# Homework 2, STAT 5241

Zongyi Liu

Mar 1, 2025

# 1 Questions

## 1.1 Question 1

To address the challenge of predicting autism in a European population using a model trained on US data, the team must carefully design their machine learning process to account for the significant differences in key features such as ethnicity, demographics, and environmental factors. The first step involves thorough data preprocessing to align and normalize features across the two datasets, ensuring consistency and handling missing values through imputation. Feature engineering should focus on creating universal patterns rather than region-specific ones. Given the distribution shift between the training and testing datasets, the team should employ domain adaptation techniques such as transfer learning, domain-adversarial neural networks (DANN), or covariate shift correction to improve generalization. Robust models like random forests, gradient boosting machines (e.g., XGBoost, LightGBM), or regularized linear models are recommended due to their ability to handle mixed data types and resist overfitting to region-specific patterns. Ensemble methods can further enhance robustness by combining multiple models.

Hyperparameter tuning should be conducted using k-fold cross-validation on the US dataset, with careful attention to ensuring the validation splits are representative. If a small subset of European data is available, it can be used for fine-tuning to improve generalization. Bayesian optimization or grid search can efficiently explore the hyperparameter space, focusing on key parameters such as regularization terms, learning rates, and model complexity. For evaluation, the primary metric should be the F1 score, which balances precision and recall and is particularly suitable for imbalanced datasets. Secondary metrics like AUC-ROC, precision, recall, and confusion matrices provide additional insights into model performance, while calibration metrics such as the Brier score ensure probabilistic predictions are reliable.

During testing, the model should be evaluated on the European dataset without further tuning, and performance metrics should be reported alongside cross-validation results from the US dataset. Error analysis is crucial to identify misclassified cases and understand whether errors stem from feature distribution differences or other factors. Additionally, the team must consider bias and fairness, ensuring the model does not disproportionately misclassify specific subgroups, and prioritize interpretability using methods like SHAP or LIME to align predictions with clinical knowledge. Continuous monitoring of the model's performance post-deployment is also recommended to detect any data drift over time. By integrating these steps, the team can develop a model that generalizes effectively to the European population while maintaining robustness, fairness, and transparency.

## 1.2 Question 2

### 1.2.1 Part A

We need to pre-process the data. There are columns with many missing values, which should be considered to drop first. Then we can drop rows with missing values. If we simply drop rows with missing data, or drop the rows with missing data before deleting columns with significant number of missing data, it will distorted the regression results.

The result I got using self-coded $K$-fold CV from scratch is as below:

```
Best RBF Kernel Parameters: lambda=0.01, gamma=0.1
Best Polynomial Kernel Parameters: lambda=0.01, degree=2
RBF Kernel Test MSE: 0.014702546829905011
Polynomial Kernel Test MSE: 0.009525269150252024
```

### 1.2.2 Part B

And comparing the results of scratch code and `sklearn` code, the hyperparameters are the same:

```
        Custom RBF Kernel Parameters: lambda=0.01, gamma=0.1
        Custom Polynomial Kernel Parameters: lambda=0.01, degree=2
        scikit-learn RBF Kernel Parameters: {'alpha': 0.01, 'gamma': 0.1}
        scikit-learn Polynomial Kernel Parameters: {'alpha': 0.01, 'degree': 2}
```
The tuned parameters are the same for the built-in model of `sklearn` and the cross-validation I coded up by myself. Plugging those hyper-parameters back, we would get the same MSE as we did in Part A.

### 1.2.3 Part C

I redid the OLS model in Homework 1, and calculated the MSE using code below

```
1    residuals = y_new - y_pred
2
3    # Calculate MSE
4    mse = np.mean(residuals**2)
```

The result is as below:

<div align="center">Mean Squared Error (MSE): 0.0165</div>

Here the non-linear model performed better than the linear model. First of all, linear model is a simple way to fit the data, which is highly likely to be influenced by the noise in the dataset, even though it might has lower MSE. It might cause the problem of overfitting. Whereas the CV can split the dataset into several subsets, and evaluate the model on multiple test and training sets, which would make it much generalizable, and overcome the problem of overfitting. And as we mentioned before, linear models just build the model based on data in the dataset, if the training set is not representative, it might not be a good model; but the CV is more generalizable, giving it ability to comprehensively reflect model performance.

Secondly, CV has hyper-parameters tuning, from which the optimal hyper-parameters can be selected, improving model performance.

Another aspect is that OLS lacks a regularization mechanism, making it susceptible to issues like high-dimensional data or multicollinearity. Where as we implemented regularization here, which helps control model complexity, prevent overfitting, and improve generalization. As for the efficiency to use data, the OLS just use one set of data, whereas CV divides the data multiple times, ensuring that all data points are used for both training and validation.

Finally, if we simply compare the MSE, which is the prediction errors given by different models, the MSE of cross-validated models are still smaller than those of the OLS model.

## 1.3 Question 3

### 1.3.1 Part A

Here I listed the confusion matrices, test accuracies given by python for 5 methods below. I did not generate the confusion matrix picture due to the restriction of page-count. Full reports, including `recall` rate, `f1-score`, etc, were inserted in the Appendix by me.

Firstly, for Logistic Regression in One-vs-Rest setting, I got the report
```
        Test Accuracy: 93.54%
        Confusion Matrix:
        [[1328   39   52]
         [  54 1195   45]
         [  39   50 1254]]
```
For Multinomial Regression
```
        Test Accuracy: 93.12%
        Confusion Matrix:
        [[1324   43   52]
         [  51 1201   42]
         [  39   52 1252]]
```
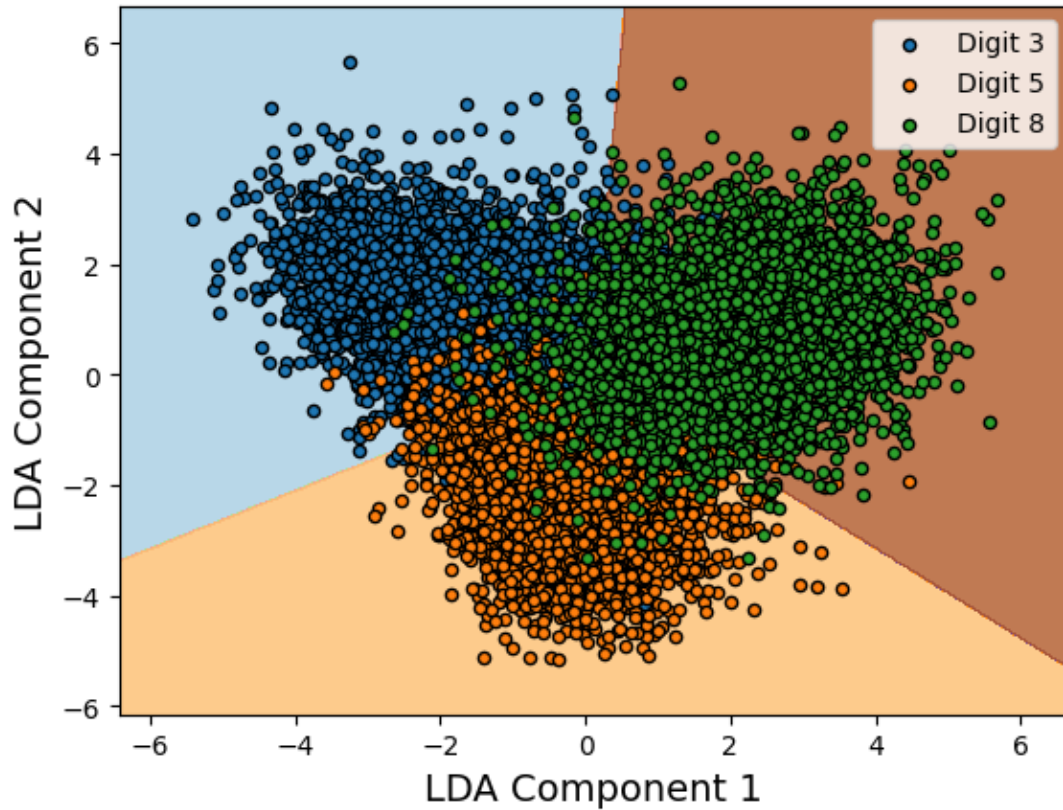For Naive Bayes
```
        Test Accuracy: 50.22%
        Confusion Matrix:
        [[ 595   26  798]
         [  95  138 1061]
         [  18   21 1304]]
```
For Linear Discriminant Analysis

```
Test Accuracy (LDA): 91.69%
Confusion Matrix (LDA):
[[1285   69   65]
 [  47 1202   45]
 [  34   77 1232]]
```
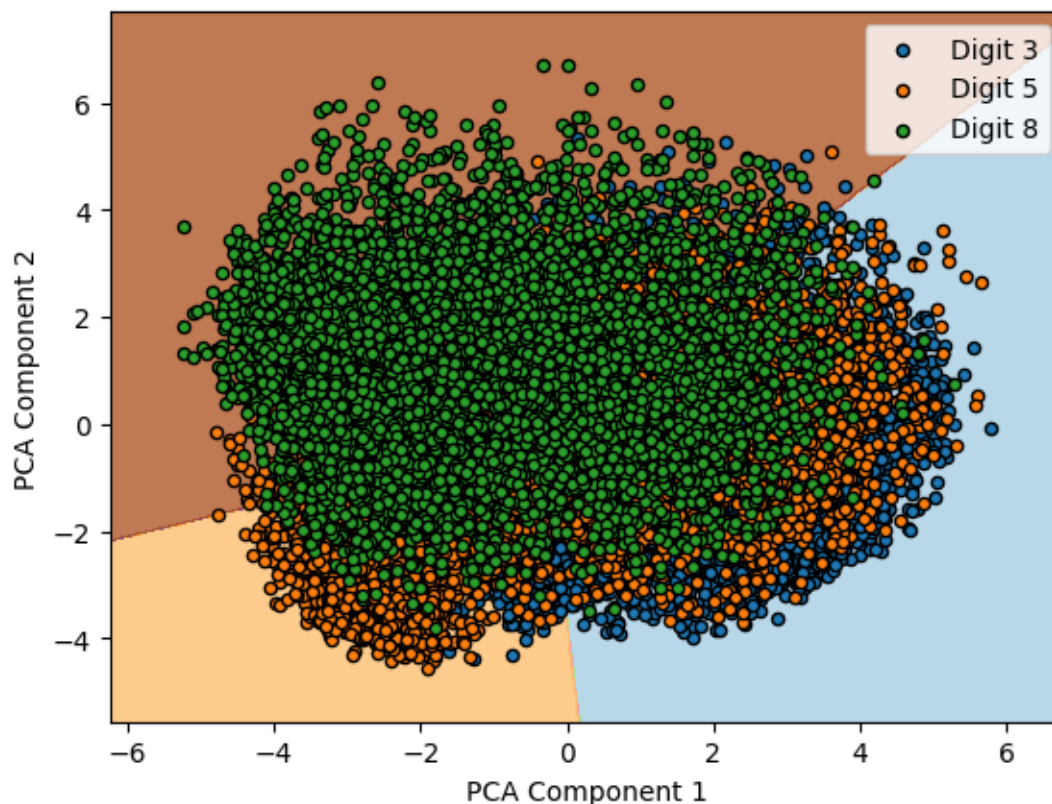
Figure 1: Decision Boundary after LDA



Finally, for Linear Support Vector Machine:

```
Test Accuracy (Linear SVM): 92.73%
Confusion Matrix (Linear SVM):
[[1319   44   56]
 [  57 1193   44]
 [  39   55 1249]]
```

If we want to plot the decision boundary of Linear SVM, we must first do a Principle Component Analysis (PCA).

Figure 2: Decision Boundary after PCA

Overall, if we directly compare the test accuracy evaluated by python, the logistic regression has the highest accuracy rate, whereas the Naive Bayes has the least accuracy.

In the python built-in function `confusion matrix()`, the vertical axis represents the true values, whereas the horizontal axis represents the predictive value. Thus we can added the misclassified terms accordingly. However, I noticed that there was a large proportion of data misclassified in Naive Bayes case, so I splitted them into two parts, firstly I calculated the total misinterpreted values for 4 methods other than Naive Bayes; here the most often misclassified digit is 3, and has 420 counts. And for 5 and 8, the numbers are equal, which is 385.

If we add up the NB method, overall, digit 5 is more likely to be misclassified, and this trend towards error reaches to peak when we use Naive Bayes method. The total number of misclassified 5 is 1542, and 3 is 1244, and 8 is 424.
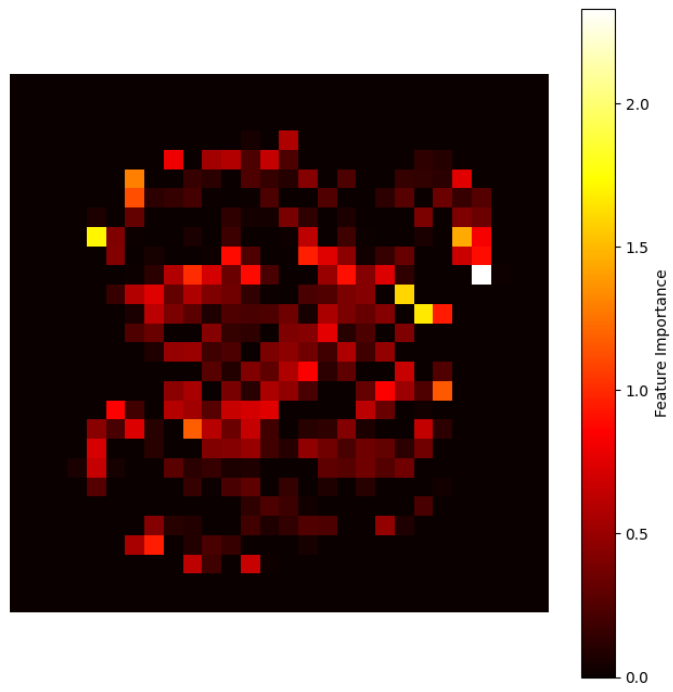
Figure 3: Count of Misclassified Digits

| | A | B | C | D | E | F | G | H | I |
|---|---|---|---|---|---|---|---|---|---|
| 1 | | miclassified | Logistic | Multinomial | LDA | LSVM | Total, except NB | Naïve Bayes | Total, with NB |
| 2 | | mis-classified as 5 | 39 | 43 | 69 | 44 | 195 | 26 | 221 |
| 3 | 3 | mis-classified as 8 | 52 | 52 | 65 | 56 | 225 | 798 | 1023 |
| 4 | | total | 91 | 95 | 134 | 100 | 420 | 824 | 1244 |
| 5 | | mis-classified as 3 | 54 | 51 | 47 | 57 | 209 | 96 | 305 |
| 6 | 5 | mis-classified as 8 | 45 | 42 | 45 | 44 | 176 | 1061 | 1237 |
| 7 | | total | 99 | 93 | 92 | 101 | 385 | 1157 | 1542 |
| 8 | | mis-classified as 3 | 39 | 39 | 34 | 39 | 151 | 18 | 169 |
| 9 | 8 | mis-classified as 5 | 50 | 52 | 77 | 55 | 234 | 21 | 255 |
| 10 | | total | 89 | 91 | 111 | 94 | 385 | 39 | 424 |
| 11 | | | | | | | | | |

### 1.3.2 Part B

The group-lasso regularized multinomial logistic regression will help us to select features that separate the three digits, and we can end up getting a heatmap. In the heatmap, bright regions represent pixels that are highly important for distinguishing between the digits 3, 5, and 8. The model relies heavily on these pixels for making predictions.

And dark regions represent pixels that contribute little to the classification task. The model considers these pixels unimportant.

Figure 4: Headmap showing Features that Separate the Three Digits

5

# 2  Appendix

Supplementary  Github Repo

## 2.1  Question 2

### 2.1.1  Setup

```python
1      # Same as Homework 1
2
3      from ucimlrepo import fetch_ucirepo
4
5
6      # fetch dataset
7      communities_and_crime = fetch_ucirepo(id=183)
8
9      # data (as pandas dataframes)
10     X = communities_and_crime.data.features
11     y = communities_and_crime.data.targets
12
13     # metadata
14     print(communities_and_crime.metadata)
15
16     # variable information
17     print(communities_and_crime.variables)
18
19     # Inspect the shape of X and y
20     print(X.shape)  # Should be (1994, 127)
21     print(y.shape)  # Should be (1994, 1)
22
23     # Check for missing values
24     print(X.isnull().sum())  # Count of missing values per feature
25
26     # Inspect the first few rows of X and y
27     print(X.head())
28     print(y.head())
29     X = X.iloc[:, 5:]
30     print(X.dtypes) # There are object columns within the data. The object data type is
           the default type for columns containing text (strings) in a pandas DataFrame.
```

Drop columns and rows with missing values:

```python
1      # Convert all values to numeric (force non-numeric to NaN)
2      X = X.applymap(pd.to_numeric, errors='coerce')
3
4      # Replace "?" with NaN
5      X.replace("?", np.nan, inplace=True)
6
7      # Check the number of missing values in each column
8      missing_counts = X.isnull().sum()
9
10     # Determine threshold for column removal (e.g., remove if >50% missing)
11     threshold = 0.5 * len(X)  # Adjust this threshold as needed
12     cols_to_drop = missing_counts[missing_counts > threshold].index
13
14     # Drop those columns
15     X_cleaned = X.drop(columns=cols_to_drop)
16
17     X_with_y = pd.concat([X_cleaned, y], axis=1)
18
19     # Remove rows containing missing values in the combined DataFrame
20     X_with_y_cleaned = X_with_y.dropna()
21
```

```
22        X_with_y_cleaned
23
24        # Separate X and y after cleaning
25        X_final = X_with_y_cleaned.drop(columns=['ViolentCrimesPerPop'])
26        y_final = X_with_y_cleaned['ViolentCrimesPerPop']
```

Standardize it:

```
1         Y = y_final.to_numpy()
2
3         # Standardize X (centered and scaled)
4         scaler = StandardScaler(with_mean=True, with_std=True)
5         X_standardized = scaler.fit_transform(X_final)
```

### 2.1.2  Split into Different Sets

```
1         # Set a random seed for reproducibility
2         rng = default_rng(1)
3
4         # Define the proportions for the split
5         train_prop = 0.6
6         validation_prop = 0.2
7         test_prop = 0.2
8
9         # Calculate the number of observations for each split
10        total_samples = X_with_y_cleaned.shape[0]
11        train_size = int(train_prop * total_samples)
12        validation_size = int(validation_prop * total_samples)
13        test_size = total_samples - train_size - validation_size
14
15        # Create a random permutation of row indices
16        indices = rng.choice(np.arange(total_samples), size=(total_samples), replace=False)
17
18        # Split the dataset into train, validation, and test sets
19        y_train = Y[indices[:train_size]]
20        y_val = Y[indices[(train_size + 1):(train_size + validation_size)]]
21        y_test = Y[indices[(train_size + validation_size + 1):]]
22
23        X_train = X_standardized[indices[:train_size]]
24        X_val = X_standardized[indices[(train_size + 1):(train_size + validation_size)]]
25        X_test = X_standardized[indices[(train_size + validation_size + 1):]]
```

### 2.1.3  Using K-fold CV from scratch to tune Kernel Ridge Regularization

```
1         import numpy as np
2         from sklearn.metrics import mean_squared_error
3
4         # Define the RBF kernel function
5         def rbf_kernel(X1, X2, gamma):
6         pairwise_dists = np.sum(X1**2, axis=1).reshape(-1, 1) + np.sum(X2**2, axis=1) - 2 *
              np.dot(X1, X2.T)
7         return np.exp(-gamma * pairwise_dists)
8
9         # Define the Polynomial kernel function
10        def polynomial_kernel(X1, X2, degree, c=1):
11        return (np.dot(X1, X2.T) + c) ** degree
12
13        # Kernel Ridge Regression
14        def kernel_ridge_regression(X_train, y_train, X_test, kernel, lambda_, **
              kernel_params):
15        """
```

```
16        Perform Kernel Ridge Regression.
17
18        Parameters:
19        X_train: Training data (n_samples, n_features).
20        y_train: Target values (n_samples,).
21        X_test: Test data (n_samples_test, n_features).
22        kernel: Kernel function (rbf_kernel or polynomial_kernel).
23        lambda_: Regularization parameter.
24        **kernel_params: Kernel-specific parameters (e.g., gamma for RBF, degree for
              Polynomial).
25
26        Returns:
27        y_pred: Predicted values for X_test.
28        """
29        # Compute the kernel matrix
30        K_train = kernel(X_train, X_train, **kernel_params)
31        K_test = kernel(X_test, X_train, **kernel_params)
32
33        # Solve for the dual coefficients alpha
34        n_samples = X_train.shape[0]
35        alpha = np.linalg.inv(K_train + lambda_ * np.eye(n_samples)) @ y_train
36
37        # Predict on the test set
38        y_pred = K_test @ alpha
39        return y_pred
40
41        # K-Fold Cross-Validation
42        def k_fold_cross_validation(X_standardized, Y, k, kernel, lambda_range, gamma_range=
              None, degree_range=None):
43        """
44        Perform K-fold cross-validation to select the best hyperparameters.
45
46        Parameters:
47        X: Input data (n_samples, n_features).
48        y: Target values (n_samples,).
49        k: Number of folds.
50        kernel: Kernel function (rbf_kernel or polynomial_kernel).
51        lambda_range: List of regularization parameters to try.
52        gamma_range: List of gamma values for RBF kernel (optional).
53        degree_range: List of degree values for Polynomial kernel (optional).
54
55        Returns:
56        best_lambda: Best regularization parameter.
57        best_gamma: Best gamma value (for RBF kernel).
58        best_degree: Best degree value (for Polynomial kernel).
59        """
60        fold_size = len(X_standardized) // k
61        best_lambda = None
62        best_gamma = None
63        best_degree = None
64        best_score = float('inf')
65
66        # Grid search over hyperparameters
67        for lambda_ in lambda_range:
68        if kernel == rbf_kernel:
69        # RBF kernel: only gamma is needed
70        for gamma in (gamma_range if gamma_range is not None else [1.0]):
71        scores = []
72        for i in range(k):
73        # Split into training and validation sets
74        val_indices = range(i * fold_size, (i + 1) * fold_size)
75        train_indices = np.setdiff1d(range(len(X)), val_indices)
76
```

```python
        X_train, y_train = X[train_indices], y[train_indices]
        X_val, y_val = X[val_indices], y[val_indices]

        # Train and predict
        y_pred = kernel_ridge_regression(X_train, y_train, X_val, kernel, lambda_, gamma=
            gamma)

        # Compute the validation score (e.g., mean squared error)
        score = mean_squared_error(y_val, y_pred)
        scores.append(score)

    # Average score across folds
    avg_score = np.mean(scores)
    if avg_score < best_score:
    best_score = avg_score
    best_lambda = lambda_
    best_gamma = gamma

    elif kernel == polynomial_kernel:
    # Polynomial kernel: degree and optionally c are needed
    for degree in (degree_range if degree_range is not None else [2]):
    scores = []
    for i in range(k):
    # Split into training and validation sets
    val_indices = range(i * fold_size, (i + 1) * fold_size)
    train_indices = np.setdiff1d(range(len(X)), val_indices)

    X_train, y_train = X[train_indices], y[train_indices]
    X_val, y_val = X[val_indices], y[val_indices]

    # Train and predict
    y_pred = kernel_ridge_regression(X_train, y_train, X_val, kernel, lambda_, degree=
        degree)

    # Compute the validation score (e.g., mean squared error)
    score = mean_squared_error(y_val, y_pred)
    scores.append(score)

    # Average score across folds
    avg_score = np.mean(scores)
    if avg_score < best_score:
    best_score = avg_score
    best_lambda = lambda_
    best_degree = degree

    return best_lambda, best_gamma, best_degree

    # Example usage
    if __name__ == "__main__":
    # Generate synthetic data
    np.random.seed(42)
    X = np.random.rand(100, 2)  # 100 samples, 2 features
    y = 3 * X[:, 0] + 5 * X[:, 1] + np.random.randn(100) * 0.1  # Linear relationship
        with noise

    # Split data into train, validation, and test sets
    train_size = int(0.6 * len(X))
    val_size = int(0.2 * len(X))
    test_size = len(X) - train_size - val_size

    indices = np.random.permutation(len(X))
    X_train, y_train = X[indices[:train_size]], y[indices[:train_size]]
```

```
136    X_val, y_val = X[indices[train_size:train_size + val_size]], y[indices[train_size:
           train_size + val_size]]
137    X_test, y_test = X[indices[train_size + val_size:]], y[indices[train_size + val_size
           :]]
138
139    # Define hyperparameter ranges
140    lambda_range = [0.01, 0.1, 1, 10]
141    gamma_range = [0.01, 0.1, 1]  # For RBF kernel
142    degree_range = [2, 3, 4, 5, 6, 7]      # For Polynomial kernel
143
144    # Perform K-fold cross-validation for RBF kernel
145    best_lambda_rbf, best_gamma_rbf, _ = k_fold_cross_validation(
146    X_train, y_train, k=5, kernel=rbf_kernel, lambda_range=lambda_range, gamma_range=
           gamma_range
147    )
148    print(f"Best RBF Kernel Parameters: lambda={best_lambda_rbf}, gamma={best_gamma_rbf}"
           )
149
150    # Perform K-fold cross-validation for Polynomial kernel
151    best_lambda_poly, _, best_degree_poly = k_fold_cross_validation(
152    X_train, y_train, k=5, kernel=polynomial_kernel, lambda_range=lambda_range,
           degree_range=degree_range
153    )
154    print(f"Best Polynomial Kernel Parameters: lambda={best_lambda_poly}, degree={
           best_degree_poly}")
155
156    # Evaluate on the test set with the best parameters
157    y_pred_rbf = kernel_ridge_regression(X_train, y_train, X_test, rbf_kernel,
           best_lambda_rbf, gamma=best_gamma_rbf)
158    y_pred_poly = kernel_ridge_regression(X_train, y_train, X_test, polynomial_kernel,
           best_lambda_poly, degree=best_degree_poly)
159
160    print(f"RBF Kernel Test MSE: {mean_squared_error(y_test, y_pred_rbf)}")
161    print(f"Polynomial Kernel Test MSE: {mean_squared_error(y_test, y_pred_poly)}")
```

### 2.1.4 Redo the OLS as in Homework 1

```
1     # Step 1: Separate y and X
2     X_new = X_final
3     y_new = y_final
4
5     # Step 2: Add a constant to X (for the intercept term)
6     X_new = sm.add_constant(X_new)
7
8     # Step 3: Fit the OLS model
9     model = sm.OLS(y_new, X_new)
10
11    results = model.fit()
12
13    # Step 4: View the results
14    print(results.summary())
```

```
                          OLS Regression Results
================================================================================
Dep. Variable:     ViolentCrimesPerPop   R-squared:                       0.696
Model:                            OLS    Adj. R-squared:                  0.680
Method:                 Least Squares    F-statistic:                     43.26
Date:                Fri, 07 Mar 2025    Prob (F-statistic):               0.00
Time:                        13:15:31    Log-Likelihood:                 1261.1
No. Observations:                1993    AIC:                            -2320.
Df Residuals:                    1892    BIC:                            -1755.
Df Model:                         100
Covariance Type:            nonrobust
================================================================================
          coef     std err         t      P>|t|      [0.025      0.975]
--------------------------------------------------------------------------------
const          0.5504      0.203      2.712      0.007       0.152       0.948
population     0.1840      0.397      0.463      0.643      -0.595       0.963
householdsize -0.0223      0.086     -0.259      0.796      -0.191       0.147
racepctblack   0.2049      0.051      4.008      0.000       0.105       0.305
racePctWhite  -0.0492      0.059     -0.837      0.403      -0.164       0.066
racePctAsian  -0.0144      0.034     -0.420      0.674      -0.082       0.053
racePctHisp    0.0609      0.053      1.139      0.255      -0.044       0.166
agePct12t21    0.1104      0.106      1.043      0.297      -0.097       0.318
agePct12t29   -0.2292      0.156     -1.467      0.143      -0.536       0.077
agePct16t24   -0.1302      0.164     -0.793      0.428      -0.452       0.192
agePct65up     0.0497      0.103      0.481      0.630      -0.153       0.253
numbUrban     -0.2964      0.387     -0.766      0.444      -1.055       0.462
pctUrban       0.0467      0.016      2.989      0.003       0.016       0.077
medIncome     -0.1998      0.173     -1.158      0.247      -0.538       0.139
pctWWage      -0.2016      0.089     -2.259      0.024      -0.377      -0.027
pctWFarmSelf   0.0488      0.020      2.422      0.016       0.009       0.088
pctWInvInc    -0.1731      0.068     -2.563      0.010      -0.306      -0.041
pctWSocSec     0.0762      0.107      0.712      0.477      -0.134       0.286
pctWPubAsst    0.0050      0.046      0.108      0.914      -0.085       0.095
pctWRetire    -0.0900      0.037     -2.445      0.015      -0.162      -0.018
medFamInc      0.2880      0.160      1.797      0.073      -0.026       0.602
perCapInc      0.0955      0.189      0.506      0.613      -0.274       0.465
whitePerCap   -0.3510      0.152     -2.303      0.021      -0.650      -0.052
blackPerCap   -0.0288      0.025     -1.131      0.258      -0.079       0.021
indianPerCap  -0.0357      0.019     -1.841      0.066      -0.074       0.002
AsianPerCap    0.0216      0.019      1.145      0.252      -0.015       0.059
OtherPerCap    0.0438      0.019      2.341      0.019       0.007       0.081
HispPerCap     0.0357      0.025      1.435      0.151      -0.013       0.085
NumUnderPov    0.1112      0.138      0.805      0.421      -0.160       0.382
PctPopUnderPov -0.1721     0.063     -2.745      0.006      -0.295      -0.049
PctLess9thGrade -0.0999    0.068     -1.474      0.141      -0.233       0.033
PctNotHSGrad   0.0525      0.096      0.548      0.584      -0.136       0.241
PctBSorMore    0.0504      0.077      0.651      0.515      -0.101       0.202
PctUnemployed  0.0045      0.041      0.111      0.911      -0.075       0.084
PctEmploy      0.2485      0.079      3.151      0.002       0.094       0.403
PctEmplManu   -0.0658      0.032     -2.054      0.040      -0.129      -0.003
PctEmplProfServ -0.0267     0.041     -0.654      0.513      -0.107       0.053
PctOccupManu   0.0723      0.055      1.318      0.188      -0.035       0.180
PctOccupMgmtProf 0.1226     0.086      1.419      0.156      -0.047       0.292
MalePctDivorce 0.4585      0.248      1.851      0.064      -0.027       0.944
MalePctNevMarr 0.2267      0.068      3.339      0.001       0.094       0.360
FemalePctDiv   0.1627      0.309      0.526      0.599      -0.444       0.770
TotalPctDiv   -0.5619      0.519     -1.084      0.279      -1.579       0.455
PersPerFam    -0.1405      0.168     -0.834      0.404      -0.471       0.190
PctFam2Par     0.0186      0.160      0.117      0.907      -0.294       0.331
PctKids2Par   -0.3227      0.155     -2.080      0.038      -0.627      -0.018
PctYoungKids2Par -0.0323    0.048     -0.670      0.503      -0.127       0.062
PctTeen2Par   -0.0029      0.043     -0.069      0.945      -0.087       0.081
```

## 2.2 Question 3

### 2.2.1 Read the Data

```python
import numpy as np
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt
from sklearn.decomposition import PCA

# Load MNIST dataset and filter digits 3, 5, and 8
from sklearn.datasets import fetch_openml

# Fetch MNIST from openml
mnist = fetch_openml('mnist_784', version=1, as_frame=False)
X, y = mnist["data"], mnist["target"]

mnist_df = pd.DataFrame(np.concatenate((mnist['target'].reshape(-1, 1), mnist['data'
    ]), axis=1),
columns= ['target'] +  mnist['feature_names'])

mnist_df
```

Then we keep only 3, 5, and 8.

```python
import numpy as np
import pandas as pd
from sklearn.datasets import fetch_openml

# Fetch MNIST from OpenML
mnist = fetch_openml('mnist_784', version=1, as_frame=False)
X, y = mnist["data"], mnist["target"]

# Convert labels to integers
y = y.astype(int)

# Correct filtering: Keep only digits 3, 5, and 8
selected_digits = {3, 5, 8}
filter_mask = np.isin(y, list(selected_digits))

X_filtered = X[filter_mask]  # Keep only selected digits
y_filtered = y[filter_mask]  # Keep corresponding labels

# Print dataset size
print(f"Original dataset size:{X.shape[0]}")
print(f"Filtered dataset size (only 3, 5, 8):{X_filtered.shape[0]}")

# Convert to DataFrame (Optional)
mnist_df = pd.DataFrame(np.column_stack((y_filtered, X_filtered)),
columns=['target'] + mnist.feature_names)

# Display the dataframe
mnist_df
```

### 2.2.2 Logistic Regression with OvR

```python
from sklearn.datasets import fetch_openml
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
import numpy as np
```

```python
      import matplotlib.pyplot as plt

      # 1. Load MNIST
      mnist = fetch_openml('mnist_784', version=1, as_frame=False)
      X, y = mnist.data, mnist.target
      y = y.astype(int)

      # filter 3, 5, 8
      selected_digits = [3, 5, 8]
      mask = np.isin(y, selected_digits)
      X_filtered, y_filtered = X[mask], y[mask]

      X_filtered = X_filtered / 255.0

      # Divide the sets
      X_train, X_test, y_train, y_test = train_test_split(X_filtered, y_filtered, test_size
          =0.2, random_state=42)

      # 2. Train the Logistic Regression (One-vs-Rest)
      log_reg_ovr = LogisticRegression(multi_class='ovr', solver='liblinear', max_iter=100,
          random_state=42)
      log_reg_ovr.fit(X_train, y_train)

      # 3. Evaluate
      y_pred = log_reg_ovr.predict(X_test)
      accuracy = accuracy_score(y_test, y_pred)
      print(f"Test Accuracy: {accuracy * 100:.2f}%")

      # Confusion Matrix
      conf_matrix = confusion_matrix(y_test, y_pred)
      print("Confusion Matrix:")
      print(conf_matrix)

      # Class Report
      class_report = classification_report(y_test, y_pred)
      print("Classification Report:")
      print(class_report)

      # 4. Visualize
      num_images = 5
      indices = np.random.choice(len(X_test), num_images, replace=False)

      plt.figure(figsize=(10, 5))
      for i, index in enumerate(indices):
      plt.subplot(1, num_images, i + 1)
      plt.imshow(X_test[index].reshape(28, 28), cmap='gray')
      plt.title(f"Pred: {y_pred[index]}\nTrue: {y_test[index]}")
      plt.axis('off')
      plt.show()
```

```
Test Accuracy: 93.12%
Confusion Matrix:
[[1328   39   52]
 [  54 1195   45]
 [  39   50 1254]]
Classification Report:
precision    recall  f1-score   support

3        0.93       0.94      0.94      1419
5        0.93       0.92      0.93      1294
8        0.93       0.93      0.93      1343

accuracy                             0.93      4056
macro avg        0.93       0.93      0.93      4056
weighted avg        0.93       0.93      0.93      4056
```

### 2.2.3   Multinomial Regression

```python
from sklearn.datasets import fetch_openml
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
import numpy as np
import matplotlib.pyplot as plt

# Similar steps

mnist = fetch_openml('mnist_784', version=1, as_frame=False)
X, y = mnist.data, mnist.target
y = y.astype(int)

selected_digits = [3, 5, 8]
mask = np.isin(y, selected_digits)
X_filtered, y_filtered = X[mask], y[mask]

X_filtered = X_filtered / 255.0

X_train, X_test, y_train, y_test = train_test_split(X_filtered, y_filtered, test_size
    =0.2, random_state=42)

log_reg_multinomial = LogisticRegression(multi_class='multinomial', solver='lbfgs',
    max_iter=100, random_state=42)
log_reg_multinomial.fit(X_train, y_train)

y_pred = log_reg_multinomial.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Test Accuracy: {accuracy * 100:.2f}%")

conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(conf_matrix)

class_report = classification_report(y_test, y_pred)
print("Classification Report:")
print(class_report)

num_images = 5
indices = np.random.choice(len(X_test), num_images, replace=False)

plt.figure(figsize=(10, 5))
for i, index in enumerate(indices):
plt.subplot(1, num_images, i + 1)
```

```
43        plt.imshow(X_test[index].reshape(28, 28), cmap='gray')
44        plt.title(f"Pred:␣{y_pred[index]}\nTrue:␣{y_test[index]}")
45        plt.axis('off')
46        plt.show()
```

```
Test Accuracy: 93.12%
Confusion Matrix:
[[1324   43   52]
 [  51 1201   42]
 [  39   52 1252]]
Classification Report:
precision    recall  f1-score   support

3        0.94      0.93      0.93      1419
5        0.93      0.93      0.93      1294
8        0.93      0.93      0.93      1343

accuracy                          0.93      4056
macro avg        0.93      0.93      0.93      4056
weighted avg        0.93      0.93      0.93      4056
```

Do the Learning Curve

```
1     import numpy as np
2     import matplotlib.pyplot as plt
3     from sklearn.datasets import fetch_openml
4     from sklearn.model_selection import train_test_split, learning_curve
5     from sklearn.linear_model import LogisticRegression
6
7     # 1. Load and filter the MNIST dataset
8     mnist = fetch_openml('mnist_784', version=1, as_frame=False)
9     X, y = mnist.data, mnist.target
10    y = y.astype(int)
11
12    # Filter out samples with labels 3, 5, 8
13    selected_digits = [3, 5, 8]
14    mask = np.isin(y, selected_digits)
15    X_filtered, y_filtered = X[mask], y[mask]
16
17    # Scale pixel values from [0, 255] to [0, 1]
18    X_filtered = X_filtered / 255.0
19
20    # Split the dataset into training and testing sets
21    X_train, X_test, y_train, y_test = train_test_split(X_filtered, y_filtered, test_size
          =0.2, random_state=42)
22
23    # 2. Define the Logistic Regression model
24    log_reg = LogisticRegression(solver='lbfgs', max_iter=100, random_state=42)
25
26    # 3. Compute the learning curve
27    train_sizes, train_scores, val_scores = learning_curve(
28    log_reg, X_train, y_train, cv=5, scoring='accuracy', train_sizes=np.linspace(0.1,
          1.0, 10)
29    )
30
31    # Calculate mean and standard deviation of training and validation scores
32    train_scores_mean = np.mean(train_scores, axis=1)
33    train_scores_std = np.std(train_scores, axis=1)
34    val_scores_mean = np.mean(val_scores, axis=1)
35    val_scores_std = np.std(val_scores, axis=1)
36
37    # 4. Plot the learning curve
38    plt.figure(figsize=(10, 6))
```

```
39    plt.plot(train_sizes, train_scores_mean, label='Training␣Accuracy', color='blue',
         marker='o')
40    plt.fill_between(train_sizes, train_scores_mean - train_scores_std, train_scores_mean
         + train_scores_std, alpha=0.15, color='blue')
41    plt.plot(train_sizes, val_scores_mean, label='Validation␣Accuracy', color='green',
         marker='o')
42    plt.fill_between(train_sizes, val_scores_mean - val_scores_std, val_scores_mean +
         val_scores_std, alpha=0.15, color='green')
43
44    # plt.title('Learning Curve for Logistic Regression', fontsize=16)
45    plt.xlabel('Training␣Set␣Size', fontsize=14)
46    plt.ylabel('Accuracy', fontsize=14)
47    plt.legend(loc='best')
48    plt.grid(True)
49    plt.show()
```

Do the Validation Curve

```
1     from sklearn.model_selection import validation_curve
2
3     # Define the range of hyperparameter C (inverse of regularization strength)
4     param_range = np.logspace(-4, 4, 10)
5
6     # Compute validation curve
7     train_scores, val_scores = validation_curve(
8     log_reg, X_train, y_train, param_name='C', param_range=param_range, cv=5, scoring='
         accuracy'
9     )
10
11    # Calculate mean and standard deviation of training and validation scores
12    train_scores_mean = np.mean(train_scores, axis=1)
13    train_scores_std = np.std(train_scores, axis=1)
14    val_scores_mean = np.mean(val_scores, axis=1)
15    val_scores_std = np.std(val_scores, axis=1)
16
17    # Plot the validation curve
18    plt.figure(figsize=(10, 6))
19    plt.semilogx(param_range, train_scores_mean, label='Training␣Accuracy', color='blue',
         marker='o')
20    plt.fill_between(param_range, train_scores_mean - train_scores_std, train_scores_mean
         + train_scores_std, alpha=0.15, color='blue')
21    plt.semilogx(param_range, val_scores_mean, label='Validation␣Accuracy', color='green'
         , marker='o')
22    plt.fill_between(param_range, val_scores_mean - val_scores_std, val_scores_mean +
         val_scores_std, alpha=0.15, color='green')
23
24    # plt.title('Validation Curve for Logistic Regression', fontsize=16)
25    plt.xlabel('Regularization␣Strength␣(C)', fontsize=14)
26    plt.ylabel('Accuracy', fontsize=14)
27    plt.legend(loc='best')
28    plt.grid(True)
29    plt.show()
```

### 2.2.4 Naive Bayes

```
1     from sklearn.datasets import fetch_openml
2     from sklearn.model_selection import train_test_split
3     from sklearn.naive_bayes import GaussianNB
4     from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
5     import numpy as np
6     import matplotlib.pyplot as plt
7
```

```python
        # 1. Load and filter the MNIST dataset
        mnist = fetch_openml('mnist_784', version=1, as_frame=False)
        X, y = mnist.data, mnist.target
        y = y.astype(int)

        # Filter out samples with labels 3, 5, 8
        selected_digits = [3, 5, 8]
        mask = np.isin(y, selected_digits)
        X_filtered, y_filtered = X[mask], y[mask]

        # Scale pixel values from [0, 255] to [0, 1]
        X_filtered = X_filtered / 255.0

        # Split the dataset into training and testing sets
        X_train, X_test, y_train, y_test = train_test_split(X_filtered, y_filtered, test_size
            =0.2, random_state=42)

        # 2. Train a Gaussian Naive Bayes model
        naive_bayes = GaussianNB()
        naive_bayes.fit(X_train, y_train)

        # 3. Evaluate the model
        y_pred = naive_bayes.predict(X_test)
        accuracy = accuracy_score(y_test, y_pred)
        print(f"Test Accuracy: {accuracy * 100:.2f}%")

        # Confusion matrix
        conf_matrix = confusion_matrix(y_test, y_pred)
        print("Confusion Matrix:")
        print(conf_matrix)

        # Classification report
        class_report = classification_report(y_test, y_pred)
        print("Classification Report:")
        print(class_report)

        # 4. Visualize the results
        num_images = 5
        indices = np.random.choice(len(X_test), num_images, replace=False)

        plt.figure(figsize=(10, 5))
        for i, index in enumerate(indices):
        plt.subplot(1, num_images, i + 1)
        plt.imshow(X_test[index].reshape(28, 28), cmap='gray')
        plt.title(f"Pred: {y_pred[index]}\nTrue: {y_test[index]}")
        plt.axis('off')
        plt.show()
```

```
Test Accuracy: 50.22%
Confusion Matrix:
[[ 595    26  798]
 [  95   138 1061]
 [  18    21 1304]]
Classification Report:
precision    recall  f1-score   support

3        0.84      0.42      0.56      1419
5        0.75      0.11      0.19      1294
8        0.41      0.97      0.58      1343

accuracy                           0.50      4056
macro avg       0.67      0.50      0.44      4056
weighted avg       0.67      0.50      0.45      4056
```

### 2.2.5 Linear Discriminant Analysis

```python
from sklearn.datasets import fetch_openml
from sklearn.model_selection import train_test_split
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
import numpy as np
import matplotlib.pyplot as plt

# 1. Load and filter the MNIST dataset
mnist = fetch_openml('mnist_784', version=1, as_frame=False)
X, y = mnist.data, mnist.target
y = y.astype(int)

# Filter out samples with labels 3, 5, 8
selected_digits = [3, 5, 8]
mask = np.isin(y, selected_digits)
X_filtered, y_filtered = X[mask], y[mask]

# Scale pixel values from [0, 255] to [0, 1]
X_filtered = X_filtered / 255.0

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X_filtered, y_filtered, test_size
    =0.2, random_state=42)

# 2. Train a Linear Discriminant Analysis (LDA) model
lda = LinearDiscriminantAnalysis()
lda.fit(X_train, y_train)

# 3. Evaluate the model
y_pred = lda.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Test Accuracy (LDA): {accuracy * 100:.2f}%")

# Confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix (LDA):")
print(conf_matrix)

# Classification report
class_report = classification_report(y_test, y_pred)
print("Classification Report (LDA):")
print(class_report)

# 4. Visualize the results
```

```
44    num_images = 5
45    indices = np.random.choice(len(X_test), num_images, replace=False)
46
47    plt.figure(figsize=(10, 5))
48    for i, index in enumerate(indices):
49    plt.subplot(1, num_images, i + 1)
50    plt.imshow(X_test[index].reshape(28, 28), cmap='gray')
51    plt.title(f"Pred:␣{y_pred[index]}\nTrue:␣{y_test[index]}")
52    plt.axis('off')
53    plt.show()
```

```
Test Accuracy (LDA): 91.69%
Confusion Matrix (LDA):
[[1285   69   65]
 [  47 1202   45]
 [  34   77 1232]]
Classification Report (LDA):
precision    recall  f1-score   support

3        0.94       0.91       0.92       1419
5        0.89       0.93       0.91       1294
8        0.92       0.92       0.92       1343

accuracy                        0.92       4056
macro avg       0.92      0.92       0.92       4056
weighted avg       0.92      0.92       0.92       4056
```

### 2.2.6   Linear SVM

```
1     from sklearn.datasets import fetch_openml
2     from sklearn.model_selection import train_test_split
3     from sklearn.svm import LinearSVC
4     from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
5     import numpy as np
6     import matplotlib.pyplot as plt
7
8     # 1. Load and filter the MNIST dataset
9     mnist = fetch_openml('mnist_784', version=1, as_frame=False)
10    X, y = mnist.data, mnist.target
11    y = y.astype(int)
12
13    # Filter out samples with labels 3, 5, 8
14    selected_digits = [3, 5, 8]
15    mask = np.isin(y, selected_digits)
16    X_filtered, y_filtered = X[mask], y[mask]
17
18    # Scale pixel values from [0, 255] to [0, 1]
19    X_filtered = X_filtered / 255.0
20
21    # Split the dataset into training and testing sets
22    X_train, X_test, y_train, y_test = train_test_split(X_filtered, y_filtered, test_size
          =0.2, random_state=42)
23
24    # 2. Train a Linear SVM model with One-vs-Rest strategy
25    linear_svm = LinearSVC(multi_class='ovr', max_iter=10000, random_state=42)  # Use One
          -vs-Rest
26    linear_svm.fit(X_train, y_train)
27
28    # 3. Evaluate the model
29    y_pred = linear_svm.predict(X_test)
30    accuracy = accuracy_score(y_test, y_pred)
```

```
31        print(f"Test␣Accuracy␣(Linear␣SVM):␣{accuracy␣*␣100:.2f}%")
32
33        # Confusion matrix
34        conf_matrix = confusion_matrix(y_test, y_pred)
35        print("Confusion␣Matrix␣(Linear␣SVM):")
36        print(conf_matrix)
37
38        # Classification report
39        class_report = classification_report(y_test, y_pred)
40        print("Classification␣Report␣(Linear␣SVM):")
41        print(class_report)
42
43        # 4. Visualize the results
44        num_images = 5
45        indices = np.random.choice(len(X_test), num_images, replace=False)
46
47        plt.figure(figsize=(10, 5))
48        for i, index in enumerate(indices):
49        plt.subplot(1, num_images, i + 1)
50        plt.imshow(X_test[index].reshape(28, 28), cmap='gray')
51        # plt.title(f"Pred: {y_pred[index]}\nTrue: {y_test[index]}")
52        plt.axis('off')
53        plt.show()
```

```
Test Accuracy (Linear SVM): 92.73%
Confusion Matrix (Linear SVM):
[[1319   44   56]
 [  57 1193   44]
 [  39   55 1249]]
Classification Report (Linear SVM):
precision    recall  f1-score   support

3        0.93       0.93       0.93      1419
5        0.92       0.92       0.92      1294
8        0.93       0.93       0.93      1343

accuracy                         0.93      4056
macro avg        0.93      0.93       0.93      4056
weighted avg       0.93      0.93       0.93       4056
```

### 2.2.7   Plot the Confusion Matrix

```
1        from sklearn.metrics import confusion_matrix
2        import seaborn as sns
3
4        # Logistic Regression Confusion Matrix
5        log_reg_pred = log_reg.predict(X_test)
6        log_reg_cm = confusion_matrix(y_test, log_reg_pred)
7
8        # Multinomial Regression Confusion Matrix
9        multinomial_pred = multinomial_reg.predict(X_test)
10        multinomial_cm = confusion_matrix(y_test, multinomial_pred)
11
12        # Naive Bayes Confusion Matrix
13        naive_bayes_pred = naive_bayes.predict(X_test)
14        naive_bayes_cm = confusion_matrix(y_test, naive_bayes_pred)
15
16        # LDA Confusion Matrix
17        lda_pred = lda.predict(X_test)
18        lda_cm = confusion_matrix(y_test, lda_pred)
19
```

```
20        # SVM Confusion Matrix
21        svm_pred = svm.predict(X_test)
22        svm_cm = confusion_matrix(y_test, svm_pred)
23
24        def plot_confusion_matrix(cm, class_names):
25        plt.figure(figsize=(6, 5))
26        sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=class_names,
              yticklabels=class_names)
27        plt.xlabel('Predicted')
28        plt.ylabel('True')
29        plt.title('Confusion␣Matrix')
30        plt.show()
31
32        class_names = ['3', '5', '8']
33
34        # Plot confusion matrices
35        plot_confusion_matrix(log_reg_cm, class_names)
36        plot_confusion_matrix(multinomial_cm, class_names)
37        plot_confusion_matrix(naive_bayes_cm, class_names)
38        plot_confusion_matrix(lda_cm, class_names)
39        plot_confusion_matrix(svm_cm, class_names)
```