

# ISYE 6740 HW 4

Mohammed Khalid Siddiqui

September 30, 2024

## 1 Conceptual Questions

1. Derive the gradient of the cost-function:

## Homework 4

Given  $m$  training samples  $(x_i, y_i)$   $i=1 \dots m$ ,

$$x_i \in \mathbb{R}, y_i \in \{0, 1\}$$

we must solve the following to fit logistic regression model for classification

$$\max_{\theta} l(\theta)$$

$$\text{where } l(\theta) = \sum_{i=1}^m \left\{ -\log(1 + \exp\{-\theta^T x_i\}) + (y_i - 1)\theta^T x_i \right\}$$

To find gradient, differentiate  $l(\theta)$  w.r.t.  $\theta$

$$\frac{\partial l(\theta)}{\partial \theta} = \frac{1}{1 + \exp(-\theta^T x_i)} \cdot \exp(-\theta^T x_i) + (y_i - 1)x_i$$

$$\nabla_{\theta} l(\theta) = \sum_{i=1}^m \left( \frac{x_i \exp(-\theta^T x_i)}{1 + \exp(-\theta^T x_i)} + x_i(y_i - 1) \right)$$

$$\text{Sigmoid function: } g(z) = \frac{1}{1 + \exp(-z)}$$

$$\text{therefore } g(\theta^T x_i) = \frac{1}{1 + \exp(-\theta^T x_i)} = 1 - \frac{\exp(-\theta^T x_i)}{1 + \exp(-\theta^T x_i)}$$

$$\text{so that } \nabla_{\theta} l(\theta) = \sum_{i=1}^m ((1 - g(\theta^T x_i))x_i + (y_i - 1)x_i)$$

$$\sum_{i=1}^m ((1 - g(\theta^T x_i)) + (y_i - 1))x_i$$

$$\frac{1 - g(\theta^T x_i)}{1 - g(\theta^T x_i)} = -(1 - g(\theta^T x_i))$$

$$+ g(\theta^T x_i) = - (1 - g(\theta^T x_i))$$

~~$$\text{Therefore, } = \sum_{i=1}^m \left\{ (g(\theta^T x_i) - y_i) \frac{(y_i - g(\theta^T x_i))x_i}{1 - g(\theta^T x_i)} \right\}$$~~

$$= \sum_{i=1}^m \left( y_i - \frac{1}{1 + \exp(-\theta^T x_i)} \right) x_i$$

- Pseudo-code to run gradient descent:

Pseudo-code for gradient descent to find  $\theta^*$

1. Set  $\theta$  to some random value (or 0)

2. Set learning rate  $\gamma$  to some value

3. Set max iterations  $n$  ~~to infinity~~ (if needed)

4. Set a threshold for convergence  $G$

5. For each iteration  $n$  is:

6. Initialize gradient = 0

7. For each training sample  $(x_i, y_i)$ :

8. Compute predicted probability

$$G(\theta^\top x_i) = \frac{1}{1 + \exp(-\theta^\top x_i)}$$

9. Compute gradient for sample

$$\nabla_{\theta_i} l(\theta_i) = (G(\theta^\top x_i) - y_i) \cdot x_i$$

10. Update  $\theta$ :

$$\theta = \sum_{i=1}^{n+1} \nabla_{\theta_i} l(\theta) \approx \theta_n - \gamma \nabla_{\theta} l(\theta)$$

11. Check for convergence ~~out~~:

$$\text{if } |\theta_{n+1} - \theta_n| < G$$

then break

12. Return  $\theta$

3. Pseudo-code to run stochastic gradient descent and the difference between SGD and GD:

### 3. Pseudo-code for Stochastic Gradient Descent

1. Initialize:  $\theta$ ,  $\gamma$ , number of epochs, 6

2. For each epoch:

3. Shuffle training data for randomness

4. Divide data into  $S_n$  samples

5. For each  $s_k$  in  $S_n$ :

6. Calculate For each training sample  $(x_i, y_i)$

7. Compute  $\ell(\theta^T x_i)$

8. Compute  $\nabla_{\theta} \ell(\theta) = (\ell(\theta^T x_i) - y_i) x_i$

9. Update  $\theta_{n+1} = \theta_n - \gamma \cdot \nabla_{\theta} \ell(\theta)$

11. Check for convergence  $|\theta_{n+1} - \theta_n| < \epsilon$

if True, then break

The primary difference between GD and SGD is in SGD, we find the gradient for sample  $S$  to save on computing power, whereas in GD, we use the whole dataset. One of the drawbacks is that path to convergence may be noisier and less accurate for SGD as it uses one sample and sometimes even just one training example.

4. Using the Hessian to prove concavity for logistic regression training problem:

4. Convexity:  $f(\theta_x + (1-\theta)y) \geq \theta f(x) + (1-\theta)f(y)$   
 for all  $0 \leq \theta \leq 1$

$$\text{Given } l(\theta) = \sum_{i=1}^m (-\log(1+\exp(-\theta^T x_i)) + (y_i - 1)\theta^T x_i)$$

- First-Order Derivative (gradient):

$$\nabla_\theta l(\theta) = \sum_{i=1}^m (\theta(\theta^T x_i) - y_i)(y_i - \theta^T x_i) x_i$$

- Second-Order Derivative (Hessian)

↳ matrix of order-2 partial derivatives

$$\begin{aligned} H(\theta) &= \nabla_\theta^2 l(\theta) \\ &= \frac{\partial}{\partial \theta} \sum_{i=1}^m (y_i - \frac{1}{1+\exp(-\theta^T x_i)}) x_i \\ &\stackrel{\text{# (Chain rule)}}{=} \sum_{i=1}^m \left( \frac{1}{(1+\exp(-\theta^T x_i))^2} \exp(-\theta^T x_i) \right) x_i \end{aligned}$$

$$\begin{aligned}
 &= \sum_{i=1}^m \frac{x_i \cdot \exp(-\theta^T x_i) \cdot x}{(1 + \exp(-\theta^T x_i))^2} \\
 &= \sum_{i=1}^m x_i \left( \frac{1}{1 + \exp(-\theta^T x_i)} \cdot \left( \frac{1}{1 + \exp(-\theta^T x_i)} - 1 \right) \right) \cdot x^T \\
 &= \sum_{i=1}^m x_i \cdot (6(\exp(-\theta^T x_i)) \cdot (6(\exp(-\theta^T x_i)) - 1)) \cdot x^T
 \end{aligned}$$

$$0 < \exp(-\theta^T x_i) < 1$$

therefore

$$6(\exp(-\theta^T x_i)) \cdot (6(\exp(-\theta^T x_i)) - 1) \leq 0$$

and of the form

$$H = \nabla^2 l(\theta) = x^T D \cdot x \leq 0$$

- This means  $H$  is ~~semi~~ negative semi-definite indicating all second-derivatives wrt  $\theta$  curves downward or flat.
- Due to concavity,  $l(\theta)$  has only global maximum which means with the right learning rate, gradient descent will always converge.

## 2 Bayes Classifier for spam filtering

1. Displaying all feature vectors and calculating priors:

## Problem 2 HW 4

1. Calculate  $P(y=0)$  and  $P(y=1)$  from training data
- $\hookrightarrow$  spam       $\hookrightarrow$  not spam

- There are 3 spam messages and 4 non-spam messages.

Therefore  $P(y=0) = \frac{3}{3+4}$ ,  $P(y=1) = \frac{4}{3+4}$

$$P(y=0) = \frac{3}{7} \quad P(y=1) = \frac{4}{7}$$

2. Feature Vectors:

$$\mathbb{V} = \{\text{free, money, no, cop, crypto, real, cash, prize, for, you, big, chance, pizza, w, vibe}\}$$

$$\Rightarrow \text{free money no cop} = [1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]$$

$$\Rightarrow \text{crypto real money} = [0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0]$$

$$\Rightarrow \text{real cash money prize for you} =$$

$$[0, 1, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0]$$

$$\Rightarrow \text{big free chance} = [1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0]$$

$$\Rightarrow \text{money, money, money} = [0, 3, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]$$

$$\Rightarrow \text{pizza is big vibe for real} = [0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 1]$$

2. Calculating MLE for Theta at four different combinations of class and word occurrences in respective class:

2. Likelihood of sentence  $w$  its feature vector  $X$  given a

class  $c$  is given by

$$P(x|y=c) = \frac{n!}{x_1! \cdots x_d!} \prod_{k=1}^d \theta_{c,k}^{x_k} \quad c = \{0, 1\}$$

where  $n = x_1 + \cdots + x_d$   $0 \leq \theta_{c,k} \leq 1$  is probability

of word  $k$  appearing in class  $c$

$$\text{Constraint} \rightarrow \sum_{k=1}^d \theta_{c,k} = 1 \quad c = \{0, 1\}$$

$$l(\theta_{0,1}, \dots, \theta_{0,d}, \dots, \theta_{1,1}, \dots, \theta_{1,d}) = \sum_{i=1}^m \sum_{k=1}^d x_k^{(i)} \log \theta_{y_i, k}$$

Calculate MLE  $\theta_{0,1}, \theta_{0,d}, \theta_{1,1}, \theta_{1,d}$

## 2. MLE Calculation

→ Using Lagrange Multiplier

$$L(P(x|y)) = l(\theta) + \lambda \left( \sum_{k=1}^d \theta_{c,k} - 1 \right)$$

$$\Rightarrow L(\theta, \lambda_c) = \sum_{i=1}^m \sum_{k=1}^d x_k^{(i)} \log \theta_{y_i, k} + \lambda_c \left( 1 - \sum_{k=1}^d \theta_{c,k} \right)$$

$$\frac{\partial L}{\partial \theta_{c,k}} = \sum_{i=1}^m \frac{x_k^{(i)}}{\theta_{c,k}} - \lambda_c \quad \leftarrow \text{Set to 0}$$

$$\sum_{i=1}^m \frac{x_k^{(i)}}{\theta_{c,k}} - \lambda_c = 0$$

$$\theta_{c,k} = \frac{1}{\lambda_c} \sum x_k^{(i)}$$

$$\sum_{k=1}^d \theta_{c,k} = 1 \rightarrow \sum_{k=1}^d \pi_c \sum_{i=1}^m x_k^{(i)} = 1$$

$$\pi_c = \sum_{k=1}^d \sum_{i=1}^m x_k^{(i)}$$

MLG:

$$\theta_{c,k} = \frac{\sum_{i=1}^m x_k^{(i)}}{\sum_{k=1}^d \sum_{i=1}^m x_k^{(i)}} \quad \left\{ \text{for } c \in \{0,1\} \right.$$

$c=1 \rightarrow$  non-spam  $c=0 \rightarrow$  spam

- $\theta_{0,1}$   $\leftarrow$  Probability of word "free" appearing in spam messages

$\rightarrow$  free appears once in spam messages

$\rightarrow$  There are 13 word

$$\text{MLG } \theta_{0,1} = \frac{1}{13} = 0.0769$$

- $\theta_{0,6}$   $\leftarrow$  Probability of "real" in spam messages

$\rightarrow$  "real" appears twice in spam messages

$\rightarrow$  13 word occurrences in spam

$$\text{MLG } \theta_{0,6} = \frac{2}{13} = 0.1538$$

- $\theta_{1,2}$   $\leftarrow$  "money" prob. in non-spam messages

$\rightarrow$  money appears 3 times

$\rightarrow$  15 total word occurrences in non-spam

$$\text{MLG } \theta_{1,2} = \frac{3}{15} = \frac{1}{5} = 0.2$$

- $\theta_{1,15}$   $\leftarrow$  Prob. "vibe" in non-spam messages

$\rightarrow$  "vibe" appears once in non-spam

$\rightarrow$  15 total "non-spam" word occurrences

$$\text{MLG } \theta_{1,15} = \frac{1}{15} = 0.0667$$

3. Determining the class using our Naive Bayes implementation for a given test message:

3. Test message: "money for real"

↳ Use NB to calculate posterior and classify as spam or non-spam

$$x_{\text{test}} = [0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0]$$

$$P(y=0 | x_{\text{test}}) \in \text{spam}$$

$$P(y=1 | x_{\text{test}}) \in \text{non-spam}$$

Bayes Theorem:

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

$$P(y=0 | x_{\text{test}}) = \frac{P(x_{\text{test}} | y=0) P(y=0)}{P(x_{\text{test}})}$$

$$P(y=1 | x_{\text{test}}) = \frac{P(x_{\text{test}} | y=1) P(y=1)}{P(x_{\text{test}})}$$

$$P(x_{\text{test}}) = P(x_{\text{test}} | y=0) P(y=0) + P(x_{\text{test}} | y=1) P(y=1)$$

$$P(y=0) = \frac{3}{7} \quad P(y=1) = \frac{4}{7}$$

$$P(x_{\text{test}} | y=0) = P(\text{"money"}) \cdot P(\text{"for"}) \cdot P(\text{"real"})$$

→ Use the MLEs

$$\theta_{0,2} = \frac{3}{13} \quad \theta_{0,6} = \frac{2}{13} \quad \theta_{0,9} = \frac{1}{13}$$

$$P(x_{\text{test}} | y=0) = \frac{6}{2197}$$

$$P(x_{\text{test}} | y=1) = P(\text{"money"}) \cdot P(\text{"for"}) \cdot P(\text{"real"})$$

$$\theta_{1,2} = \frac{1}{5} \quad \theta_{1,6} = \frac{1}{13} \quad \theta_{1,9} = \frac{1}{13}$$

$$P(x_{\text{test}} | y=1) = \frac{1}{3375} \quad ; P(x_{\text{test}}) =$$

- We want to calculate whether  $P(y=0 | x_{\text{test}})$  or

$$P(y=1 | x_{\text{test}}) \gg \text{larger.}$$

- Since  $P(x_{\text{test}})$  is the same denominator for

both, we can choose to ignore it as proportionality relationship will be maintained, or we can calculate posterior.

$$P(y=0|x_{text}) = \frac{\frac{6}{2197} \cdot \frac{3}{7}}{\left(\frac{6}{2197}\right)\left(\frac{3}{7}\right) + \left(\frac{1}{3375}\right)\left(\frac{4}{7}\right)}$$

$$= \frac{\frac{18}{15879}}{\frac{18}{15879} + \frac{4}{23625}}$$

$$\approx 0.8736$$

$$P(y=1|x_{text}) = \frac{\frac{1}{3375} \cdot \frac{4}{7}}{\left(\frac{6}{2197}\right)\left(\frac{3}{7}\right) + \left(\frac{1}{3375}\right)\left(\frac{4}{7}\right)}$$

$$\approx 0.12637$$

Given that the posterior for  $P(y=0|x_{text})$  is greater than  $P(y=1|x_{text})$ , we can conclude that our model predicts the message to be classified as spam.

### 3 Comparing classifiers: Divorce classification/prediction

- For this problem, we predict whether a participant in a personal information form is divorced based on 54 predictors and 1 ground truth label where 1 means divorce and 0 represents not divorced.

The three classifiers employed in this task are the standard sklearn implementations of Naive Bayes, Logistic Regression, and K-Nearest Neighbors.

The success metric for our prediction is the testing accuracy, measured using the accuracy\_score implementation within sklearn.metrics applied on the 20% test-set. For robustness, I applied Monte Carlo Cross Validation to split the dataset and run the algorithms 1000 times returning the mean test accuracy and variance for each classifier as well as the variance in the testing accuracy across the 1000 iterations. The data was also scaled prior to implementation as we are using KNN with a euclidean distance metric.

Below, you'll see the distribution of test accuracy along with the mean accuracy and variance of the accuracy for each of the classifiers.

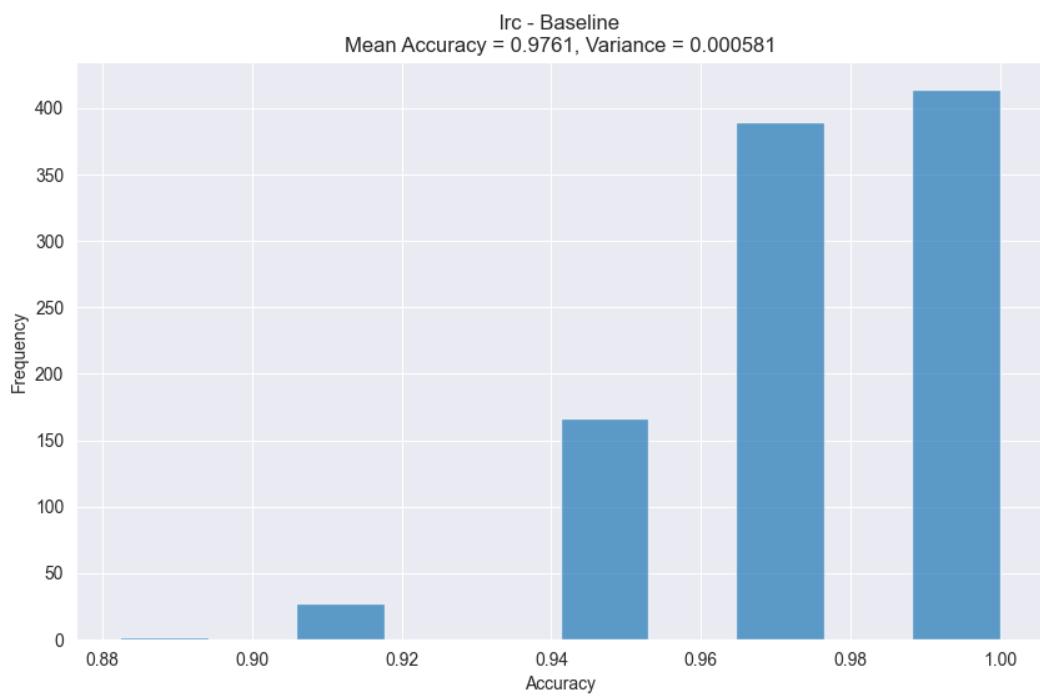


Figure 1: Logistic Regression

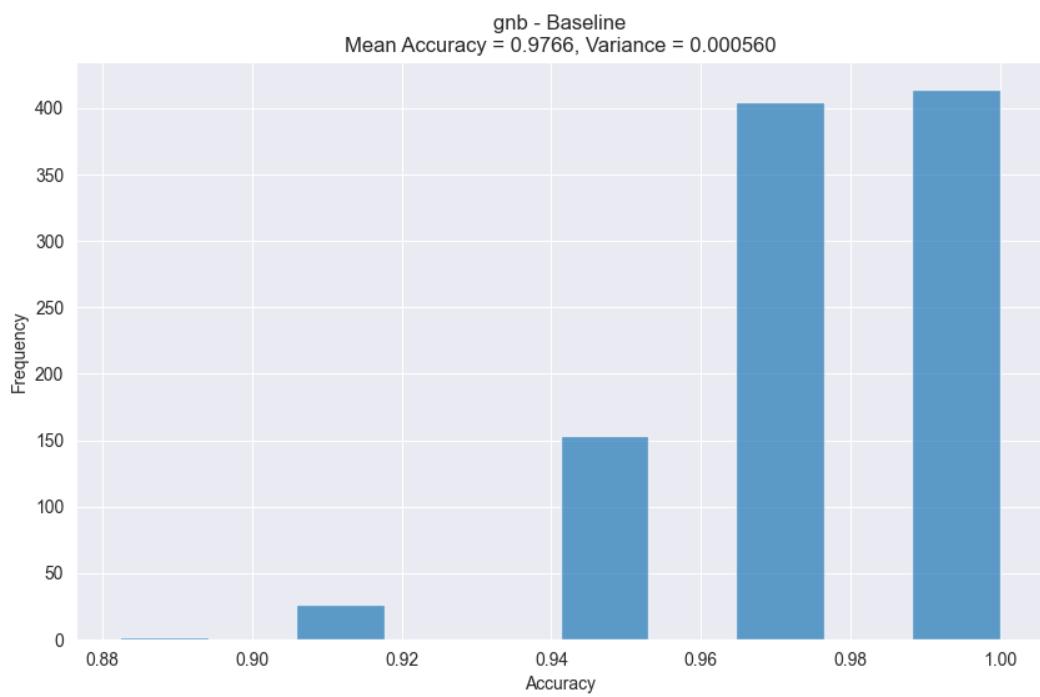


Figure 2: Naive Bayes

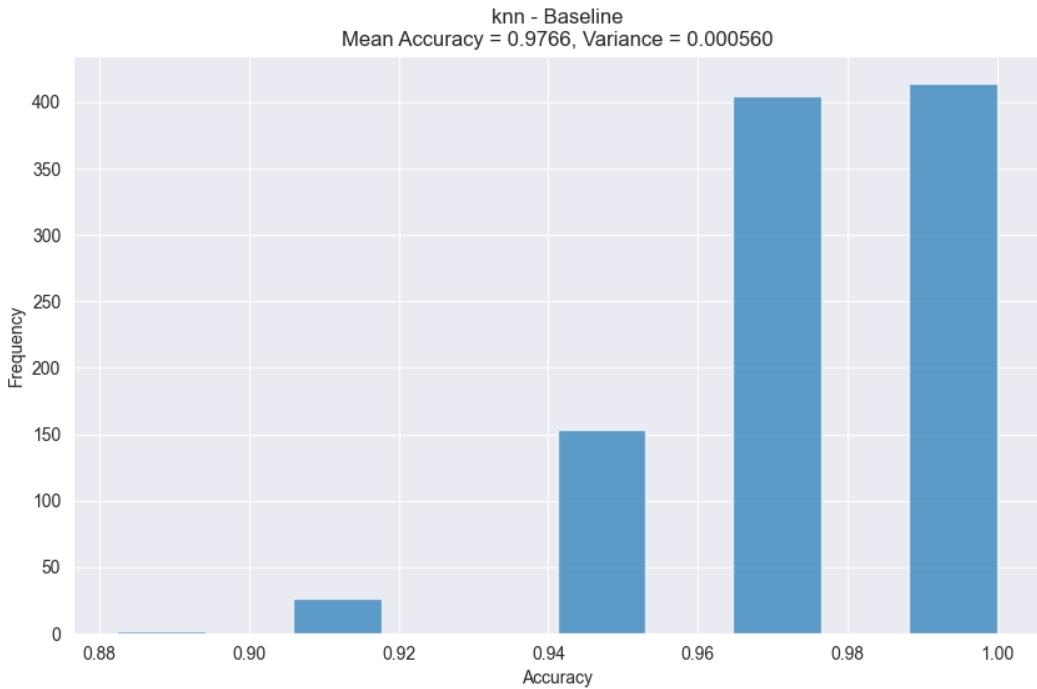


Figure 3: K-Nearest Neighbors

It seems the standard implementation of KNN (with five neighbors and using euclidean distance metric) and Naive Bayes (with a suggested smoothing parameter of 1e-3 to prevent divide by 0 errors arising where there is no variance among predictors) are tied for the mean accuracy and variance performing slightly better than our logistic regression model.

Metric	Logistic Regression	Naive Bayes	KNN
Mean Accuracy	0.9761	0.9766	0.9766
Variance	0.000581	0.000560	0.000560

KNN and Naive Bayes may perform better than Logistic Regression when decision boundaries are non-linear. Naive Bayes performs better when features can be assumed (almost) independent while KNN performs best when local relationships dominate and an appropriate distance metric can be used to map relationships between data points. Naive Bayes and KNN both perform well in low-dimensional, well-separated spaces.

2. The next step in this problem was to use PCA (I used the sklearn implementation with SVD) to reduce the dimension spaces to 2 before running the same Monte Carlo CV and returning the testing accuracy and variance for each.

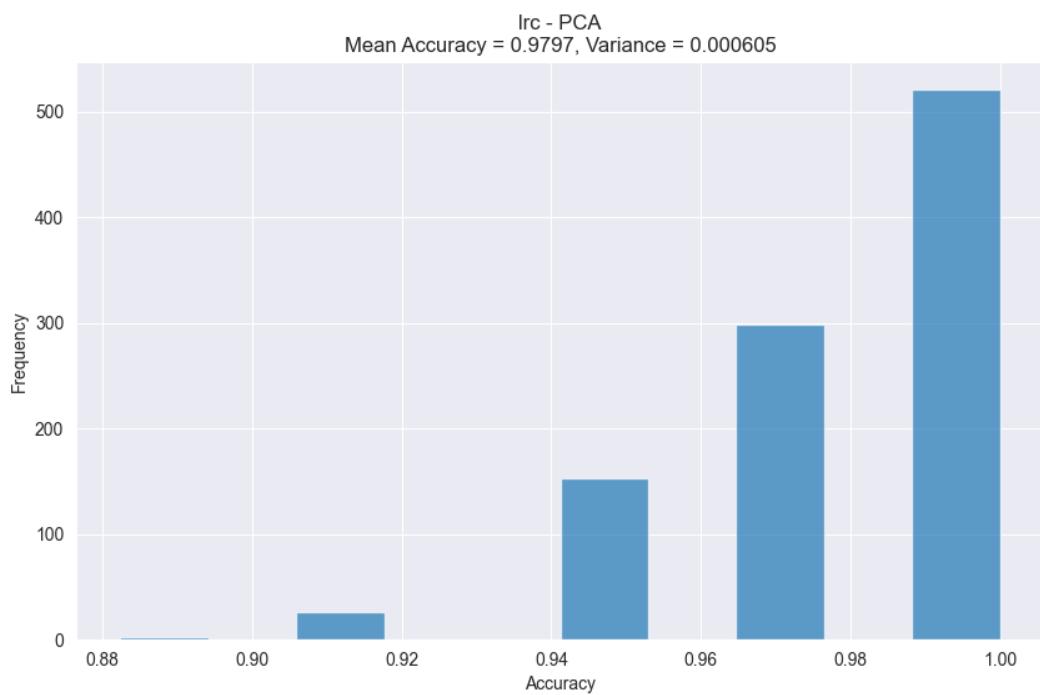


Figure 4: Logistic Regression

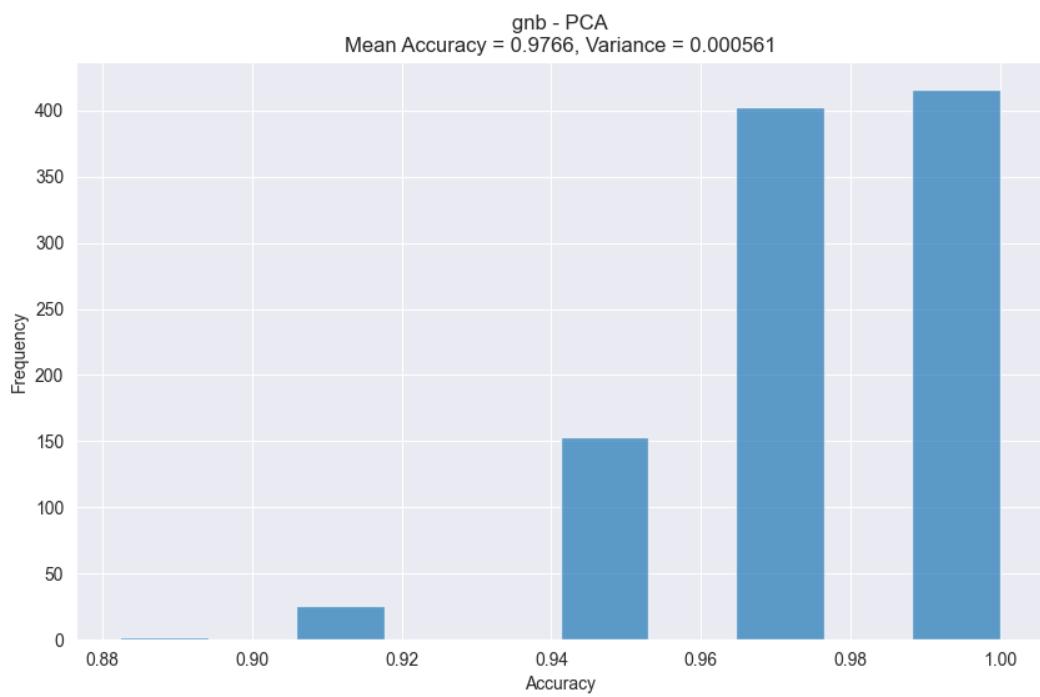


Figure 5: Naive Bayes

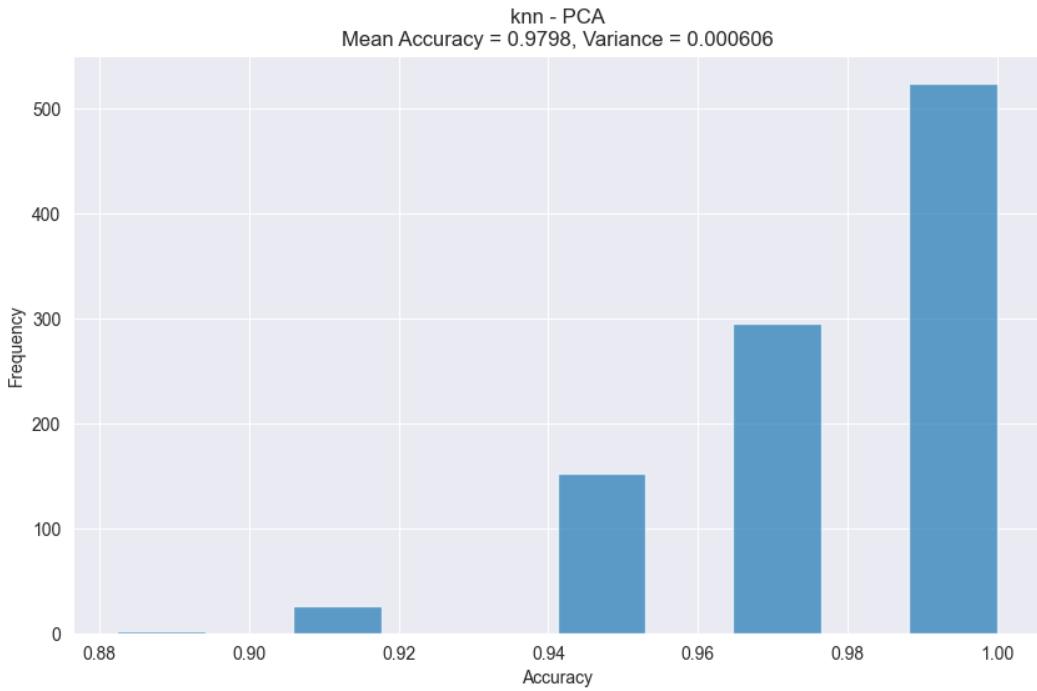


Figure 6: KNN

Metric	Logistic Regression	Naive Bayes	KNN
Mean Accuracy	0.9797	0.9766	0.9798
Variance	0.000605	0.000561	0.000606

After reducing the dimension spaces to 2, it appears KNN has the highest mean accuracy with prediction for our dataset while Naive Bayes has the lowest variance, being more sure of its predictions over different data splits. The decision boundaries and data points are plotted below through the use of an sklearn implementation to display model decision boundaries.

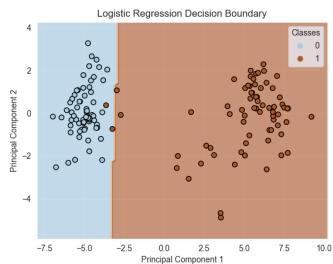


Figure 7: Logistic Regression

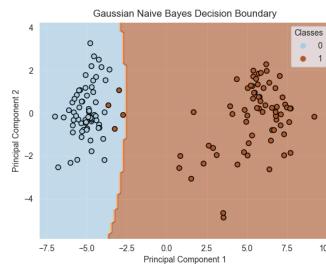


Figure 8: Naive Bayes

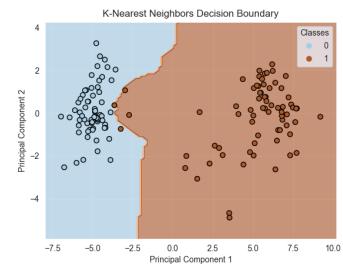


Figure 9: K-Nearest Neighbors

The decision for logistic regression appears to be a straight line indicating a linear decision boundary. Within 2 dimensions, this is a line, and in higher dimensions, a hyperplane. The decision boundary for Naive Bayes appears to be curved as it assumes Gaussian distribution and calculated the probability

of each data point belonging to a class. K-Nearest Neighbors is non-parametric and uses the distance between data points to classify them. A majority vote of the k nearest neighbors (in this case 5) determines which class a given point belongs to. As a result, the decision boundary for KNN can be quite irregular, sometimes even fragmented.

## Sources

1. <https://scikit-learn.org/dev/modules/generated/sklearn.neighbors.KNeighborsClassifier.html>
2. <https://scikit-learn.org/dev/modules/generated/sklearn.inspection.DecisionBoundaryDisplay.html>
3. [https://scikit-learn.org/dev/modules/generated/sklearn.naive\\_bayes.GaussianNB.html](https://scikit-learn.org/dev/modules/generated/sklearn.naive_bayes.GaussianNB.html)
4. [https://scikit-learn.org/1.5/modules/generated/sklearn.linear\\_model.LogisticRegression.html](https://scikit-learn.org/1.5/modules/generated/sklearn.linear_model.LogisticRegression.html)