# Digit Recognizer Report

## Overview

The MNIST ("Modified National Institute of Standards and Technology") dataset is a classic benchmark in computer vision. Since its release in 1999, this dataset has served as a foundational resource for evaluating the performance of classification algorithms. The dataset comprises grey-scale images of hand-drawn digits ranging from zero to nine, each represented as a 28x28 pixel matrix. The task is to accurately classify these digits based on their pixel values.

## Experimental Results

### Random Forest Classifier

A Random Forest classifier was first applied to the full set of explanatory variables. The model fitting and prediction times, along with the validation accuracy, were recorded as follows:

- **Time taken to fit the Random Forest model:** 0:00:11.763256 sec
- **Time taken to make predictions:** 0:00:00.115075
- **Validation Accuracy:** 0.9655
- **Kaggle Score:** 0.96364 (rfc_prediction.csv)

The Random Forest classifier achieved a high validation accuracy of 96.55%, demonstrating its effectiveness in capturing the patterns in the MNIST dataset.

### Principal Component Analysis (PCA)

PCA was applied to reduce the dimensionality of the data while retaining 95% of the variability in the explanatory variables. The following results were obtained:

- **Time taken to perform PCA:** 0:00:02.425129 sec
- **Number of PCA components at 95% variability:** 153

- **Time for PCA fit_transform and inverse_transform:** 0:00:01.682767 sec
- **Reconstruction error:** 220.1751

The PCA reduced the number of dimensions from 784 to 153, significantly decreasing the computational load while preserving most of the information.

## PCA with Random Forest Classifier

The Random Forest classifier was then retrained using the principal components identified by PCA. The results were as follows:

- **Time taken to fit the Random Forest model:** 0:00:34.847865 sec
- **Validation Accuracy:** 0.9414
- **Kaggle Score:** 0.93967 (pca_prediction.csv)

While the PCA reduced the computational complexity, the validation accuracy dropped slightly to 94.14%, indicating some loss of information during dimensionality reduction.

## K-Means Clustering

Initially, K-Means clustering was applied with 10 clusters, resulting in the following performance:

- **Time taken to fit the K-Means model:** 0:00:12.272523 sec
- **Training Accuracy:** 0.59
- **Validation Accuracy:** 0.60
- **Kaggle Score:** 0.59114 (kmeans_prediction.csv)

The low accuracy indicated a need to reconsider the number of clusters.

## Fixing the Design Flaw: K-Means with 2000 Clusters

To address the design flaw, the number of clusters was increased to 2000, leading to improved performance:

- **Time taken to fit the K-Means model with 2000 clusters:** 0:00:29.022648 sec
- **Training Accuracy:** 0.94
- **Validation Accuracy:** 0.94
- **Kaggle Score:** 0.93792 (kmeans_v2_prediction.csv)

By increasing the number of clusters, the K-Means model achieved a validation accuracy of 94%, comparable to the PCA-augmented Random Forest classifier.

# Kaggle Details

**Username**: sachinsharma03

**Submissions**

| All | Successful | Errors | | Recent ▾ |
|---|---|---|---|---|

| Submission and Description | | Public Score ⓘ |
|---|---|---|
| ✓ | **kmeans_v2_prediction.csv**<br>Complete · 26m ago | **0.93792** |
| ✓ | **kmeans_prediction.csv**<br>Complete · 26m ago | **0.59114** |
| ✓ | **pca_prediction.csv**<br>Complete · 26m ago | **0.93967** |
| ✓ | **rfc_prediction.csv**<br>Complete · 27m ago | **0.96364** |

# Appendix

**Data Characteristics:**

- **Training Data:** 42,000 samples
- **Test Data:** 28,000 samples
- **Image Dimensions:** 28x28 pixels
- **Pixel Range:** 0-255

**Model Evaluation Metric:**

- **Accuracy:** Percentage of correctly classified images

**Important Tools and Libraries Used:**

- **Pandas**
- **NumPy**
- **Scikit-learn**
- **Seaborn**
- **Matplotlib**
- **Datetime**

# Module 6 Assignment 3: Digit Recognizer

**Sachin Sharma**

**MSDS-422**

**07/26/2022**

## Management/Research Question

In layman's terms, what is the management/research question of interest, and why would anyone care?

Requirements

- Fit a random forest classifier using the full set of explanatory variables and the model training set (csv).
- Record the time it takes to fit the model and then evaluate the model on the csvdata by submitting to Kaggle.com. Provide your Kaggle.com score and user ID.
- Execute principal components analysis (PCA) on the combined training and test set data together, generating principal components that represent 95 percent of the variability in the explanatory variables. The number of principal components in the solution should be substantially fewer than the explanatory variables.
- Record the time it takes to identify the principal components.
- Using the identified principal components from step (2), use thecsvto build another random forest classifier.
- Record the time it takes to fit the model and to evaluate the model on the csvdata by submitting to Kaggle.com. Provide your Kaggle.com score and user ID.
- Use k-means clustering to group MNIST observations into 1 of 10 categories and then assign labels. (Follow the example here if needed: kmeans mnist.pdf). kmeans mnist-2.pdf Download kmeans mnist-2.pdf
- Submit the RF Classifier, the PCA RF, and k-means estimations to Kaggle.com, and provide screen snapshots of your scores as well as your Kaggle.com user name.
- The experiment we have proposed has a major design flaw. Identify the flaw. Fix it. Rerun the experiment in a way that is consistent with a training-and-test regimen, and submit this to Kaggle.com.
- Report total elapsed time measures for the training set analysis. It is sufficient to run a single time-elapsed test for this assignment. In practice, we might consider the possibility of repeated executions of the relevant portions of the programs, much as the Benchmark Example programs do. Some code that might help you with reporting elapsed total time follows.

## Libraries to be loaded

```
import matplotlib as mpl
import matplotlib.pyplot as plt
```

```python
import numpy as np
import pandas as pd
import seaborn as sns
from datetime import datetime
from sklearn.cluster import KMeans
from sklearn.cluster import MiniBatchKMeans
from sklearn.decomposition import PCA
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix
import warnings
# Ignore all FutureWarnings
warnings.filterwarnings('ignore')
```

## Ingest

```python
train_df = pd.read_csv("digit-recognizer/train.csv")
train_df.shape
```

```
(42000, 785)
```

```python
train_df.describe()
```

```
              label    pixel0    pixel1    pixel2    pixel3    pixel4
pixel5  \
count  42000.000000   42000.0   42000.0   42000.0   42000.0   42000.0
42000.0
mean       4.456643       0.0       0.0       0.0       0.0       0.0
0.0
std        2.887730       0.0       0.0       0.0       0.0       0.0
0.0
min        0.000000       0.0       0.0       0.0       0.0       0.0
0.0
25%        2.000000       0.0       0.0       0.0       0.0       0.0
0.0
50%        4.000000       0.0       0.0       0.0       0.0       0.0
0.0
75%        7.000000       0.0       0.0       0.0       0.0       0.0
0.0
max        9.000000       0.0       0.0       0.0       0.0       0.0
0.0

         pixel6    pixel7    pixel8  ...        pixel774        pixel775  \
count   42000.0   42000.0   42000.0  ...    42000.000000    42000.000000
mean        0.0       0.0       0.0  ...        0.219286        0.117095
std         0.0       0.0       0.0  ...        6.312890        4.633819
min         0.0       0.0       0.0  ...        0.000000        0.000000
25%         0.0       0.0       0.0  ...        0.000000        0.000000
```

```
50%          0.0         0.0         0.0  ...      0.000000         0.000000
75%          0.0         0.0         0.0  ...      0.000000         0.000000
max          0.0         0.0         0.0  ...    254.000000       254.000000

           pixel776      pixel777       pixel778       pixel779  pixel780  \
count  42000.000000  42000.00000  42000.000000  42000.000000   42000.0

mean       0.059024      0.02019      0.017238      0.002857       0.0

std        3.274488      1.75987      1.894498      0.414264       0.0

min        0.000000      0.00000      0.000000      0.000000       0.0

25%        0.000000      0.00000      0.000000      0.000000       0.0

50%        0.000000      0.00000      0.000000      0.000000       0.0

75%        0.000000      0.00000      0.000000      0.000000       0.0

max      253.000000    253.00000    254.000000     62.000000       0.0


       pixel781  pixel782  pixel783
count   42000.0   42000.0   42000.0
mean        0.0       0.0       0.0
std         0.0       0.0       0.0
min         0.0       0.0       0.0
25%         0.0       0.0       0.0
50%         0.0       0.0       0.0
75%         0.0       0.0       0.0
max         0.0       0.0       0.0

[8 rows x 785 columns]

test_df = pd.read_csv("digit-recognizer/test.csv")
test_df.shape

(28000, 784)
```

# EDA

```
train_df_X = train_df.copy()
train_df_y = train_df_X['label']
train_df_X.drop(['label'], axis=1, inplace=True)

train_df_X.shape

(42000, 784)

train_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 42000 entries, 0 to 41999
Columns: 785 entries, label to pixel783
dtypes: int64(785)
memory usage: 251.5 MB
```
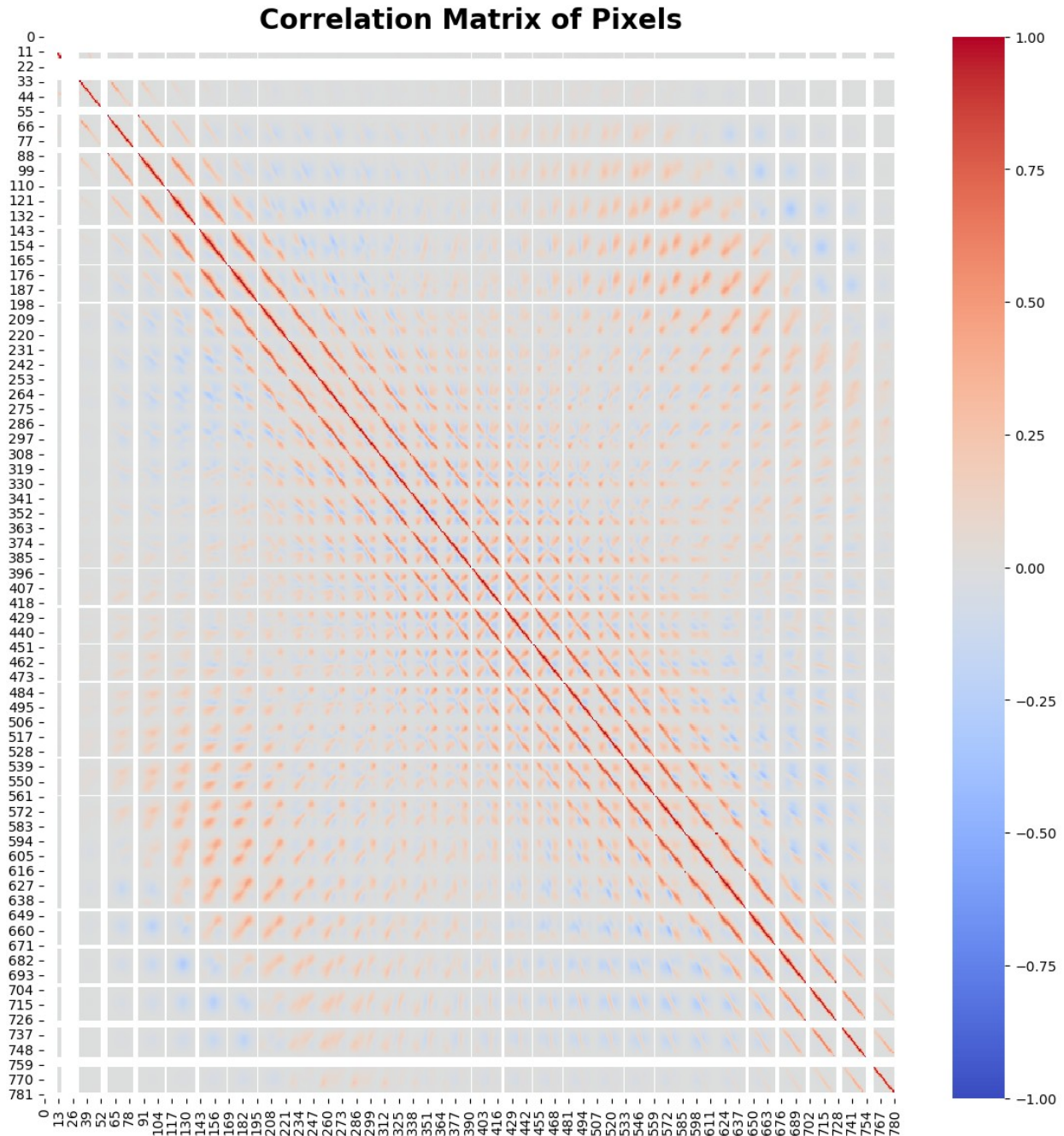
There are no null values in the training or test data.

```python
print("Null values in Train DF: ",train_df.isna().sum().sum())
print("Null values in Test DF: ",test_df.isna().sum().sum())

Null values in Train DF:  0
Null values in Test DF:  0
```

Let's output some sample digits as a 28x28 pixel image.

```python
# Compute the correlation matrix
images = train_df.drop(columns=['label']).values
labels = train_df['label'].values

flat_images = images.reshape(-1, 28*28)
correlation_matrix = np.corrcoef(flat_images.T)

# Plot the correlation matrix
plt.figure(figsize=(14, 14))
sns.heatmap(correlation_matrix, cmap='coolwarm', vmin=-1, vmax=1)
plt.title("Correlation Matrix of Pixels", size=20, fontweight='bold')
plt.show();
```

Correlation Matrix of Pixels

```python
labels = train_df['label'].values
images = train_df.drop(columns=['label']).values

# Reshape the images
images = images.reshape(-1, 28, 28)

# Combine images and labels into a single dataset
dataset = list(zip(images, labels))

# Plot images from the dataset with their corresponding labels
```
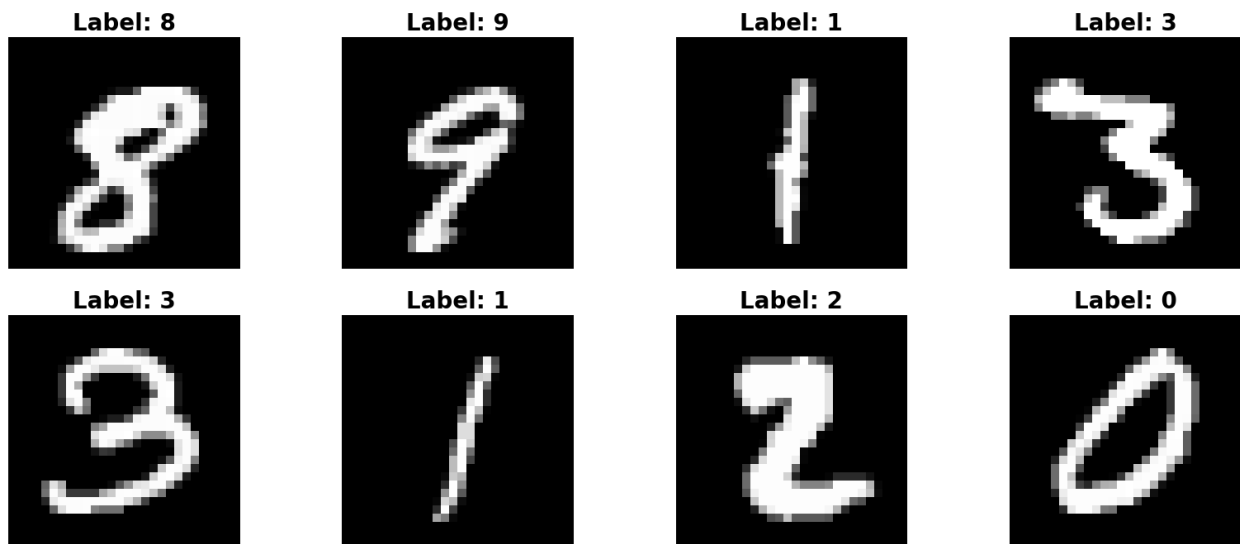
```
plt.figure(figsize=(20, 8))
for i in range(10, 18):
    image, label = dataset[i]
    plt.subplot(2, 4, i-9)  # Adjust subplot indexing to match the
desired layout (2 rows, 4 columns)
    plt.imshow(image, cmap='gray')
    plt.title('Label: ' + str(label), fontweight='bold', size=20)
    plt.axis('off')  # Turn off axis

plt.show()
```



Let's look at the distribution of digits in the training set.

```
# Set the style and color palette
sns.set(style="whitegrid", palette="pastel")

# Create the figure and axis objects
fig, ax = plt.subplots(figsize=(18, 6))

# Plot the countplot
sns.countplot(x='label', data=train_df, ax=ax)

# Customize the plot
ax.set_title("Class Distribution", size=24, fontweight='bold')
ax.set_ylabel("No. of Observations", size=20)
ax.set_xlabel("Class Name", size=20)
ax.tick_params(axis='both', which='major', labelsize=14)
ax.tick_params(axis='x', rotation=45)  # Rotate x-axis labels for
better readability

# Remove the top and right spines
sns.despine()
```
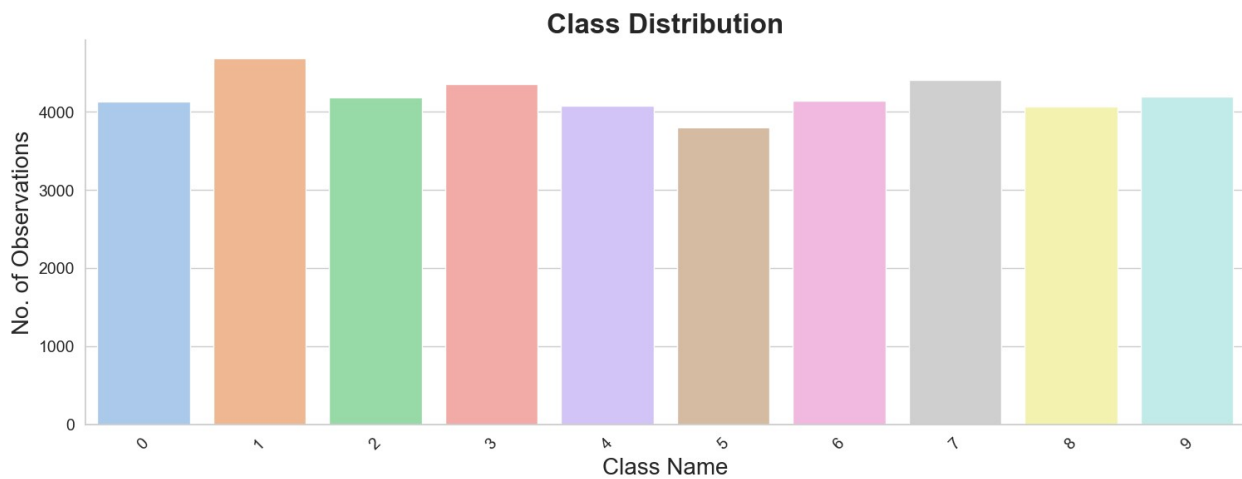
```
# Show the plot
plt.show()
```



**Class Distribution**

# MODELLING

Let's split train.csv 80/20 into training and validation datasets.

```
X_train, X_val, y_train, y_val = train_test_split(train_df_X,
train_df_y, test_size=0.2, random_state=42)
```

## RANDOM FOREST CLASSIFIER

Fit a random forest classifier using the entire set of explanatory variables from the training set (csv). Measure and record the time required to fit the model.

```
start_time = datetime.now()
rfc_model = RandomForestClassifier(random_state=42)
rfc_model.fit(X_train, y_train)
end_time = datetime.now()
# Calculate the time taken to fit the model
time_taken = end_time - start_time
print(f"Time taken to fit the Random Forest model: {time_taken}")

start_time = datetime.now()
y_pred = rfc_model.predict(X_val)
end_time = datetime.now()
# Calculate the time taken to predict
time_taken_predict = end_time - start_time
print(f"Time taken to make predictions: {time_taken_predict}")

# Evaluate the model
accuracy = accuracy_score(y_val, y_pred)
print("Validation Accuracy: ",accuracy)
```

```
print("Classification Report:")
print(classification_report(y_val, y_pred))

Time taken to fit the Random Forest model: 0:00:11.878233
Time taken to make predictions: 0:00:00.134266
Validation Accuracy:  0.9654761904761905
Classification Report:
              precision    recall  f1-score   support

           0       0.98      0.99      0.98       816
           1       0.98      0.99      0.99       909
           2       0.97      0.96      0.97       846
           3       0.96      0.95      0.96       937
           4       0.96      0.97      0.96       839
           5       0.96      0.96      0.96       702
           6       0.97      0.98      0.98       785
           7       0.97      0.95      0.96       893
           8       0.96      0.96      0.96       835
           9       0.94      0.94      0.94       838

    accuracy                           0.97      8400
   macro avg       0.97      0.97      0.97      8400
weighted avg       0.97      0.97      0.97      8400
```

## Submission

Evaluate the model on the csvdata by submitting to Kaggle.com.

```python
y = rfc_model.predict(test_df)
# create submission file
submission = pd.DataFrame({"ImageId": (test_df.index + 1),"Label": y})
submission.to_csv('rfc_prediction.csv', index=False)
```

## PCA

Perform principal components analysis (PCA) on the combined training and test set data to generate principal components that capture 95 percent of the variability in the explanatory variables. Ensure that the number of principal components is significantly fewer than the original number of explanatory variables. Measure and record the time taken to identify these principal components.

```python
# Record the start time for PCA
start_time = datetime.now()

# Perform PCA on the training data
pca = PCA()
pca.fit(X_train)
```

```python
# Calculate the cumulative sum of explained variance ratios
cumsum = np.cumsum(pca.explained_variance_ratio_)

# Find the number of components that explain at least 95% of the
variance
pca_n_components = np.argmax(cumsum >= 0.95) + 1

# Record the end time for PCA
end_time = datetime.now()

time_taken_fit = end_time - start_time

# Print the time taken and the number of components
print('Time taken to fit the PCA model: ', time_taken_fit)
print('# PCA Components at 95% variability: ', pca_n_components)

Time taken to fit the PCA model:  0:00:02.027835
# PCA Components at 95% variability:   153

# Create the plot
plt.figure(figsize=(10, 6))
plt.plot(cumsum, linewidth=2, color='blue', label='Cumulative
Explained Variance')
plt.axhline(y=0.95, color='r', linestyle='--', linewidth=1.5)
plt.axvline(x=pca_n_components, color='r', linestyle='--',
linewidth=1.5)
plt.scatter(pca_n_components, 0.95, color='red', s=50)

# Annotate the "Elbow" point
plt.annotate("Elbow", xy=(pca_n_components, 0.95),
xytext=(pca_n_components + 50, 0.8),
             arrowprops=dict(facecolor='black', shrink=0.05),
             fontsize=14, color='black')

# Set titles and labels
plt.title("Explained Variance vs. Number of Dimensions", fontsize=20,
fontweight='bold')
plt.xlabel("Number of Principal Components", fontsize=16)
plt.ylabel("Cumulative Explained Variance", fontsize=16)
plt.axis([0, 400, 0, 1])
plt.grid(True, which='both', linestyle='--', linewidth=0.7)
plt.legend(loc='lower right', fontsize=14)

# Show the plot
plt.show()
```
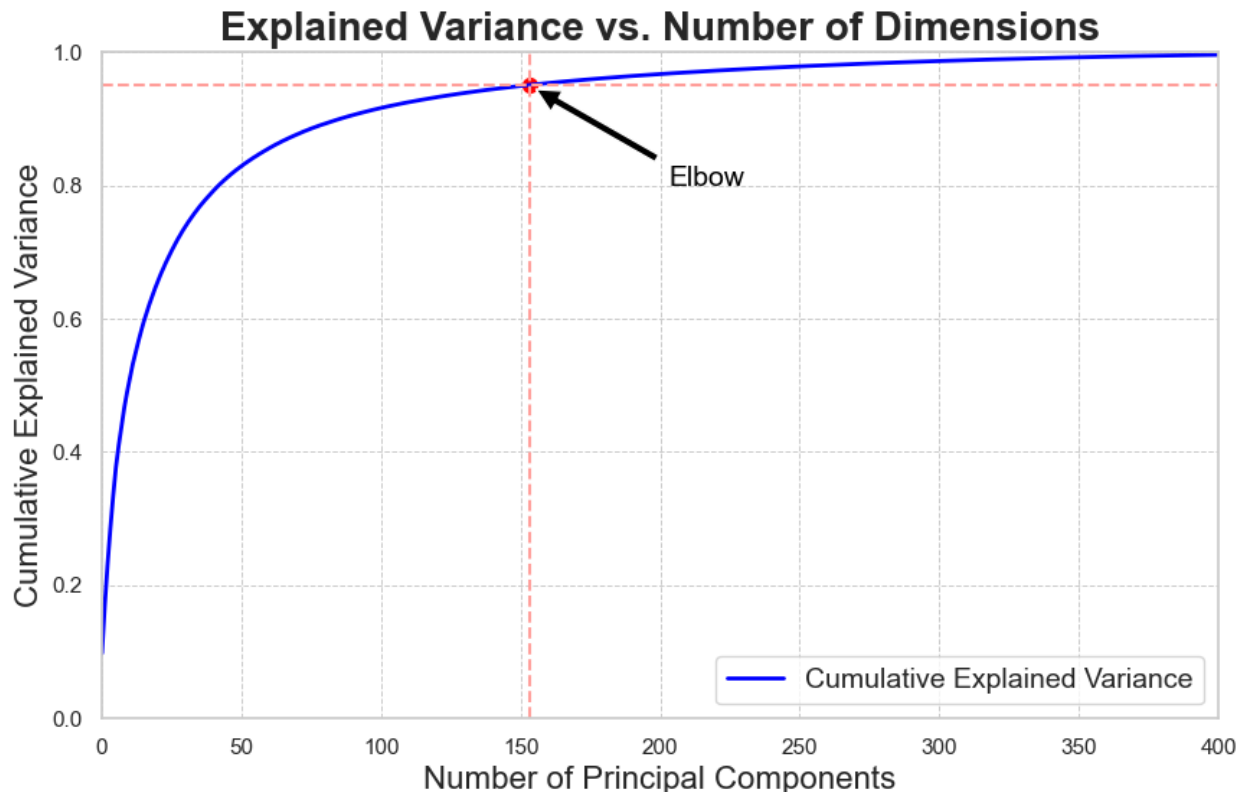
**Explained Variance vs. Number of Dimensions**

```python
# Perform PCA with the determined number of components and measure the
time
start = datetime.now()
pca = PCA(n_components=pca_n_components)
X_train_pca = pca.fit_transform(X_train)
X_recovered = pca.inverse_transform(X_train_pca)
end = datetime.now()
print('Time for PCA fit_transform and inverse_transform: ', end -
start)

# Optionally, evaluate the reconstructed data
reconstruction_error = np.mean((X_train - X_recovered) ** 2)
print(f'Reconstruction error: {reconstruction_error}')

Time for PCA fit_transform and inverse_transform:  0:00:01.801173
Reconstruction error: 220.25837900189745
```

## PCA RANDOM FOREST CLASSIFIER

Utilize the principal components identified in step (2) to construct a new random forest classifier using the data from train.csv. Measure and document the time required to fit this model.

```python
# Record the start time for fitting the Random Forest model
start_time = datetime.now()
```

```python
# Perform PCA with the determined number of components
pca = PCA(n_components=pca_n_components)
X_train_pca = pca.fit_transform(X_train)
X_val_pca = pca.transform(X_val)

# Initialize the Random Forest Classifier
rf_classifier = RandomForestClassifier(random_state=42)

# Fit the model to the PCA-transformed training data
rf_classifier.fit(X_train_pca, y_train)

# Record the end time for fitting the Random Forest model
end_time = datetime.now()

# Calculate the time taken to fit the model
time_taken = end_time - start_time
print(f"Time taken to fit the PCA Random Forest model: {time_taken}")

# Predict on the PCA-transformed validation set
y_pred = rf_classifier.predict(X_val_pca)

# Evaluate the model
accuracy = accuracy_score(y_val, y_pred)
print("Validation Accuracy: ", accuracy)
print("Classification Report:")
print(classification_report(y_val, y_pred))
```

```
Time taken to fit the PCA Random Forest model: 0:00:33.983587
Validation Accuracy:  0.9375
Classification Report:
              precision    recall  f1-score   support

           0       0.96      0.97      0.97       816
           1       0.98      0.98      0.98       909
           2       0.95      0.93      0.94       846
           3       0.91      0.91      0.91       937
           4       0.93      0.94      0.94       839
           5       0.91      0.93      0.92       702
           6       0.94      0.97      0.95       785
           7       0.95      0.93      0.94       893
           8       0.92      0.90      0.91       835
           9       0.91      0.92      0.91       838

    accuracy                           0.94      8400
   macro avg       0.94      0.94      0.94      8400
weighted avg       0.94      0.94      0.94      8400
```

## Submission

Evaluate the model on the csvdata by submitting to Kaggle.com.

```
test_df_pca = pca.transform(test_df)
y = rf_classifier.predict(test_df_pca)

# create submission file
submission = pd.DataFrame({"ImageId": (test_df.index + 1),"Label": y})
submission.to_csv('pca_prediction.csv', index=False)
```

As expected, with PCA compressing the data into less features, we would get a lower, but still fairly accurate, Kaggle score.

## K-MEANS CLASSIFIER

Use k-means clustering to group MNIST observations into 1 of 10 categories and then assign labels.

```
# Use k-means clustering to group the observations
kmeans = KMeans(n_clusters=10, random_state=42)
start_time = datetime.now()
kmeans.fit(X_train)
end_time = datetime.now()

# Calculate the time taken to fit the k-means model
time_taken = end_time - start_time
print(f"Time taken to fit the k-means model: {time_taken}")

# Predict cluster labels for the training and validation sets
train_cluster_labels = kmeans.predict(X_train)
val_cluster_labels = kmeans.predict(X_val)

# Map cluster labels to actual labels (majority voting within each
cluster)
def map_cluster_labels_to_actual_labels(cluster_labels,
actual_labels):
    label_mapping = {}
    for cluster in np.unique(cluster_labels):
        mask = cluster_labels == cluster
        most_common_label = np.bincount(actual_labels[mask]).argmax()
        label_mapping[cluster] = most_common_label
    return label_mapping

# Map the cluster labels to the actual labels for both training and
validation sets
train_label_mapping =
map_cluster_labels_to_actual_labels(train_cluster_labels, y_train)
val_label_mapping =
map_cluster_labels_to_actual_labels(val_cluster_labels, y_val)
```
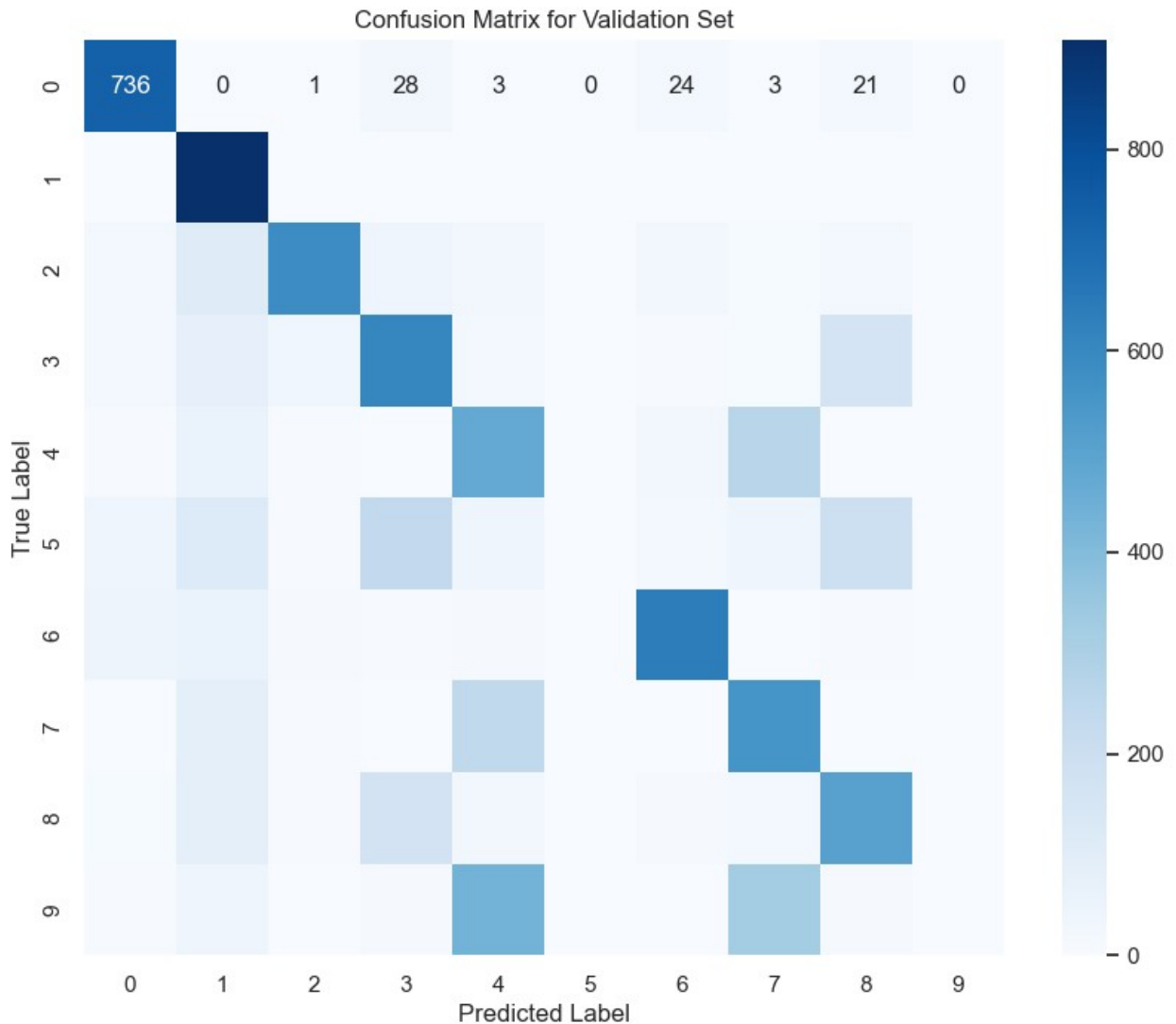
```python
# Assign the mapped labels
train_pred_labels = np.array([train_label_mapping[cluster] for cluster
in train_cluster_labels])
val_pred_labels = np.array([val_label_mapping[cluster] for cluster in
val_cluster_labels])

# Evaluate the clustering performance
train_accuracy = accuracy_score(y_train, train_pred_labels)
val_accuracy = accuracy_score(y_val, val_pred_labels)
print(f"Training Accuracy: {train_accuracy:.2f}")
print(f"Validation Accuracy: {val_accuracy:.2f}")

# Confusion matrix for validation set
conf_matrix = confusion_matrix(y_val, val_pred_labels)
plt.figure(figsize=(10, 8))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues',
xticklabels=np.unique(y_val), yticklabels=np.unique(y_val))
plt.title("Confusion Matrix for Validation Set")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()

Time taken to fit the k-means model: 0:00:11.586316
Training Accuracy: 0.59
Validation Accuracy: 0.60
```

Confusion Matrix for Validation Set

## Submission

Evaluate the model on the csvdata by submitting to Kaggle.com.

```python
# Predict cluster labels for the test set
test_cluster_labels = kmeans.predict(test_df)

# Assign the mapped labels to the test set
test_labels = [val_label_mapping[cluster] for cluster in
test_cluster_labels]

# Create submission file
submission = pd.DataFrame({"ImageId": test_df.index + 1, "Label":
test_labels})
submission.to_csv('kmeans_prediction.csv', index=False)
```

# DESIGN FLAW

The experiment we proposed has a major design flaw. Identify the flaw and fix it.

The number of K-means clusters isn't strictly limited to the number of available labels; we can increase the number of clusters to allow labels to be better represented in multiple, smaller clusters.

Although MiniBatchKMeans may sometimes be prone to memory leaks (which can be mitigated by using `batch_size=2048`), and may not be as accurate as the traditional KMeans API, increasing the number of `n_clusters` can compensate for this accuracy gap while still enabling the model to process the data in a reasonable amount of time.

```python
# Define the single number of clusters
n_clusters = 2000

# Initialize MiniBatchKMeans with the specified number of clusters
kmeans = MiniBatchKMeans(n_clusters=n_clusters, batch_size=2048,
random_state=42)

# Fit the k-means model
start_time = datetime.now()
kmeans.fit(X_train)
end_time = datetime.now()

# Calculate the time taken to fit the k-means model
time_taken = end_time - start_time
print(f"Time taken to fit the k-means model with {n_clusters}
clusters: {time_taken}")

# Predict cluster labels for the training and validation sets
train_cluster_labels = kmeans.predict(X_train)
val_cluster_labels = kmeans.predict(X_val)

# Map cluster labels to actual labels (majority voting within each
cluster)
def map_cluster_labels_to_actual_labels(cluster_labels,
actual_labels):
    label_mapping = {}
    for cluster in np.unique(cluster_labels):
        mask = cluster_labels == cluster
        if np.any(mask):
            most_common_label =
np.bincount(actual_labels[mask]).argmax()
            label_mapping[cluster] = most_common_label
    return label_mapping

# Map the cluster labels to the actual labels for training data
train_label_mapping =
map_cluster_labels_to_actual_labels(train_cluster_labels, y_train)
```

```python
# Assign the mapped labels to both training and validation sets
train_pred_labels = np.array([train_label_mapping[cluster] for cluster
in train_cluster_labels])
val_pred_labels = np.array([train_label_mapping[cluster] for cluster
in val_cluster_labels])

# Evaluate the clustering performance
train_accuracy = accuracy_score(y_train, train_pred_labels)
val_accuracy = accuracy_score(y_val, val_pred_labels)

print(f"Number of Clusters: {n_clusters}")
print(f"Training Accuracy: {train_accuracy:.2f}")
print(f"Validation Accuracy: {val_accuracy:.2f}")

# Confusion matrix for validation set
conf_matrix = confusion_matrix(y_val, val_pred_labels)
plt.figure(figsize=(10, 8))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues',
xticklabels=np.unique(y_val), yticklabels=np.unique(y_val))
plt.title(f"Confusion Matrix for Validation Set with {n_clusters}
Clusters")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()

Time taken to fit the k-means model with 2000 clusters: 0:00:29.022648
Number of Clusters: 2000
Training Accuracy: 0.94
Validation Accuracy: 0.94
```
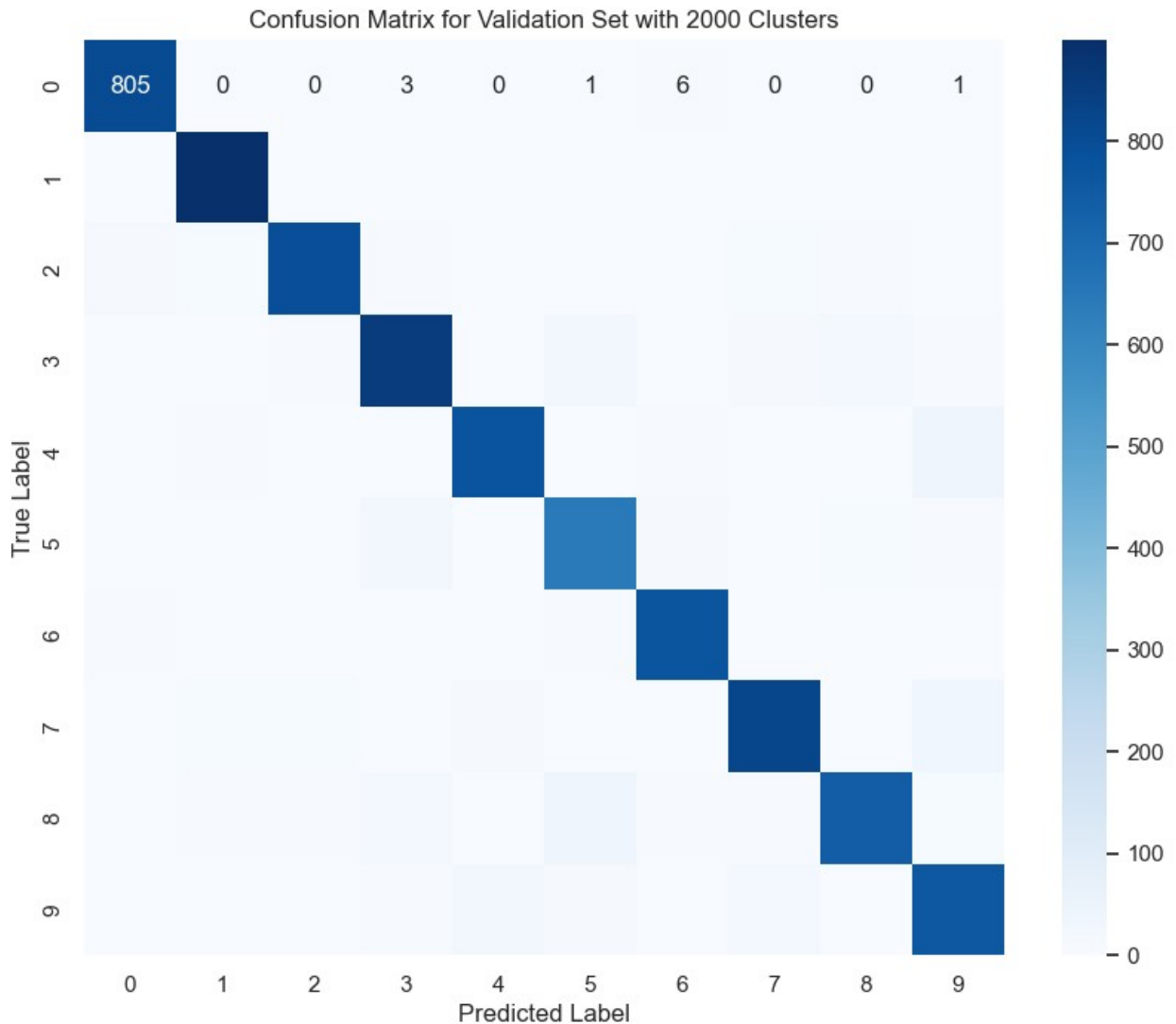
Confusion Matrix for Validation Set with 2000 Clusters

## Submission

Rerun the experiment in a manner that adheres to a proper training-and-test regimen, and then submit the results to Kaggle.com.

```python
# Predict cluster labels for the test data
kmeans_pred_fix = kmeans.predict(test_df)

# Map cluster labels to actual labels for test data
test_labels_fix = [train_label_mapping.get(cluster, -1) for cluster in
kmeans_pred_fix]

# Create submission file
submission = pd.DataFrame({"ImageId": (test_df.index + 1), "Label":
test_labels_fix})
submission.to_csv('kmeans_v2_prediction.csv', index=False)
```

**Management/Research Question Question of Interest: How can we accurately classify handwritten digits using machine learning models?**

**Why It Matters:**

Accurately classifying handwritten digits with machine learning models is crucial for automating data processing tasks and reducing human error. This technology can streamline processes in various sectors, such as postal services and banking, by automating the reading of addresses and processing of checks. Additionally, advancements in this area can improve optical character recognition (OCR) systems, enhancing their ability to digitize printed and handwritten text, which has broad applications in digitizing documents and forms.

Beyond practical applications, this research serves as a foundational learning tool in the field of machine learning. The MNIST dataset is widely used to teach and understand the basics of computer vision. Success in this domain not only pushes the boundaries of technology but also provides essential knowledge for tackling more complex data science challenges. Overall, improving handwritten digit classification can lead to greater efficiency, accuracy, and technological advancement in various fields.