**CPT_S 570 HW4**

**Yang Zhang**

**11529139**


**1.** The paper states two points: (1) The discriminative learning has lower asymptotic error than generative learning. (2) The generative learning may approach its asymptotic error faster than discriminative learning.


**2. (a)** According to the paper from question 1, it really depends on what rule to measure the performance. Logistic Regress has lower asymptotic error, but Naïve Bayes approach its asymptotic error faster.

**(b)** Logistic Regress is better, because LR doesn't require the data to be independent labelling, so its result won't be affected by independence. However, the Naïve Bayes do require independence, so the result from Naïve Bayes is unreliable is the data does not satisfy the assumption.


**3.(a) Can we compute P(X) from the learned parameters of a Naive Bayes classifier? Please explain your reasoning.**

Yes, because Naïve Bayes involve computing $P(x|y)$ and $P(y)$. we can compute $P(x)$ from the sum of $P(x|y_i)P(y_i)$


**b. Can we compute P(X) from the learned parameters of a Logistic Regression classifier? Please explain your reasoning.**

No, because Logistic Regression learns $P(y|x)$ directly.


**4. (a)**

Let $\hat{y}_k = P(y = k \mid x) = \dfrac{\exp(W_k \cdot x)}{\sum_{i=1}^{k} \exp(W_i \cdot x)}$

$P(y \mid x) = \prod_{i=1}^{k} \hat{y}_k^{\wedge k}$

$\ell(w) = \sum_{i} \log P(y \mid x) = \sum_{i} \sum_{j=1}^{k} j \log \hat{y}_j$

$\qquad = \sum_{i} [\, y_k \log \hat{y}_k + (1 - y_k) \log(1 - \hat{y}_k)\,] \Rightarrow \boxed{Likelihood}$

The gradient:

$\dfrac{\partial \ell(w)}{\partial W} = \dfrac{\partial}{\partial \omega} [\, y_k \log \hat{y}_k + (1 - y_k) \log(1 - \hat{y}_k)\,]$

$\qquad = \dfrac{y_k}{\hat{y}_k}\left(\dfrac{\partial \hat{y}_k}{y_k}\right) + \dfrac{1 - y_k}{1 - \hat{y}_k}\left(-\dfrac{\partial \hat{y}_k}{y_k}\right)$

$\qquad = \left[\dfrac{y_k - \hat{y}_k}{\hat{y}_k (1 - y_k)}\right] \dfrac{\partial \hat{y}_k}{\partial \omega}$

$\qquad = \left[\dfrac{y_k - \hat{y}_k}{\hat{y}_k (1 - y_k)}\right] \hat{y}_k (1 - \hat{y}_k) x = (y_k - \hat{y}_k) \cdot x$

Therefore $\dfrac{\partial \ell(w)}{\partial W_j} = \sum_{i=1}^{N} (y_k^i - \hat{y}_k^i) x^i$

**(b)**

with adding the regularization term the likelihood become:

$\ell(w) = \sum_{i} \left[\, \sum_{j} y_j^i \log \hat{y}_j^i - \lambda \|w\|^2 \,\right]$

by applying chain rule

$\nabla \ell(w) = \sum_{i} (y_j^i - \hat{y}_j^i) x^i - \lambda w$

**5.** This paper compared 5 tests for the propose of picking the algorithm with better performance. As a result, the test for the difference of two proportions and paired-differences t test (based on random train/test split) should never be used, because those 2 tests are shown Type I error in certain cases. While the 10 fold cross validation test somewhat elevates the Type I error. In contrast, McNemar's tests are shown to have low Type I error. The last test 5x2cv also provides the acceptable Type I error. As the matter of detecting differences between algorithms, the cross-validation test is the most powerful test.

**6.** This paper talked about the overfitting issue of machine learning. In most of the cases, the general objective function of a ML algorithm is to minimize the in sample error. However, by achieving this objective, the noise data has possibly been fitted as well, this is not the general solution we want. Under computing is one of the way to avoid overfitting.

**7.** This paper provided detailed introduction of popular ensemble methods (Adaboost, bootstrap .etc) and states that ensemble methods often perform better than any single classifier. The paper also explained why in 3 aspects: (1) statistical (2) computational (3) representational.

**8.** Result:

```
Training Accuracy: 0.953416149068323
Testing Accuracy: 0.7920792079207921
```

The training accuracy is higher than testing accuracy

**9**

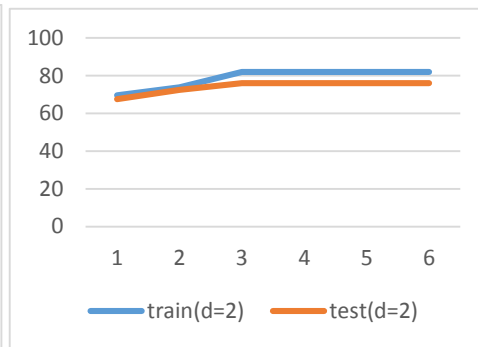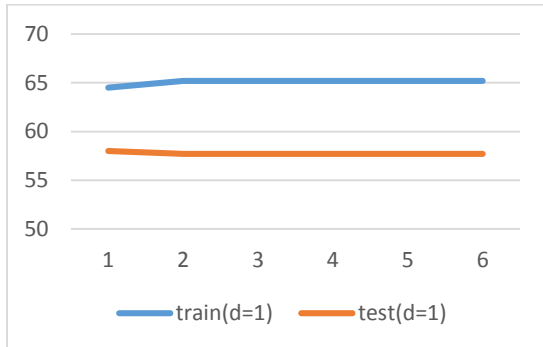**Plots for bagging (average result from two datasets):**

**X-axis mapping {1->10, 2->20, 3->40, 4->60, 5->80, 6->100}**

The accuracy depends on tree depth a lot more than the bag size and for smaller tree depth comes with lower accuracy but higher convergence speed. In contrast, bigger tree depth comes with higher accuracy but slower convergence speed.

**Plots for bagging (average result from two datasets):**

**X-axis mapping {1->10, 2->20, 3->40, 4->60, 5->80, 6->100}**



The accuracy depends on tree depth a lot more than the number of iterations and for smaller tree depth comes with lower accuracy but higher convergence speed. In contrast, bigger tree depth comes with higher accuracy but slower convergence speed.