} \sum X_R^{(i)})^2} = N$

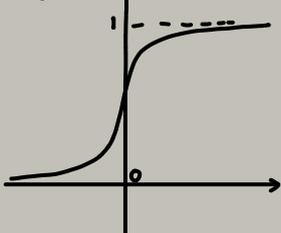
$$\begin{aligned} \beta_{12} &= \frac{1}{N} \sum_{i=1}^N \bar{x}_1^{(i)} \bar{x}_2^{(i)} = \frac{1}{N} \sum_{i=1}^N \frac{X_1^{(i)} - \frac{1}{N} \sum X_1^{(i)}}{\sigma_1} \frac{X_2^{(i)} - \frac{1}{N} \sum X_2^{(i)}}{\sigma_2} \\ &= \frac{1}{N \sigma_1 \sigma_2} \sum (X_1^{(i)} - \frac{1}{N} \sum X_1^{(i)}) (X_2^{(i)} - \frac{1}{N} \sum X_2^{(i)}) \\ &= \frac{1}{N \sigma_1 \sigma_2} \sum \left(X_1^{(i)} X_2^{(i)} - \frac{X_1^{(i)}}{N} \sum X_2^{(i)} - \frac{X_2^{(i)}}{N} \sum X_1^{(i)} + \frac{1}{N^2} \sum X_1^{(i)} \sum X_2^{(i)} \right) \\ &= \frac{1}{N \sigma_1 \sigma_2} \left(\sum X_1^{(i)} X_2^{(i)} - \frac{1}{N} \sum X_2^{(i)} \sum X_1^{(i)} \right) \end{aligned}$$

$$(b) X^T X = \begin{bmatrix} \bar{x}_1^{(1)} & \dots & \bar{x}_1^{(N)} \\ \bar{x}_2^{(1)} & \dots & \bar{x}_2^{(N)} \end{bmatrix} \begin{bmatrix} \bar{x}_1^{(1)} & \bar{x}_2^{(1)} \\ \vdots & \vdots \\ \bar{x}_1^{(N)} & \bar{x}_2^{(N)} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N \bar{x}_1^{(i)2} & \sum_{i=1}^N \bar{x}_1^{(i)} \bar{x}_2^{(i)} \\ \sum_{i=1}^N \bar{x}_2^{(i)} \bar{x}_1^{(i)} & \sum_{i=1}^N \bar{x}_2^{(i)2} \end{bmatrix} = N \begin{bmatrix} 1 & r_{12} \\ r_{12} & 1 \end{bmatrix}$$

$$(X^T X)^{-1} = \frac{N}{1-r_{12}^2} \begin{bmatrix} 1 & -r_{12} \\ -r_{12} & 1 \end{bmatrix}$$

$\det(X^T X) \neq 0 \Leftrightarrow 1-r_{12}^2 \neq 0$ to make $X^T X$ invertible, which means $r_{12} \neq \pm 1$

3. a) sigmoid function



$$(b) y(x, w) = \sigma(w_{11}^{(4)} z_1^{(3)} + w_{12}^{(4)} z_2^{(3)} + w_{10}^{(4)})$$

$$\text{where } \begin{cases} z_1^{(0)} = \sigma(x) \\ z_1^{(1)} = \sigma(w_{11}^{(1)} \sigma(x) + w_{10}^{(1)}) \\ z_2^{(1)} = \sigma(w_{21}^{(1)} \sigma(x) + w_{20}^{(1)}) \\ z_1^{(2)} = \sigma(w_{11}^{(2)} \sigma(w_{11}^{(1)} \sigma(x) + w_{10}^{(1)}) + w_{12}^{(2)} \sigma(w_{21}^{(1)} \sigma(x) + w_{20}^{(1)}) + w_{10}^{(2)}) \\ z_1^{(3)} = \sigma(w_{11}^{(3)} z_1^{(2)} + w_{10}^{(3)}) \\ z_2^{(3)} = \sigma(w_{21}^{(3)} z_1^{(2)} + w_{20}^{(3)}) \end{cases}$$

$$(c) \delta_{out} = \sigma'(a_{out}) - (1-t) \\ = \sigma'(w_{11}^{(4)} z_1^{(3)} + w_{12}^{(4)} z_2^{(3)} + w_{10}^{(4)}) - (1-t)$$

$$\delta_{out} = \delta_1^{(4)}$$

$$\delta_1^{(3)} = \sigma'(a_1^{(3)}) \sum_{j=1}^{N_2} \delta_j^{(4)} w_{j1}^{(4)} = \sigma'(a_1^{(3)}) \delta_1^{(4)} w_{11}^{(4)} \\ = \sigma(a_1^{(3)}) (1 - \sigma(a_1^{(3)})) \delta_1^{(4)} w_{11}^{(4)} \\ = \sigma(w_{11}^{(3)} z_1^{(2)} + w_{10}^{(3)}) (1 - \sigma(w_{11}^{(3)} z_1^{(2)} + w_{10}^{(3)})) \delta_1^{(4)} w_{11}^{(4)}$$

$$\delta_1^{(2)} = \sigma'(a_1^{(2)}) \sum_{j=1}^{N_2} \delta_j^{(3)} w_{j1}^{(3)} = \sigma'(a_1^{(2)}) (\delta_1^{(3)} w_{11}^{(3)} + \delta_2^{(3)} w_{21}^{(3)}) \\ = \sigma(a_1^{(2)}) (1 - \sigma(a_1^{(2)})) (\delta_1^{(3)} w_{11}^{(3)} + \delta_2^{(3)} w_{21}^{(3)}) \\ = \sigma(w_{11}^{(2)} \sigma(w_{11}^{(1)} \sigma(x) + w_{10}^{(1)}) + w_{12}^{(2)} \sigma(w_{21}^{(1)} \sigma(x) + w_{20}^{(1)}) + w_{10}^{(2)}) \cdot \\ (1 - \sigma(w_{11}^{(2)} \sigma(w_{11}^{(1)} \sigma(x) + w_{10}^{(1)}) + w_{12}^{(2)} \sigma(w_{21}^{(1)} \sigma(x) + w_{20}^{(1)}) + w_{10}^{(2)})) \cdot \\ (\delta_1^{(3)} w_{11}^{(3)} + \delta_2^{(3)} w_{21}^{(3)})$$

$$\delta_1^{(1)} = \sigma'(a_1^{(1)}) \sum_{j=1}^{N_2} \delta_j^{(2)} w_{j1}^{(2)} = \sigma'(a_1^{(1)}) \delta_1^{(2)} w_{11}^{(2)} \\ = \sigma(a_1^{(1)}) (1 - \sigma(a_1^{(1)})) \delta_1^{(2)} w_{11}^{(2)} \\ = \sigma(w_{11}^{(1)} \sigma(x) + w_{10}^{(1)}) (1 - \sigma(w_{11}^{(1)} \sigma(x) + w_{10}^{(1)})) \delta_1^{(2)} w_{11}^{(2)}$$

$$\frac{\partial L}{\partial w_{11}^{(1)}} = \delta_1^{(1)} z_1^{(0)} = \delta_1^{(1)} \sigma(x)$$

4. (a) $\beta_0 = 3, \beta_1 = 0.1, \beta_2 = 1$

In this case, all blue points have positive value and all red points have negative value

(b) The large coefficients will be shrunk. Specifically, $\sum_{j=1}^p |\beta_j|^2$ will be restricted within a range. As a result, new β will increase the bias and decrease the variance.