

# Machine Learning HW #1

Andrei-Gabriel Mitran

[andreigabrielmitran@gmail.com](mailto:andreigabrielmitran@gmail.com)

Faculty of Automatic Control and Computer Science, Polytechnic University of Bucharest

This report talks about the methods and models used for classifying two very popular datasets: Fashion-MNIST (items of clothing) and Fruits-360 (fruits and nuts). It will go through how the features were extracted, what they represent and how the best models performed for both datasets.

## 1. Feature Extraction Methods

For this homework, I have used two feature extraction methods for the datasets: **Principal Component Analysis (PCA)** and **Histogram of Oriented Gradients (HOG)**.

### a. Principal Component Analysis (PCA)

PCA is a technique used to reduce the dimensionality of the data by finding the most important features that explain the most variance in the data. In the case of images, PCA helps in identifying the key patterns (like overall shape or color), all the while reducing the number of pixels needed to represent the image.

It simplifies the data by focusing on the most important aspects, which helps speed up the classification process. For example, PCA can capture the general shape of objects like clothing or fruits.

### b. Histogram of Oriented Gradients (HOG)

HOG is a feature descriptor that focuses on capturing the edges and textures in an image. It works by calculating the gradient of pixel intensities in small regions of the image and summarizing the direction of those gradients.

It is great for detecting shapes and textures, which are important for identifying objects. For example, HOG can help capture the fine details of a clothing item in Fashion-MNIST or the contours of a fruit in Fruits-360.

### c. Combining PCA and HOG

While PCA captures the overall shape of the object, HOG captures finer details like edges and textures. By combining both, we get a richer set of features that can help the model better differentiate between classes.

Basically, PCA reduces the data to its most important features, making it easier to handle, while HOG adds important local information about the edges and textures. Together, they provide a more complete description of the image, improving classification accuracy.

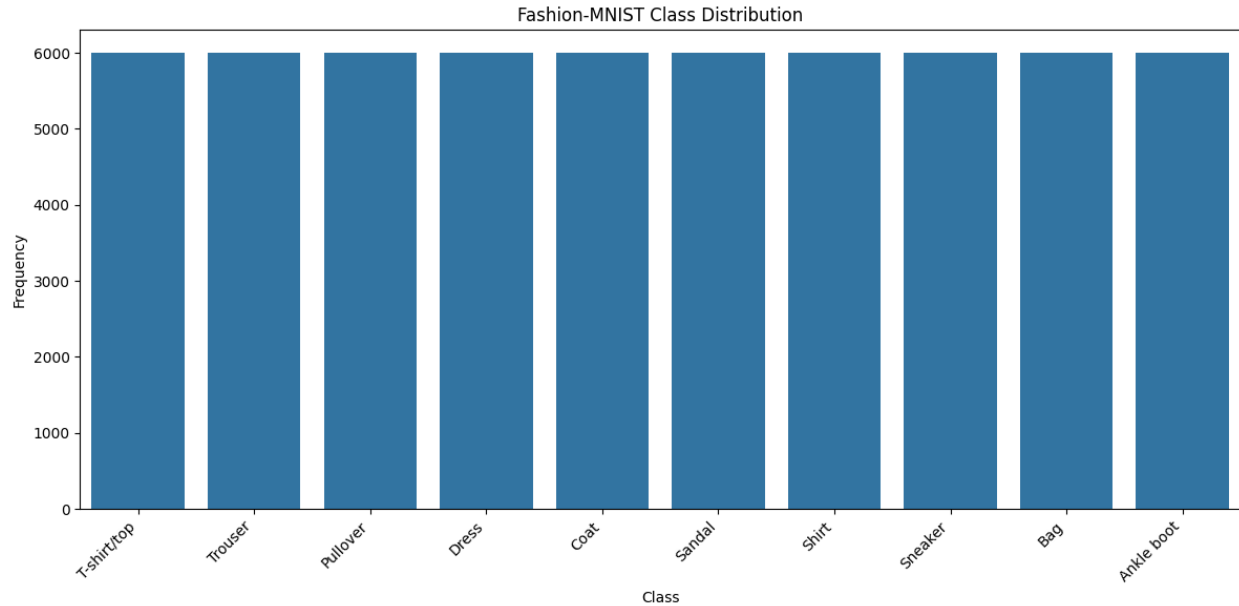
In summary, PCA gives a simpler, more efficient representation of the image, and HOG provides detailed texture and edge information. By combining both, an improvement to the model's ability to recognize objects in the Fashion-MNIST and Fruits-360 datasets is made.

## 2. Visualizations

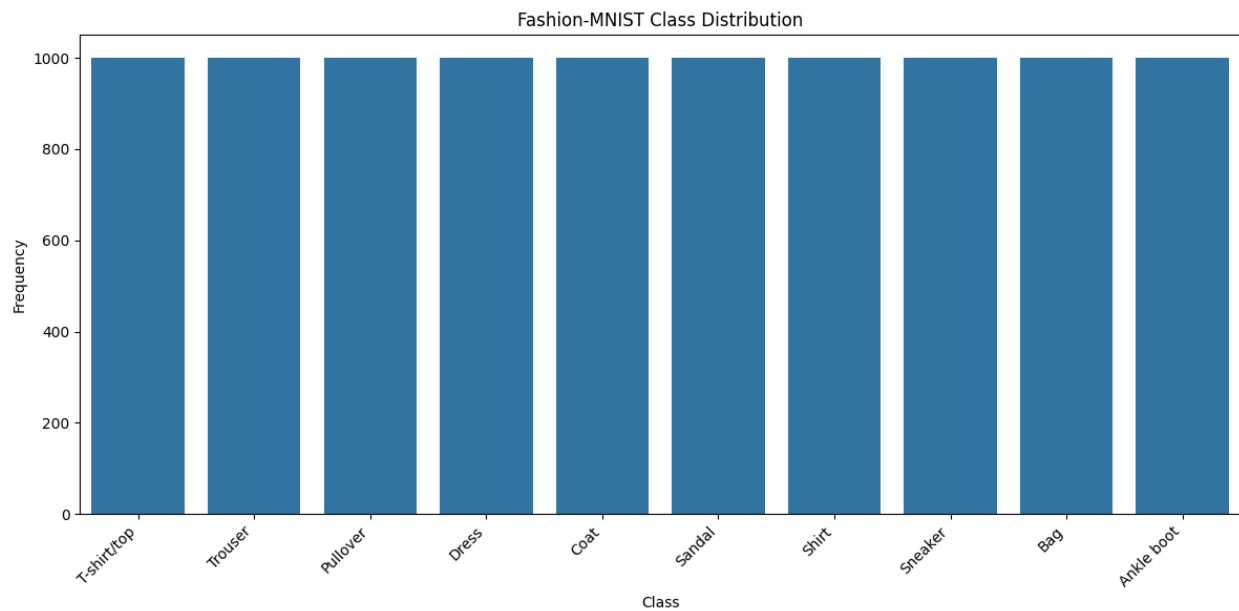
### a. Class Distributions

The Fashion-MNIST dataset is uniformly distributed, while the Fruits-360 dataset shows varying levels of frequency for each class, with some classes being more frequent than others, resulting in an imbalanced distribution across the dataset. There seems to be no difference between the train and test sets of both datasets.

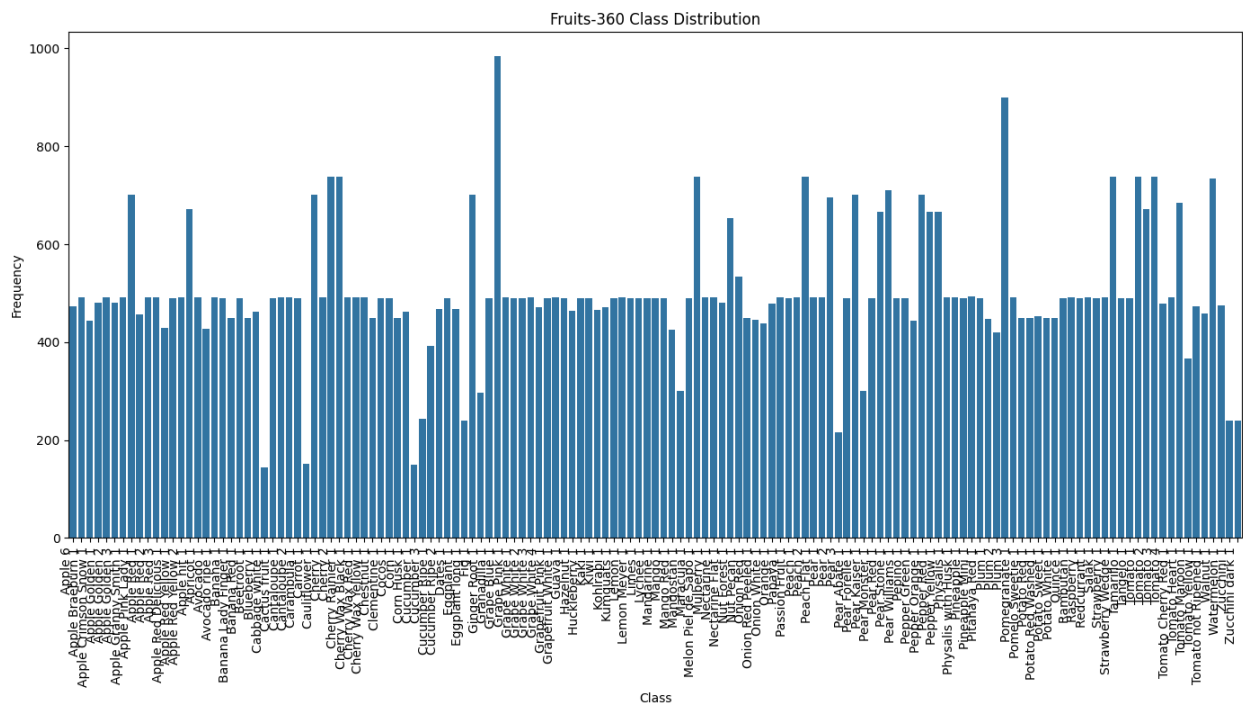
#### Fashion-MNIST Train:



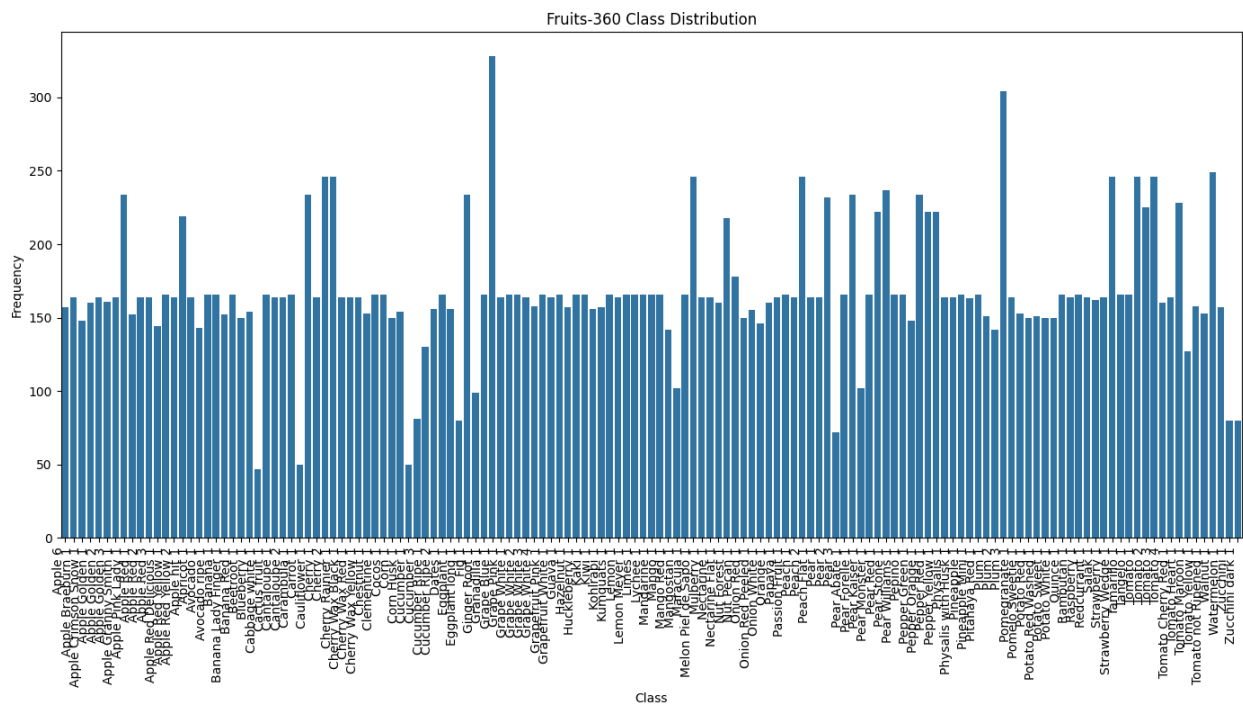
#### Fashion-MNIST Test:



Fruits-360 Train:



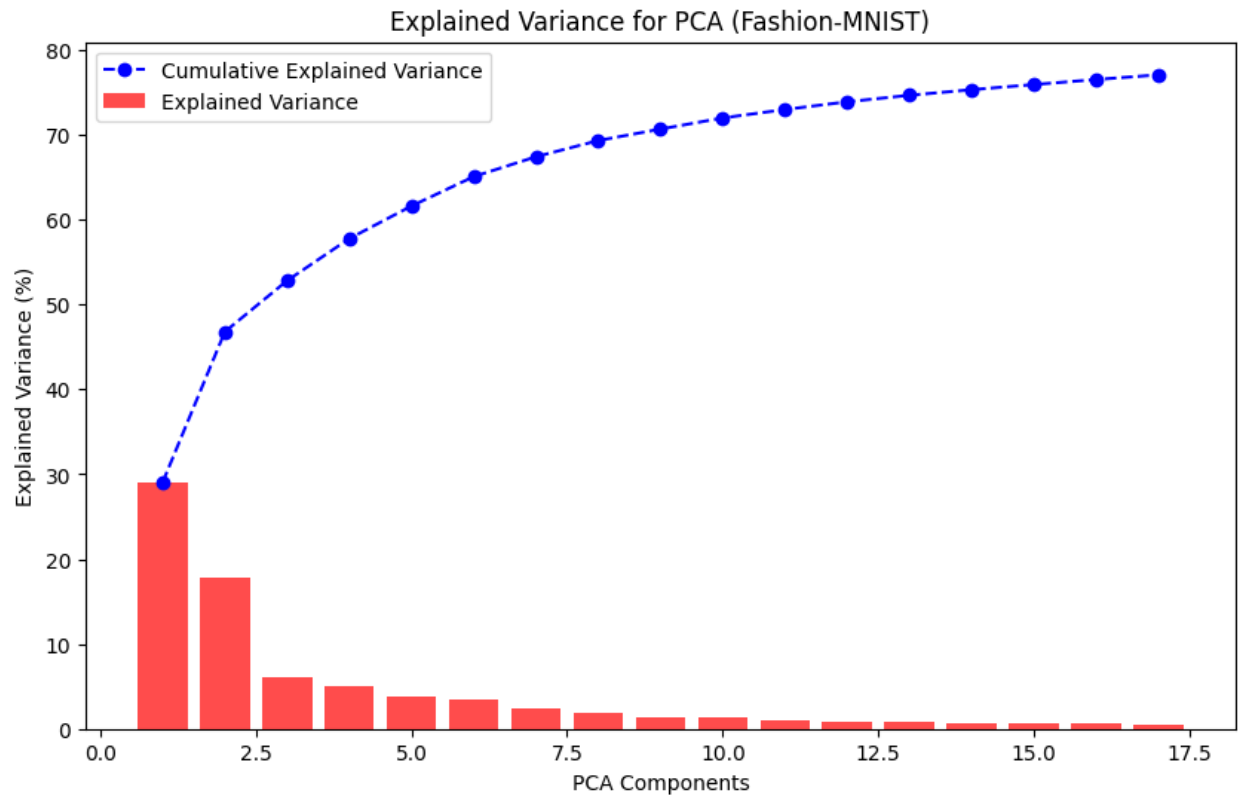
Fruits-360 Test:



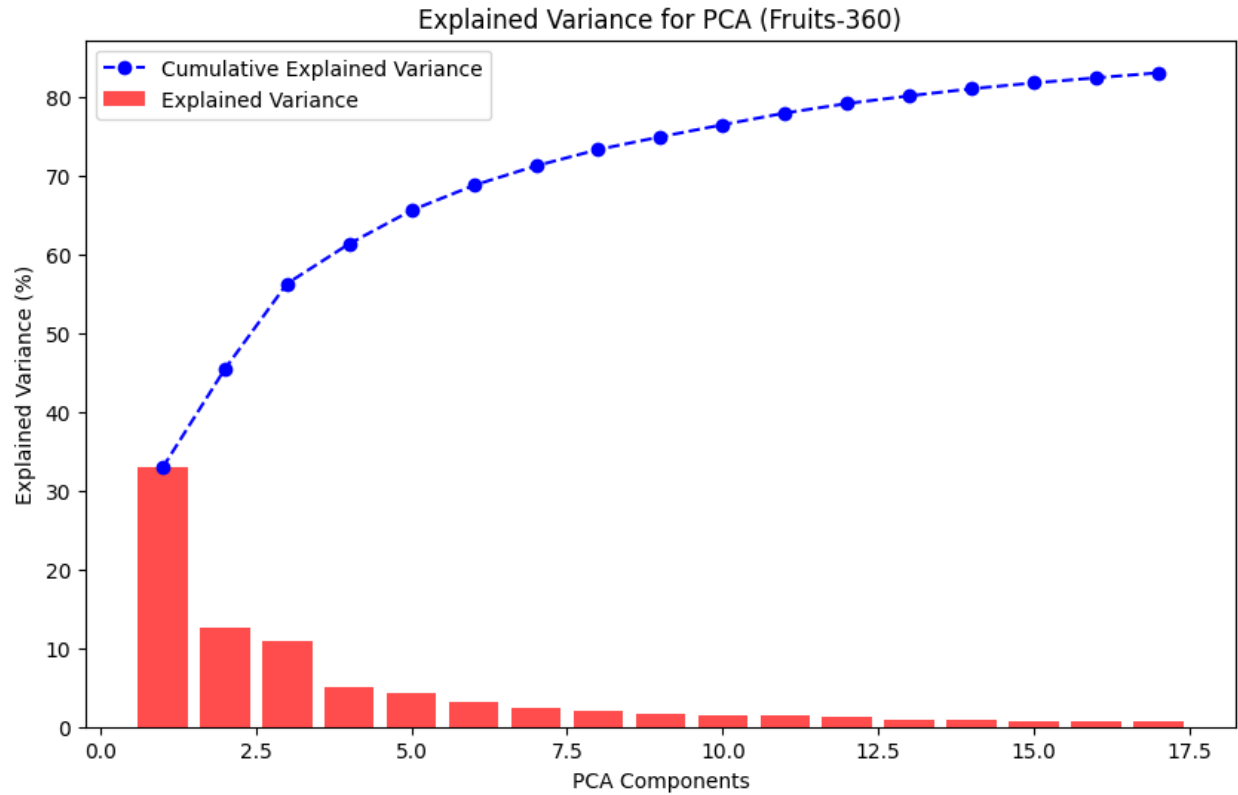
## b. Extracted Features Plots

### 1) Quantitative

#### A. PCA

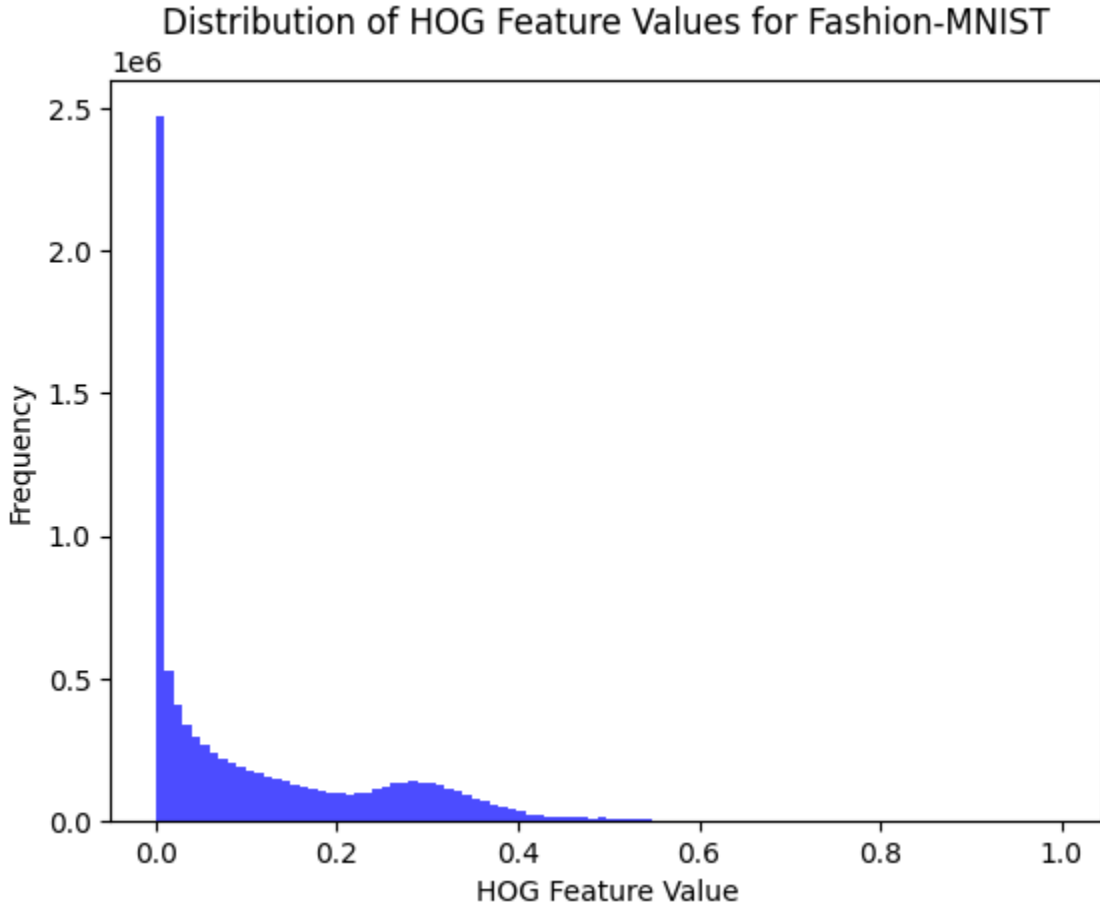


The plot shows the explained variance for the first 17 principal components of the Fashion-MNIST dataset using PCA. The red bars represent the explained variance for each component, where the first few components capture most of the variance (with the first component explaining around 30%). The blue dashed line shows the cumulative explained variance, which increases rapidly with the first few components and then levels off, indicating that most of the data's information is captured by the first components. After the first few components, the contribution to the cumulative variance becomes minimal, suggesting that fewer components are necessary to represent most of the data's structure.

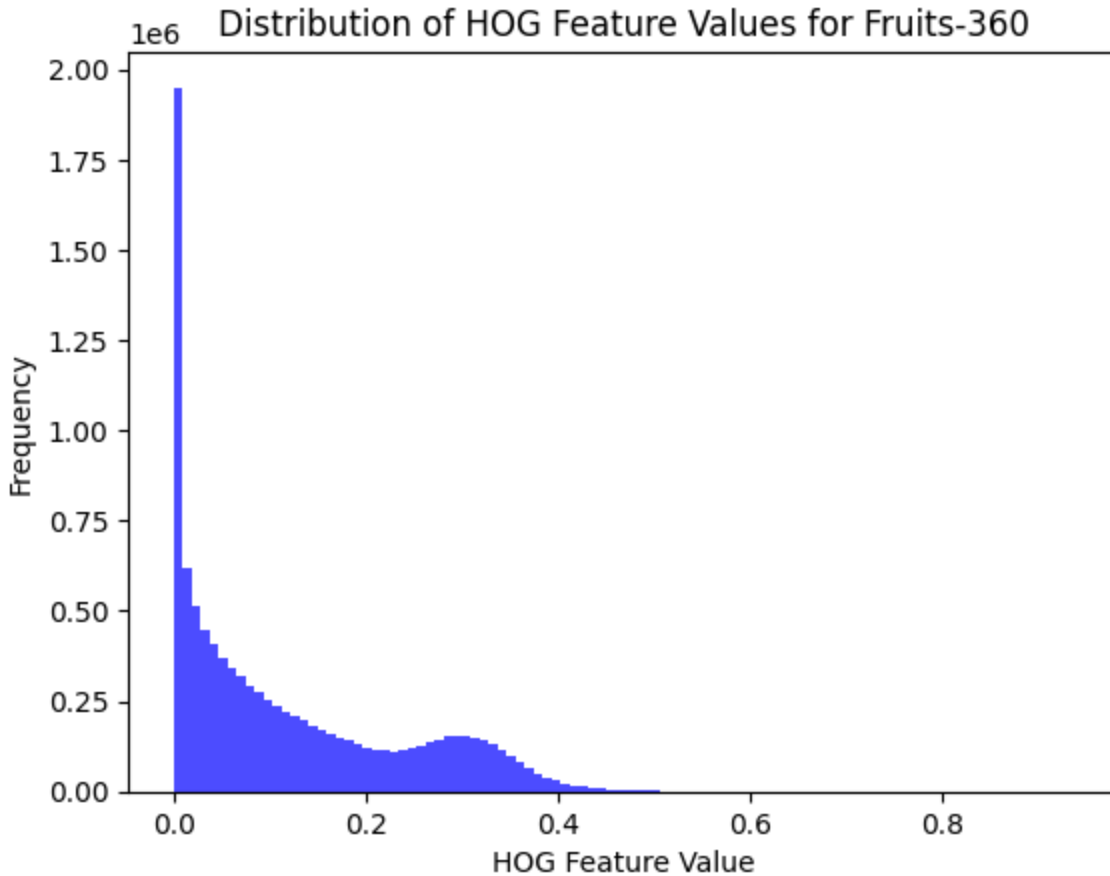


Similarly to the Fashion-MNIST dataset, the plot for Fruits-360 shows the explained variance for the first 17 principal components using PCA. The red bars represent the explained variance for each component, with the first few components capturing most of the variance, particularly the first component which explains around 30%. The blue dashed line shows the cumulative explained variance, which rises steeply at first and then gradually levels off, indicating that a large portion of the data's information is captured by the initial components. This pattern suggests that fewer components can effectively represent most of the data's structure, making PCA an efficient dimensionality reduction method for the Fruits-360 dataset as well.

## B. HOG



The distribution of HOG feature values for Fashion-MNIST shows a sharp peak at 0, indicating that most of the image regions have weak gradients, representing flat or uniform areas. As the values increase, the frequency drops significantly, suggesting that stronger gradients (probably corresponding to edges or textures) are less common, however crucial for distinguishing different classes in the dataset.

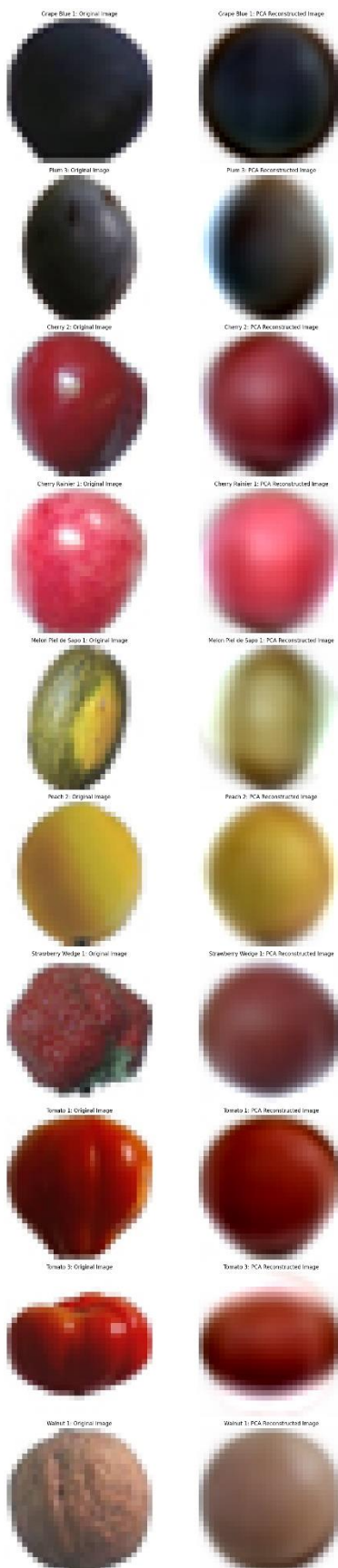
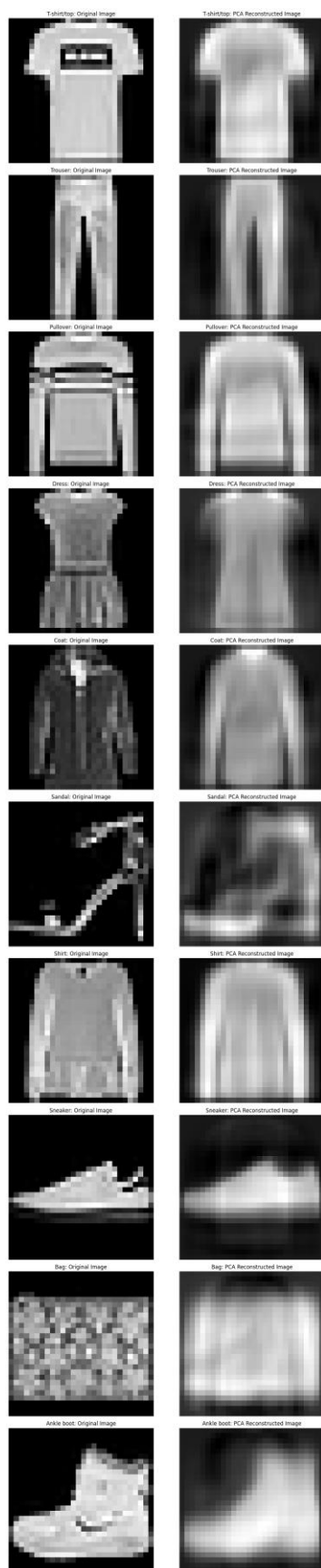


The distribution of HOG feature values for Fruits-360 shows a similar distribution to Fashion-MNIST, with a sharp peak at 0, indicating that most of the image regions have weak gradients, likely representing smooth areas of the fruit. As the values increase, the frequency decreases, with fewer regions showing stronger gradients, which are typically associated with edges and more distinct features, such as the contours of the fruit's shape or color.

## 2) Qualitative

### A. PCA

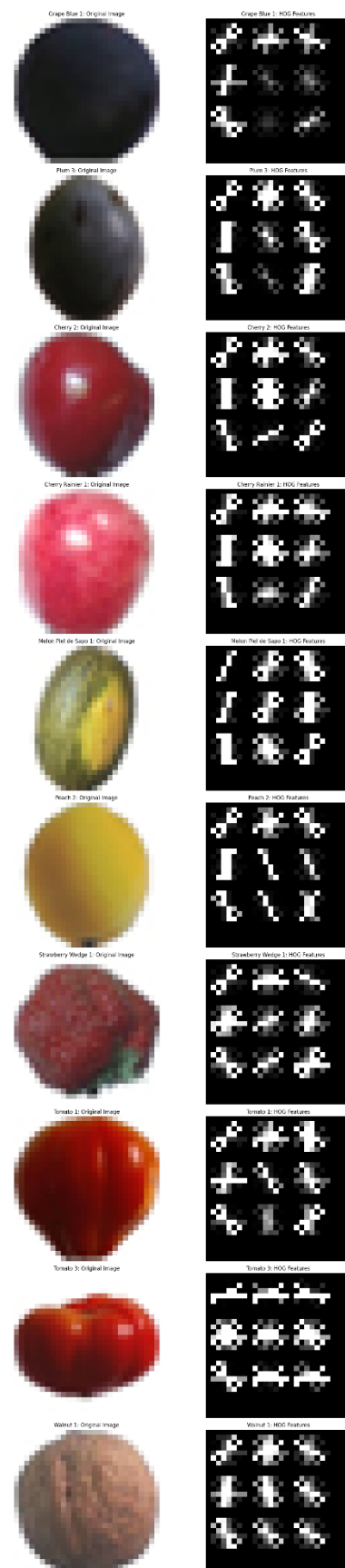
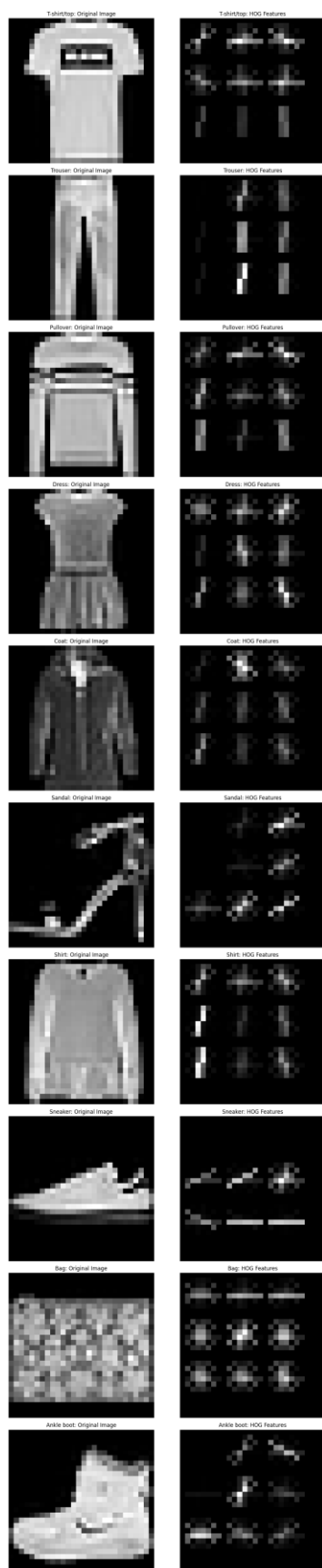
The images show the original items from both the Fashion-MNIST and Fruits-360 datasets on the left, alongside their PCA-reconstructed images on the right. The PCA reconstruction, using the first 17 principal components, captures the overall structure of the images but loses finer details. For example, in the Fashion-MNIST images, the basic shape of clothing is visible, but the intricate textures or small features like seams are not well-represented (making a T-shirt sometimes look like a shirt or dress, more on that later). Similarly, in the Fruits-360 images, the general outline of the fruits is clear, but details like the texture or specific features (e.g., stems or skin patterns) are blurred. This illustrates how PCA retains the most significant patterns but simplifies fine details.





## B. HOG

The images and their corresponding HOG features show how the local gradient information is captured for each object in the datasets. For both Fashion-MNIST and Fruits-360, the original images are paired with their HOG features, which highlight edges and textures. In the HOG features, the distinct edges of clothing items (like the seams of a T-shirt or the contours of a shoe) and the outlines of fruits (like apples or grapes) are visible. The HOG features help emphasize the shape and texture.



### 3. Standardization and Selection

The combined features for both Fashion-MNIST and Fruits-360 datasets are standardized using `StandardScaler`. This step is essential because it ensures that all features are on the same scale, making the machine learning model training more effective.

After standardization, Variance Thresholding is applied to reduce the number of features by removing those with low variance, under the assumption that low-variance features carry less information. However, in both datasets, this method does not significantly reduce the number of features, as the HOG features already exhibit relatively low variance. As a result, the variance thresholding does not contribute much to reducing dimensionality, leaving most of the features intact.

Additionally, there was also the possibility of using `SelectPercentile` with the `f_classif` score function. However, this method leads to a ~15% drop in accuracy when selecting the top 10% of features, suggesting that the selected features may not be as informative as expected. This shows that while variance-based selection methods might not have much impact, feature selection using statistical methods like `SelectPercentile` could potentially hurt the model performance in this case.

### 4. Model Evaluation

After applying Variance Thresholding to the combined features, the number of features remained the same for both datasets. The original feature set contained 161 features (17 PCA + 144 HOG), and after applying the variance threshold, it still had 161 features. The reason is explored in the section above.

Using `RandomSearchCV`, I have searched for hyperparameters for all models. For example, the top 5 for SVM (for the Fashion-MNIST dataset) were:

Parameters	Mean Test Score
{'kernel': 'linear', 'gamma': 'scale', 'C': 1}	0.661829
{'kernel': 'linear', 'gamma': 'scale', 'C': 1000}	0.661772
{'kernel': 'linear', 'gamma': 'auto', 'C': 10}	0.661772
{'kernel': 'linear', 'gamma': 'scale', 'C': 10}	0.661772
{'kernel': 'poly', 'gamma': 'auto', 'C': 100}	0.584103

However, henceforth I will focus on the best hyperparameters and showcase different metrics relating to the runs.

I. Logistic Regression  
a. Fashion-MNIST

Best Parameters: {'solver': 'lbfgs', 'penalty': 'l2', 'multi\_class': 'multinomial', 'C': 10}

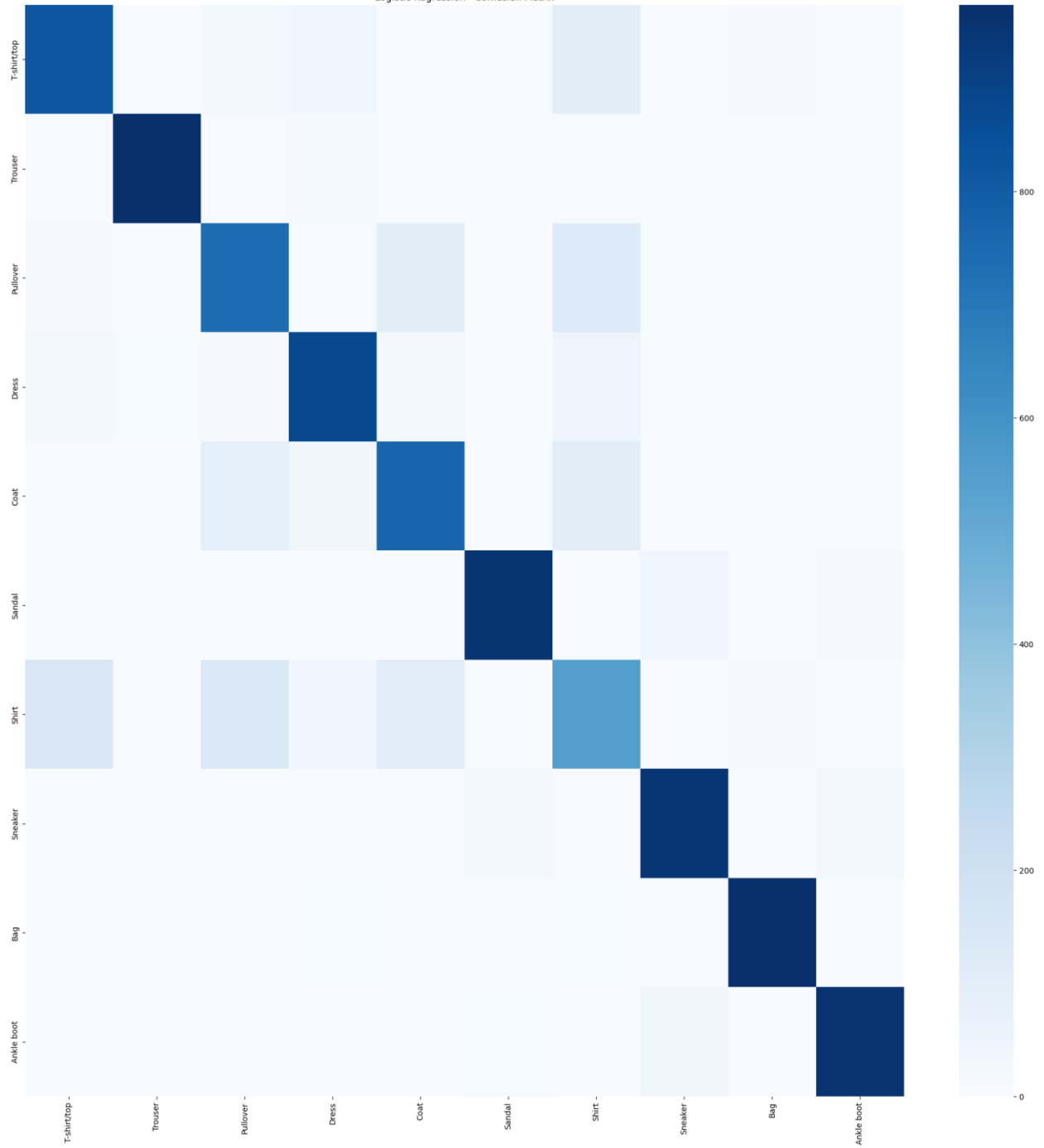
Class	Precision	Recall	F1-Score	Support
T-shirt/top	0.81	0.82	0.81	1000
Trouser	0.97	0.96	0.97	1000
Pullover	0.74	0.74	0.74	1000
Dress	0.86	0.88	0.87	1000
Coat	0.76	0.77	0.76	1000
Sandal	0.96	0.95	0.95	1000
Shirt	0.58	0.55	0.56	1000
Sneaker	0.92	0.95	0.93	1000
Bag	0.96	0.96	0.96	1000
Ankle boot	0.96	0.95	0.96	1000

**Accuracy:** 0.8524

**Macro Average:** Precision = 0.85, Recall = 0.85, F1-Score = 0.85

**Weighted Average:** Precision = 0.85, Recall = 0.85, F1-Score = 0.85

### Logistic Regression - Confusion Matrix



b. Fruits-360

Best Parameters: {'solver': 'newton-cg', 'penalty': 'l2', 'multi\_class': 'multinomial', 'C': 1}

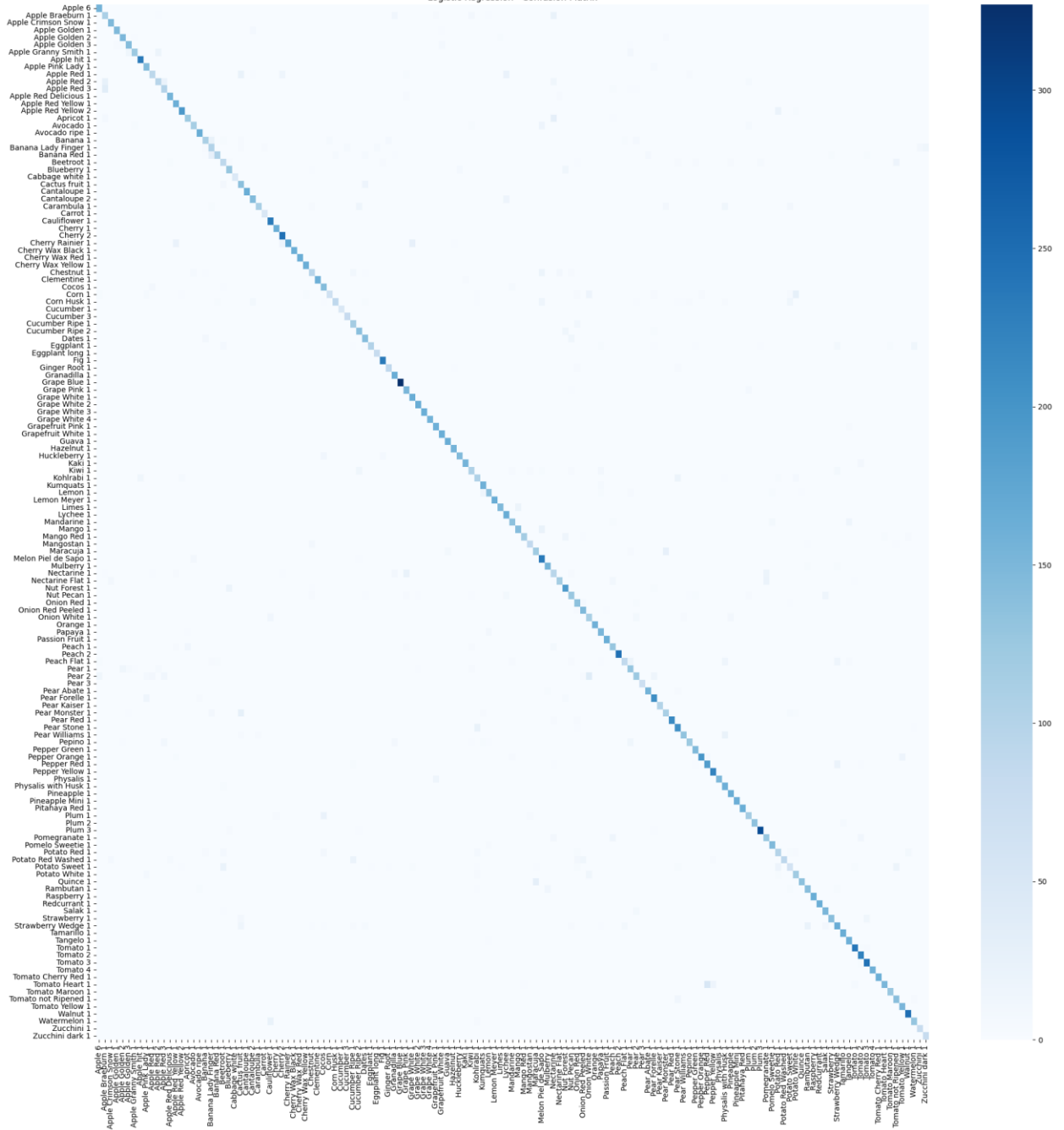
Class Name	Precision	Recall	F1-Score	Support
Apple hit 1	0.94	1.00	0.97	234
Mango Red 1	0.94	1.00	0.98	164
Apple Granny Smith 1	0.94	1.00	0.97	164
Apple Red Yellow 1	0.83	0.99	0.90	164
Pear Williams 1	0.94	1.00	0.97	164
Banana Red 1	0.81	0.56	0.66	164
Apple Red 1	0.81	0.56	0.66	164
Pear Stone 1	0.83	0.79	0.81	164
Banana Lady Finger 1	0.60	0.68	0.64	152
Pear Abate 1	0.92	0.95	0.93	164

**Accuracy:** 0.8729

**Macro Average:** Precision = 0.88, Recall = 0.87, F1-Score = 0.87

**Weighted Average:** Precision = 0.88, Recall = 0.87, F1-Score = 0.87

Logistic Regression - Confusion Matrix



## II. SVM

### a. Fashion-MNIST

Best parameters: {'kernel': 'rbf', 'gamma': 'auto', 'C': 10}

Class Name	Precision	Recall	F1-Score	Support
T-shirt/top	0.83	0.86	0.84	1000
Trouser	0.99	0.97	0.98	1000
Pullover	0.80	0.82	0.81	1000
Dress	0.89	0.90	0.89	1000
Coat	0.82	0.82	0.82	1000
Sandal	0.97	0.96	0.97	1000
Shirt	0.72	0.69	0.70	1000
Sneaker	0.94	0.96	0.95	1000
Bag	0.98	0.98	0.98	1000
Ankle boot	0.96	0.96	0.96	1000

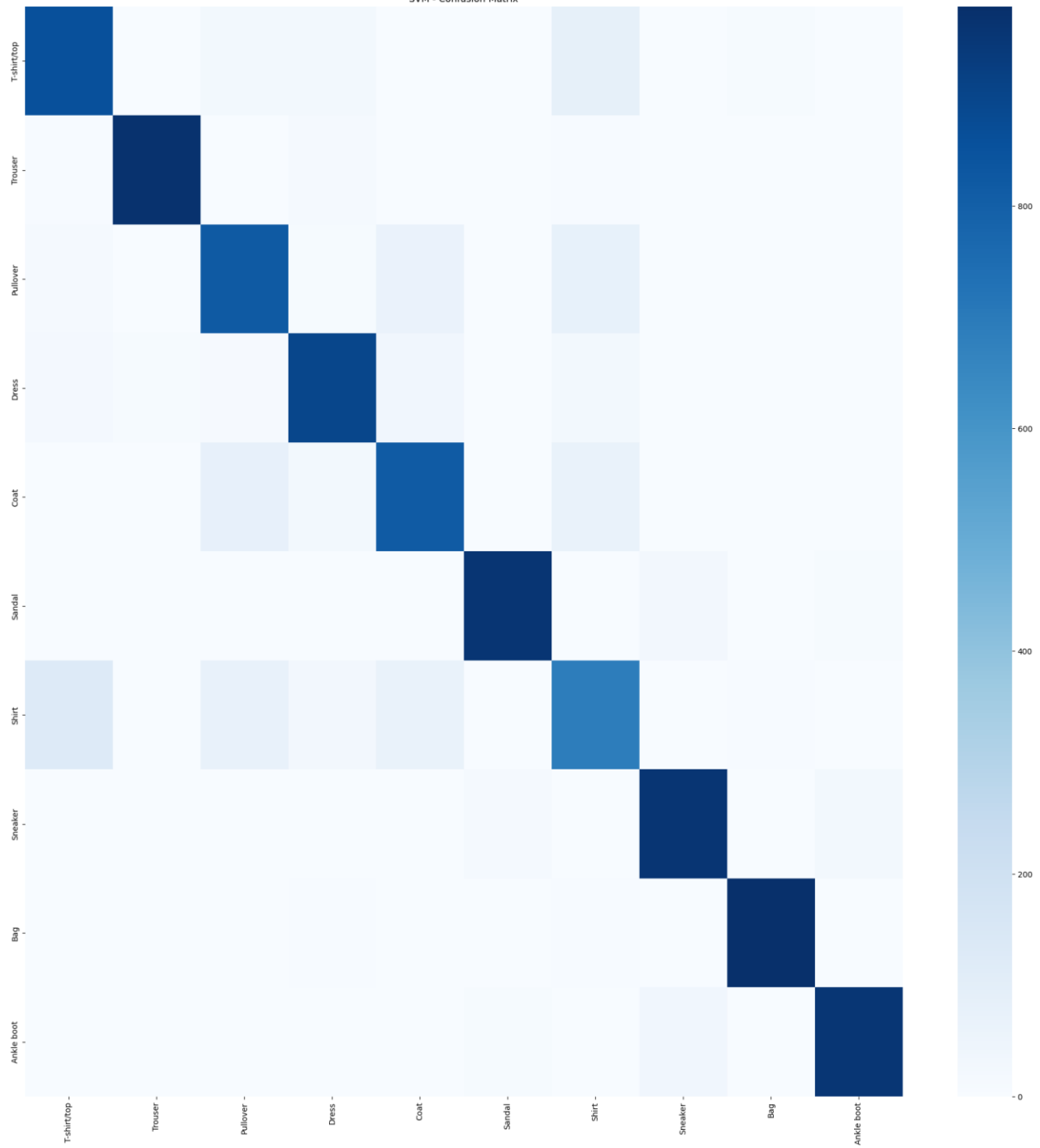
**Accuracy: 0.89**

**Macro Average:** Precision = **0.89**, Recall = 0.89, F1-Score = **0.89**

**Weighted Average:** Precision = 0.89, Recall = 0.89, F1-Score = 0.89



### SVM - Confusion Matrix



b. Fruits-360

Best parameters: {'kernel': 'linear', 'gamma': 'scale', 'C': 1}

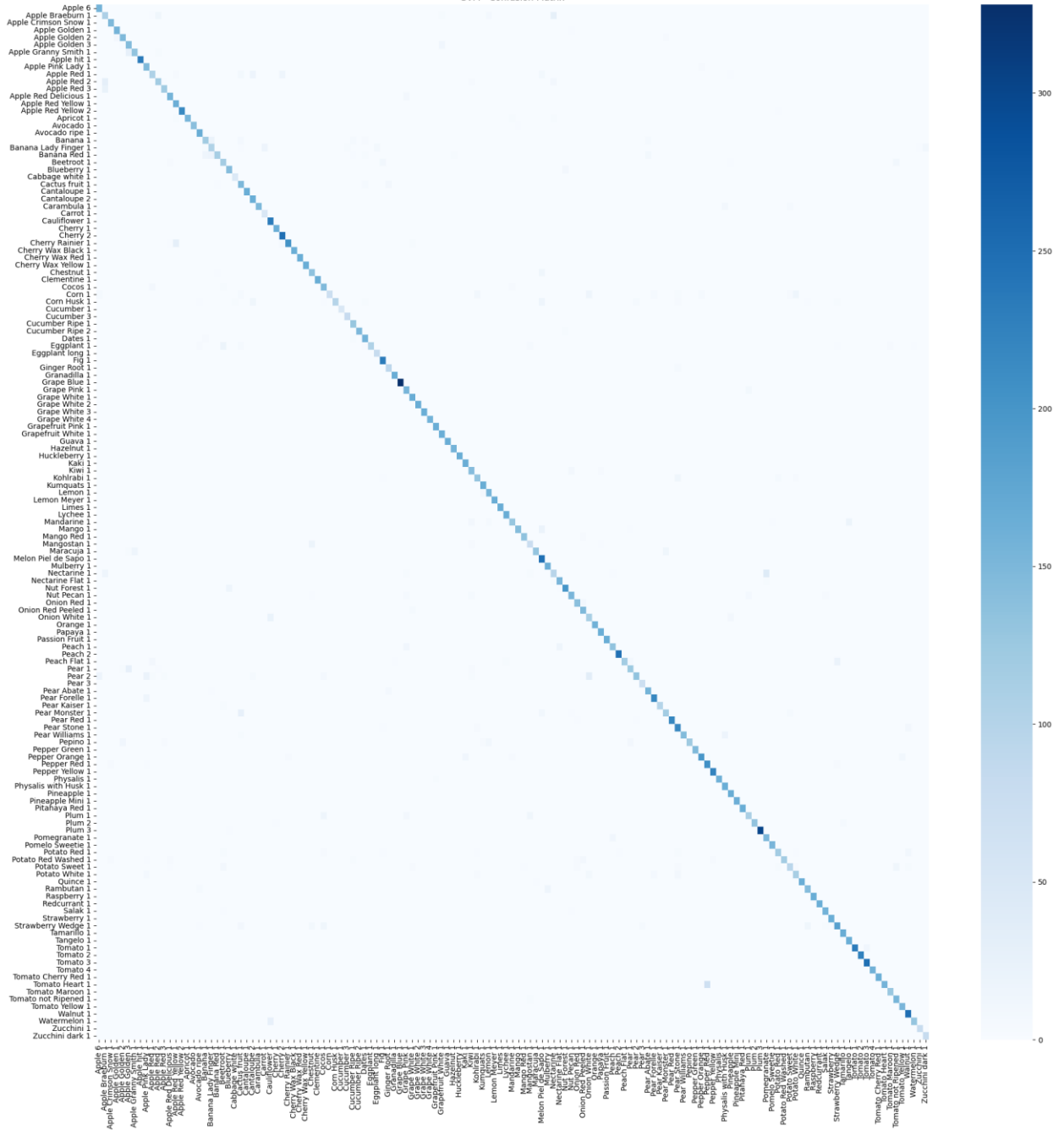
Class Name	Precision	Recall	F1-Score	Support
Apple hit 1	1.00	1.00	1.00	234
Mango Red 1	1.00	1.00	1.00	164
Pear Williams 1	1.00	1.00	1.00	164
Apple Red Yellow 1	1.00	0.99	0.99	219
Apple Red Delicious 1	1.00	1.00	1.00	164
Banana Red 1	0.87	1.00	0.93	157
Pear Stone 1	0.85	0.88	0.86	164
Pear Abate 1	0.92	0.95	0.93	164
Potato Red Washed 1	0.78	0.79	0.78	164
Banana Lady Finger 1	0.71	0.90	0.80	161

**Accuracy: 0.92**

**Macro Average:** Precision = **0.92**, Recall = **0.92**, F1-Score = **0.91**

**Weighted Average:** Precision = 0.92, Recall = 0.92, F1-Score = 0.91

SVM - Confusion Matrix



### III. Random Forest

#### a. Fashion-MNIST

Best parameters: {'n\_estimators': 500, 'min\_samples\_split': 5, 'min\_samples\_leaf': 1, 'max\_features': 'sqrt', 'max\_depth': None, 'criterion': 'entropy'}

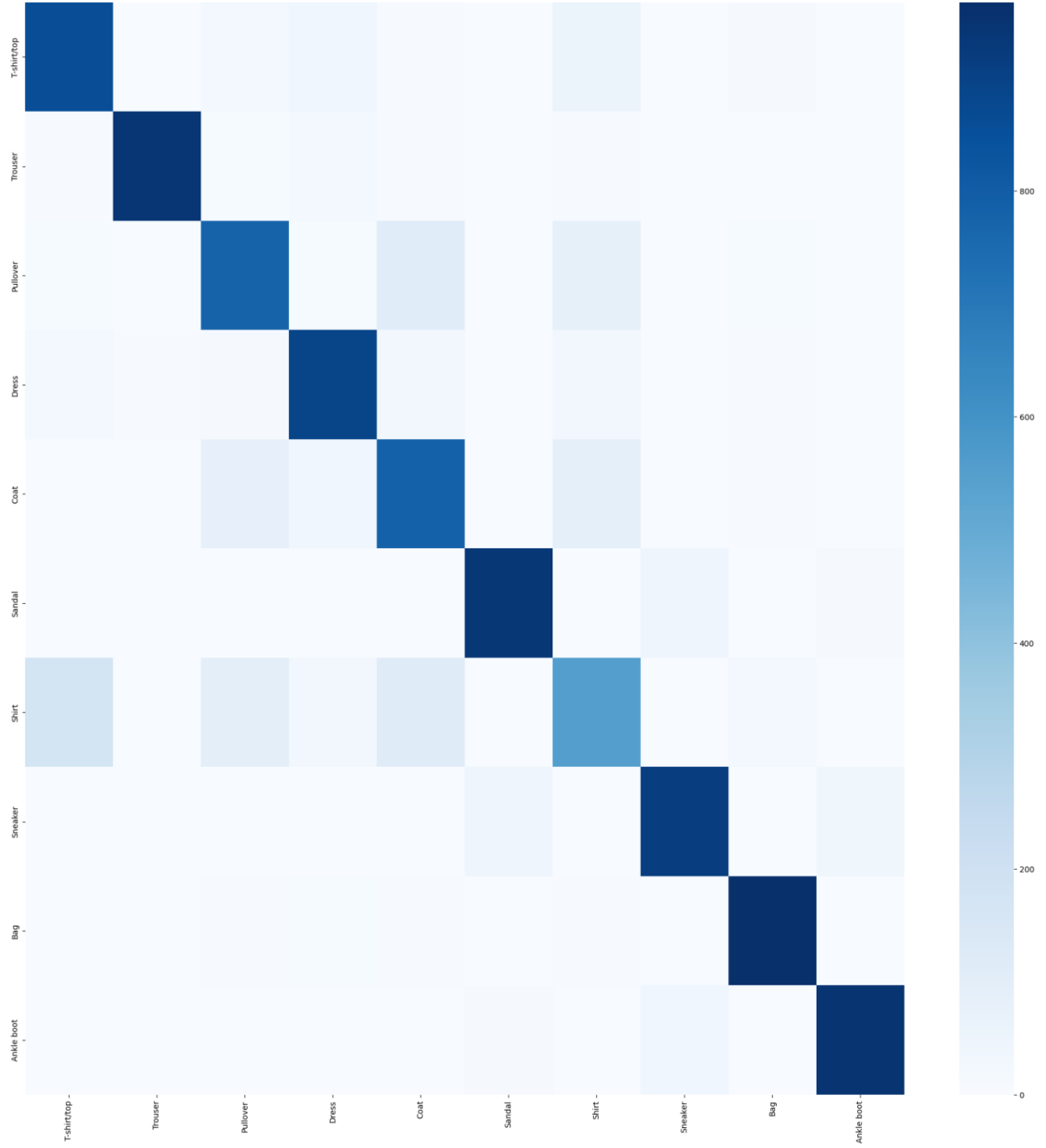
Class Name	Precision	Recall	F1-Score	Support
T-shirt/top	0.80	0.86	0.83	1000
Trouser	0.99	0.95	0.97	1000
Pullover	0.77	0.78	0.77	1000
Dress	0.86	0.89	0.87	1000
Coat	0.74	0.79	0.76	1000
Sandal	0.94	0.94	0.94	1000
Shirt	0.67	0.55	0.61	1000
Sneaker	0.91	0.92	0.92	1000
Bag	0.94	0.97	0.96	1000
Ankle boot	0.94	0.95	0.95	1000

**Accuracy:** 0.86

**Macro Average:** Precision = 0.86, Recall = 0.86, F1-Score = 0.86

**Weighted Average:** Precision = 0.86, Recall = 0.86, F1-Score = 0.86

### RF - Confusion Matrix



b. Fruits-360

Best parameters: {'n\_estimators': 200, 'min\_samples\_split': 5, 'min\_samples\_leaf': 1, 'max\_features': 'sqrt', 'max\_depth': 30, 'criterion': 'log\_loss'}

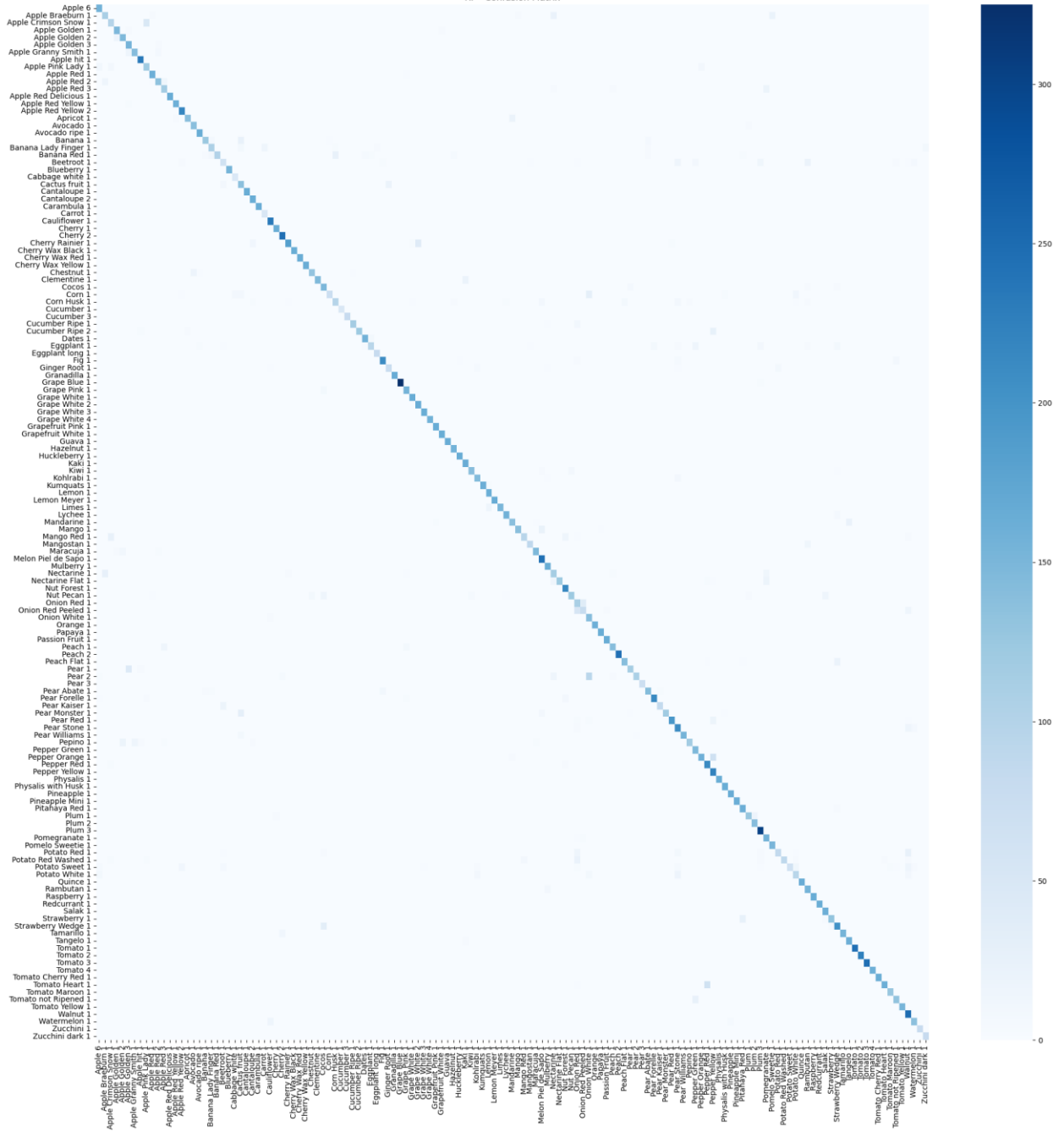
Class Name	Precision	Recall	F1-Score	Support
Apple hit 1	1.00	1.00	1.00	234
Mango Red 1	1.00	1.00	1.00	164
Pear Williams 1	1.00	1.00	1.00	164
Apple Red Yellow 1	1.00	0.99	0.99	219
Apple Red Delicious 1	1.00	1.00	1.00	164
Banana Red 1	0.84	1.00	0.91	157
Pear Stone 1	0.79	0.78	0.79	99
Pear Abate 1	0.91	1.00	0.95	164
Potato Red Washed 1	0.64	0.64	0.64	154
Banana Lady Finger 1	0.73	0.68	0.70	164

**Accuracy:** 0.90

**Macro Average:** Precision = 0.91, Recall = 0.89, F1-Score = 0.89

**Weighted Average:** Precision = 0.91, Recall = 0.90, F1-Score = 0.89

RF - Confusion Matrix



IV. Gradient Boosted Trees  
a. Fashion-MNIST

Best Parameters: {'subsample': 0.9, 'n\_estimators': 400, 'max\_depth': 4, 'learning\_rate': 0.1, 'colsample\_bytree': 0.8}

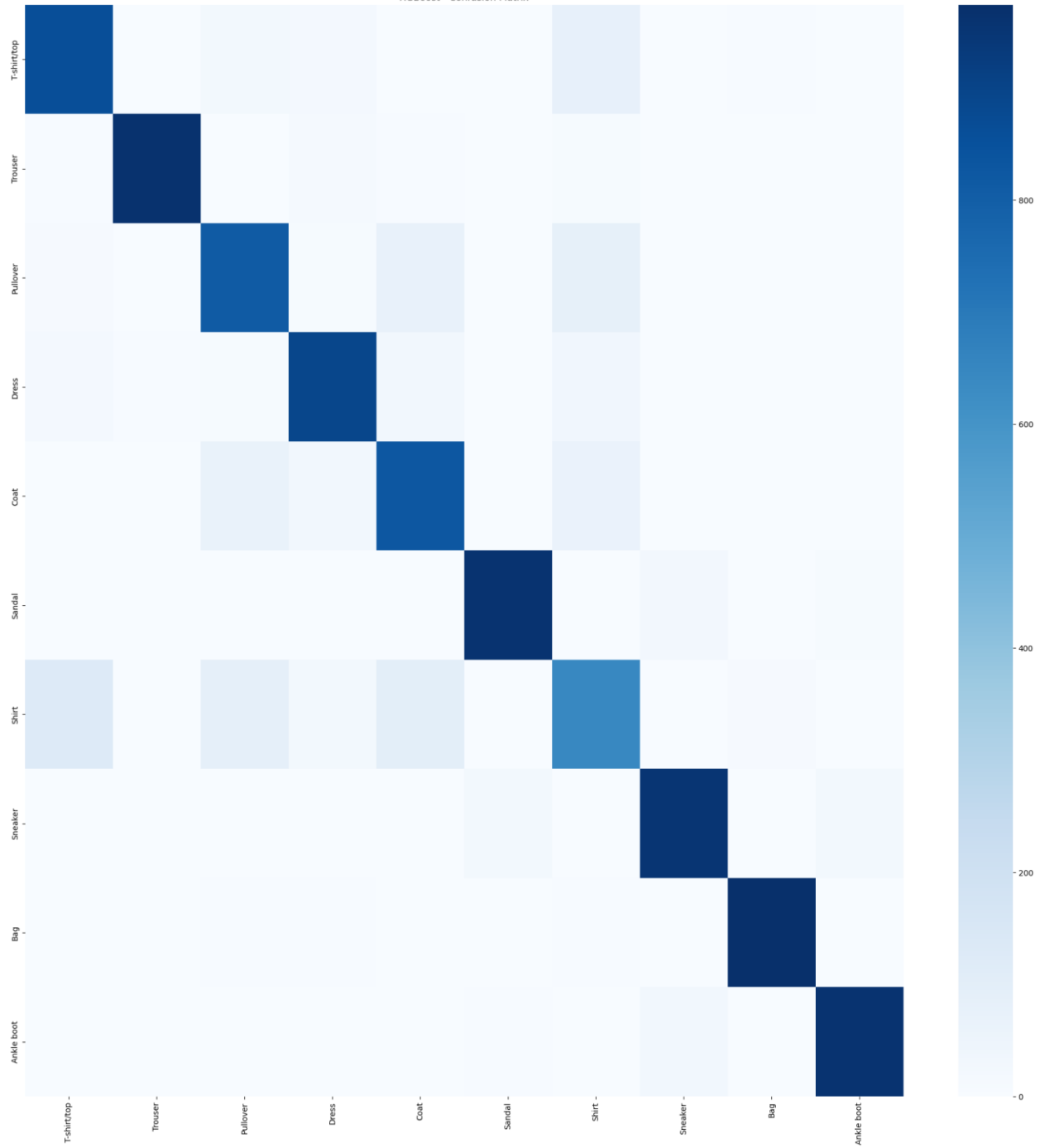
Class Name	Precision	Recall	F1-Score	Support
T-shirt/top	0.83	0.86	0.85	1000
Trouser	0.99	0.96	0.98	1000
Pullover	0.80	0.81	0.81	1000
Dress	0.89	0.89	0.89	1000
Coat	0.79	0.83	0.81	1000
Sandal	0.97	0.96	0.97	1000
Shirt	0.69	0.64	0.67	1000
Sneaker	0.94	0.95	0.95	1000
Bag	0.97	0.97	0.97	1000
Ankle boot	0.97	0.96	0.96	1000

**Accuracy: 0.89**

**Macro Average:** Precision = **0.89**, Recall = **0.89**, F1-Score = 0.88

**Weighted Average:** Precision = 0.89, Recall = 0.89, F1-Score = 0.88



[illegible]

b. Fruits-360

Best Parameters: {'subsample': 0.9, 'n\_estimators': 400, 'max\_depth': 5, 'learning\_rate': 0.05, 'colsample\_bytree': 0.9}

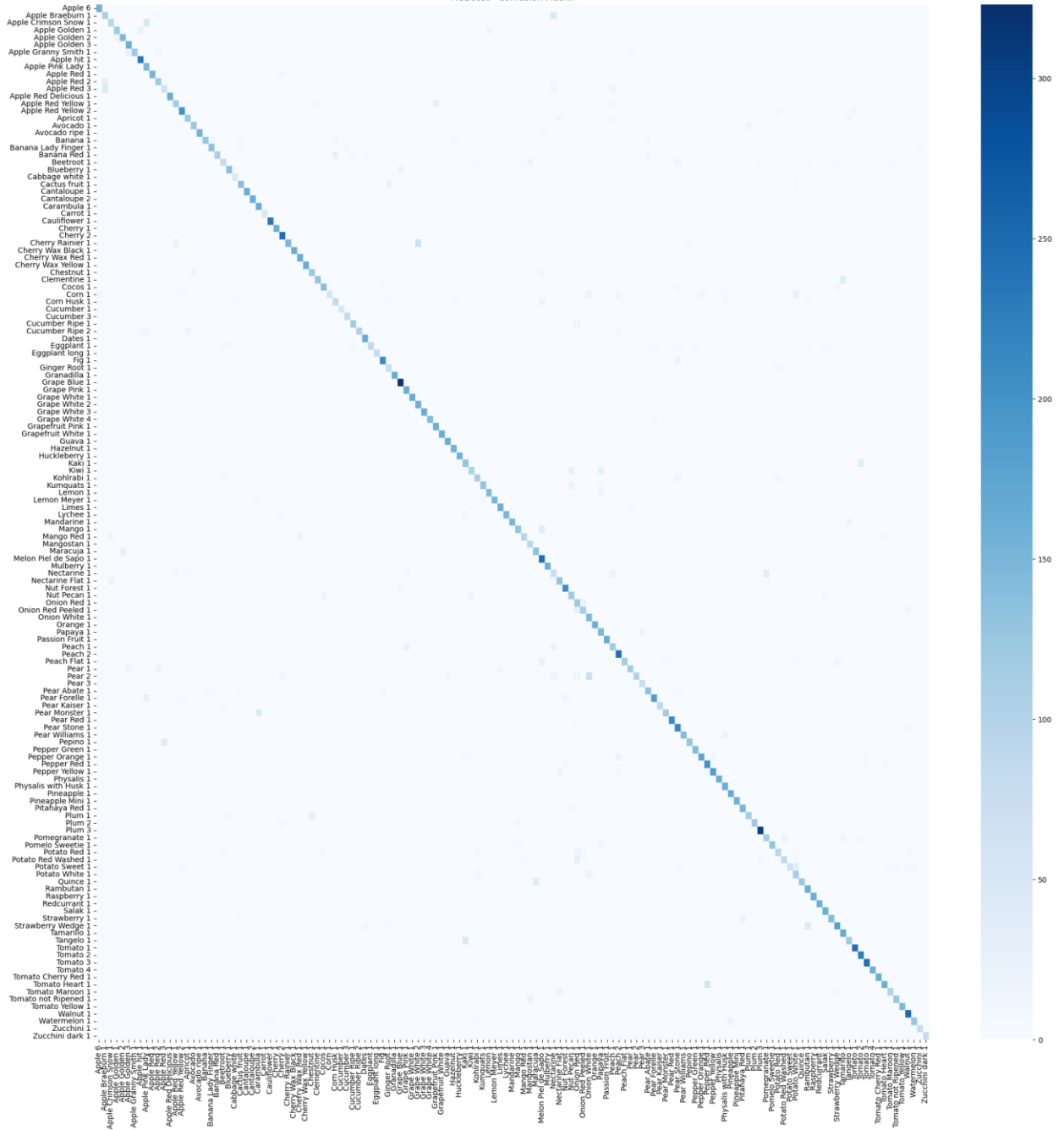
Class Name	Precision	Recall	F1-Score	Support
Physalis with Husk 1	1.00	1.00	1.00	80
Apple 6	0.88	1.00	0.93	157
Pineapple 1	1.00	1.00	1.00	164
Pear Red 1	1.00	1.00	1.00	164
Apple hit 1	0.88	1.00	0.94	234
Potato Sweet 1	0.39	0.45	0.42	164
Potato Red Washed 1	0.56	0.76	0.64	150
Beetroot 1	0.64	0.54	0.59	150
Tomato Heart 1	0.72	0.57	0.66	151
Potato Red 1	0.71	0.62	0.66	150

**Accuracy:** 0.86

**Macro Average:** Precision = 0.87, Recall = 0.86, F1-Score = 0.86

**Weighted Average:** Precision = 0.87, Recall = 0.86, F1-Score = 0.86

XGBoost - Confusion Matrix



## V. General Trends

- Overall, SVM got the best results, however it took the longest time to search for hyperparameter (around 8 hours for Fashion-MNIST).
- Logistic Regression performed the worst but was easily the fastest to calibrate and train.
- Random Forest performed well on Fruits-360 and took a lot less time to calibrate, although the training was longer.
- Gradient Boosted Trees got the same accuracy as SVM for the first dataset, so it is preferred as it is much faster to do cross-validation on.
- All models struggle with the T-shirt class in the Fashion-MNIST dataset. It is usually confused with all items of clothing that go on the upper body. As such, all accuracy is below 90% for this dataset.
- The models that performed the best in 1 metric (accuracy), performed the best or almost the best in all other metrics as well (precision, recall, f1 score).

## 5. Conclusions

The combination of PCA and HOG proved highly effective in extracting meaningful features, capturing both the overall structure and fine details of images. SVM achieved the best accuracy but was computationally expensive, while Gradient Boosted Trees offered comparable performance with faster training, making it a more practical choice. Logistic Regression was the fastest but least accurate, and Random Forest balanced accuracy and training time well, especially for Fruits-360. Challenges remained with specific classes, such as T-shirts in Fashion-MNIST, which were frequently confused with similar items. Overall, the models demonstrated strong performance across metrics, with accuracies exceeding 80% for both datasets.