Part B

1.



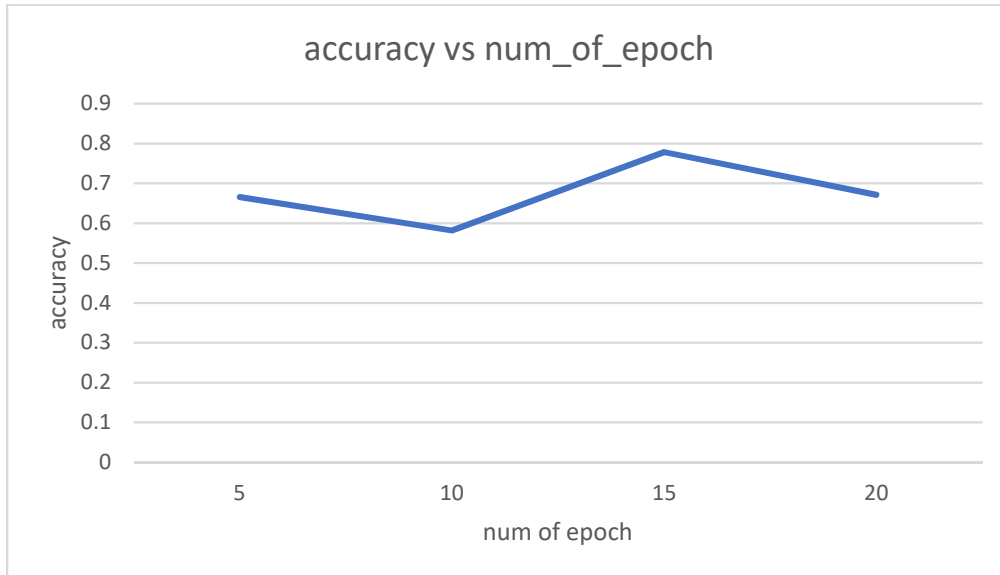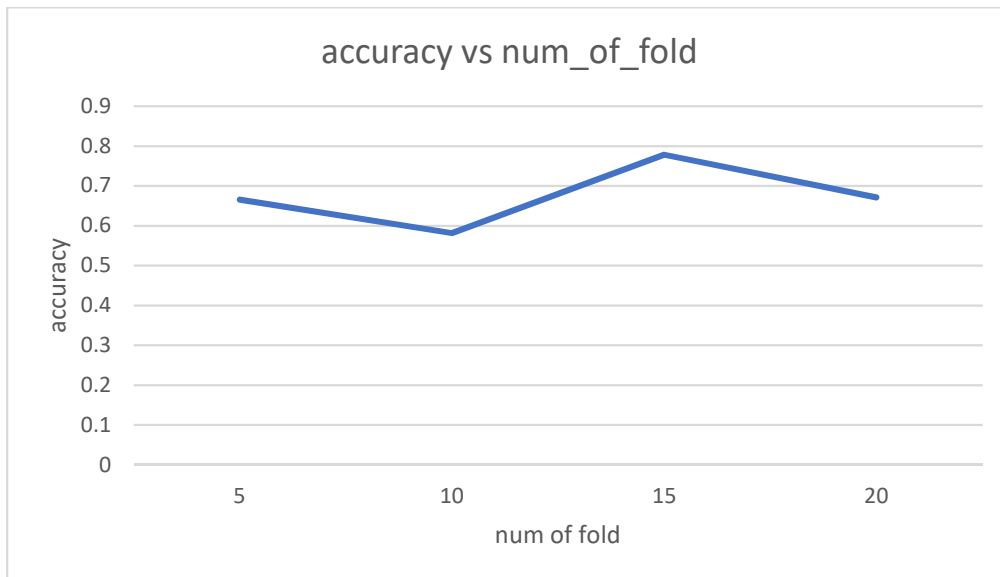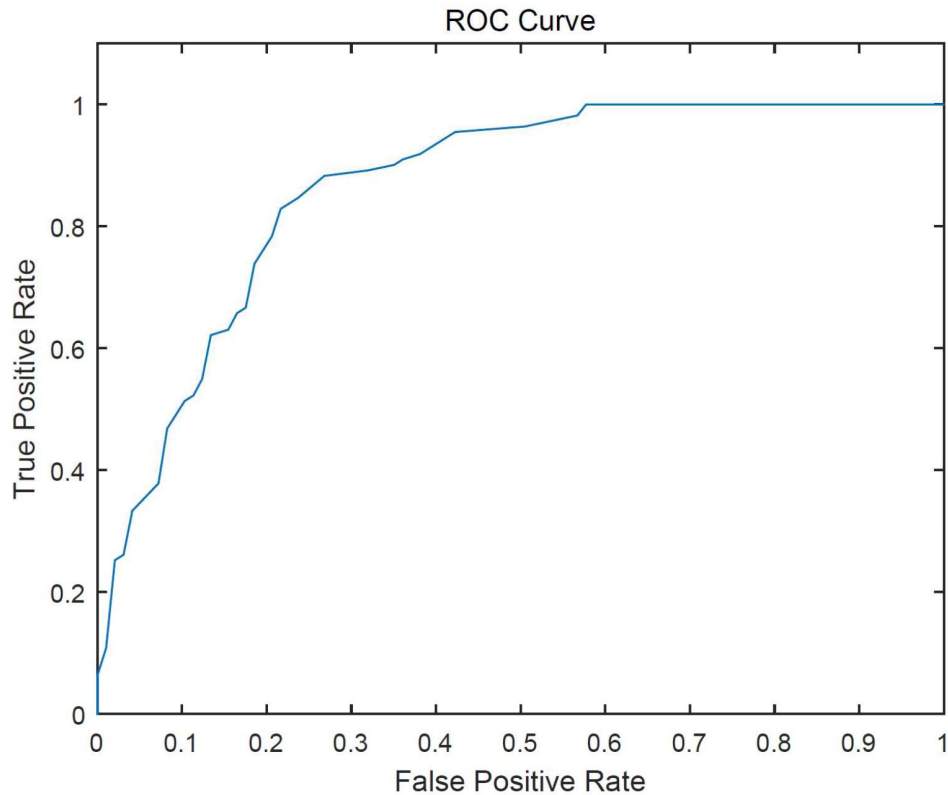accuracy vs num_of_epoch

2.



accuracy vs num_of_fold

3.

ROC Curve

Part C

1. Support : Binary vectors of length K

$x \in \{0, 1\}^K$

$Y \sim \text{Bernoulli}(\emptyset)$

$X_k \sim\sim \text{Bernoulli}(\theta_{k,Y}) \ \forall k \in \{1,\ldots\ldots,K\}$

Model: $p_{\emptyset,\theta}(x, y) = (\emptyset)^y (1 - \emptyset)^{(1-y)}$

$$\emptyset = \frac{\sum_{i=1}^{N} \prod(y^{(i)} = 1)}{N}$$

$$\theta_{k,0} = \frac{\sum_{i=1}^{N} \prod(y^{(i)} = 0 \wedge x_k^{(i)} = 0)}{\sum_{i=1}^{N} \prod(y^{(i)} = 0)}$$

$$\theta_{k,1} = \frac{\sum_{i=1}^{N} \prod(y^{(i)} = 1 \wedge x_k^{(i)} = 1)}{\sum_{i=1}^{N} \prod(y^{(i)} = 1)}$$

$\forall k \in \{1, \ldots\ldots, K\}$

2.

According to the information showed in the problem. We can define

$P(Y = y_k | X) \propto \exp(w_{k0} + \sum_{i=1}^{d} w_{ki} X_i)$ for $k = 1, \ldots, K - 1$

Since all probabilities must sum to 1, we should have

$P(Y = y_k | X) = 1 - \sum_{k=1}^{K} P(Y = y_k | X)$

In binary classification, we define

$$P(Y = y_k|X) = \frac{1}{1 + \sum_{k=1}^{K} \exp(w_{k0} + \sum_{i=1}^{d} w_{ki}X_i)}$$

and for k = 1,......,K-1

$$P(Y = y_k|X) = \frac{\exp(w_{k0} + \sum_{i=1}^{d} w_{ki}X_i)}{1 + \sum_{k=1}^{K} \exp(w_{k0} + \sum_{i=1}^{d} w_{ki}X_i)}$$

Hence, it is clearly that p(y = i|x) is a softmax