

# Handwritten-Digit-Recognition-with-Logistic-Regression

Silent0Wings

August 2, 2025

## Contents

<b>1</b>	<b>Introduction and Data Preparation</b>	<b>2</b>
1.1	Training Setup and Evaluation Approach . . . . .	5
1.2	What is Logistic Regression? . . . . .	5
<b>2</b>	<b>Parameter Tuning &amp; Testing</b>	<b>6</b>
	Final Recommendations . . . . .	10
<b>3</b>	<b>Results</b>	<b>11</b>
3.1	Model and Evaluation . . . . .	11
3.2	Classification Report . . . . .	12
3.3	Confusion Matrix . . . . .	13
<b>4</b>	<b>Sample Prediction</b>	<b>14</b>
<b>5</b>	<b>MNIST Logistic Regression: Small vs Large Dataset</b>	<b>15</b>
<b>6</b>	<b>Project Code</b>	<b>17</b>
6.1	Sklearn Dataset Code . . . . .	17
6.2	MNIST Dataset Code . . . . .	18
6.3	Utility Functions . . . . .	19
6.3.1	display_prediction_image Function . . . . .	19
6.3.2	plot_classification_report Function . . . . .	19
6.3.3	plot_confusion_matrix Function . . . . .	19
	<b>References</b>	<b>20</b>

# 1 Introduction and Data Preparation

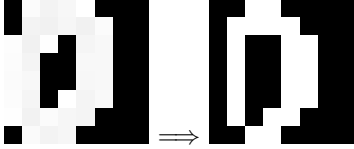
This report implements and evaluates a digit classifier using scikit-learn’s logistic regression on two different datasets:

- The “**load-digits**” dataset from scikit-learn [1], consisting of 1,797 grayscale images sized  $8 \times 8$  pixels. Pixel values range from 0 to 16 and are normalized by dividing by 16 to scale between 0 and 1 for better performance.
- The **MNIST** dataset in CSV format [2] containing 70,000 grayscale images sized  $28 \times 28$  pixels, normalized by dividing pixel values by 255.

These two datasets are distinct: **load-digits** is a smaller, low-resolution dataset often used for quick prototyping, while MNIST is a larger, widely used benchmark for handwritten digit classification.

```
1 from sklearn.datasets import load_digits
2 digits = load_digits()
3 x = digits.data / 16.0 # Normalize pixels to 0-1
4 y = digits.target
```

Normalization converts pixel values where 16 = white, 0 = black, and 8 = gray:



We export the dataset to CSV for easy inspection and manipulation:

```
1 import pandas as pd
2 data = pd.DataFrame(digits.data)
3 data['label'] = digits.target
4 data.to_csv('digits_dataset.csv', index=False)
```

This visualizes pixel intensities as digits from 0 to 16:

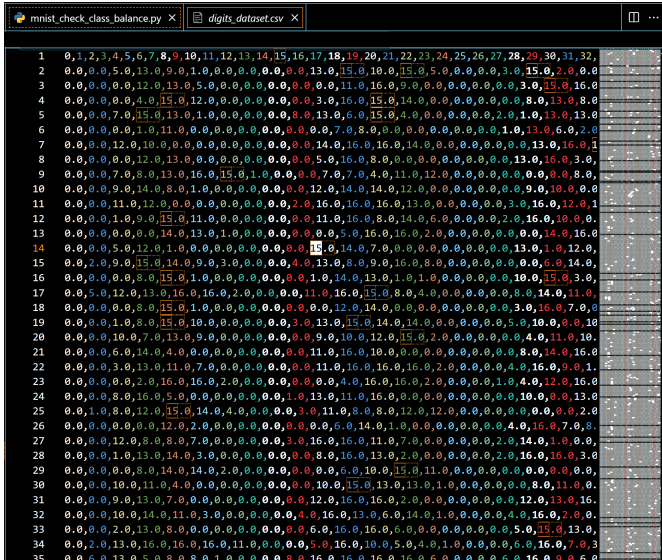


Figure 1: image of scikit-learn load-digits dataset ( $8 \times 8$ )

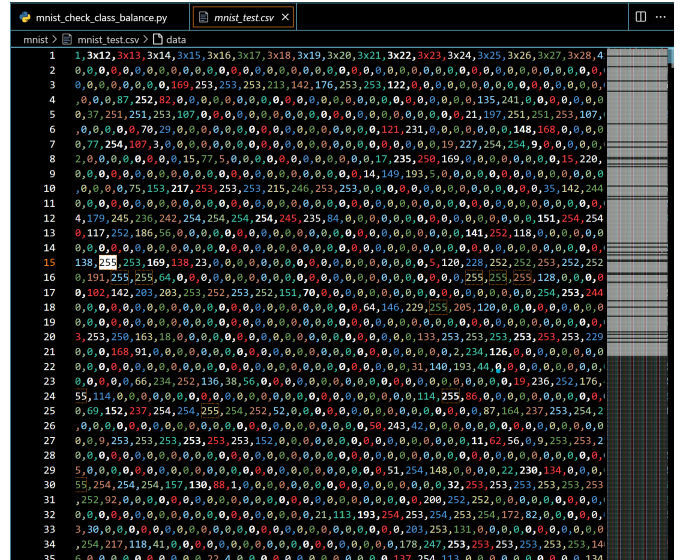


Figure 2: image of MNIST dataset ( $28 \times 28$ )

To convert CSV rows back into images saved by label folder:

```
1 import numpy as np
2 import os
3 import cv2
4
5 df = pd.read_csv('digits_dataset.csv')
6 output_dir = 'processed_images'
7 os.makedirs(output_dir, exist_ok=True)
8
9 for idx, row in df.iterrows():
10     label = row['label']
11     pixels = np.array(row[1:]).reshape(8, 8).astype(np.uint8) * 255
12     label_folder = os.path.join(output_dir, str(label))
13     os.makedirs(label_folder, exist_ok=True)
14     filename = os.path.join(label_folder, f'image_{idx}.png')
15     cv2.imwrite(filename, pixels)
```

All images are saved in folders by digit, easy to navigate:

```
\COMP472\Project1> cd .\processed_images\
PS C:\Users\ypers\OneDrive\Documents\COMP472\Project1\processed_images> tree
Folder PATH listing for volume Windows-SSD
Volume serial number is FE20-4359
C:.\
|—0.0
|—1.0
|—2.0
|—3.0
|—4.0
|—5.0
|—6.0
|—7.0
|—8.0
|—9.0
```

Here are four sample images per digit (0–9), showing balance and handwriting variation:

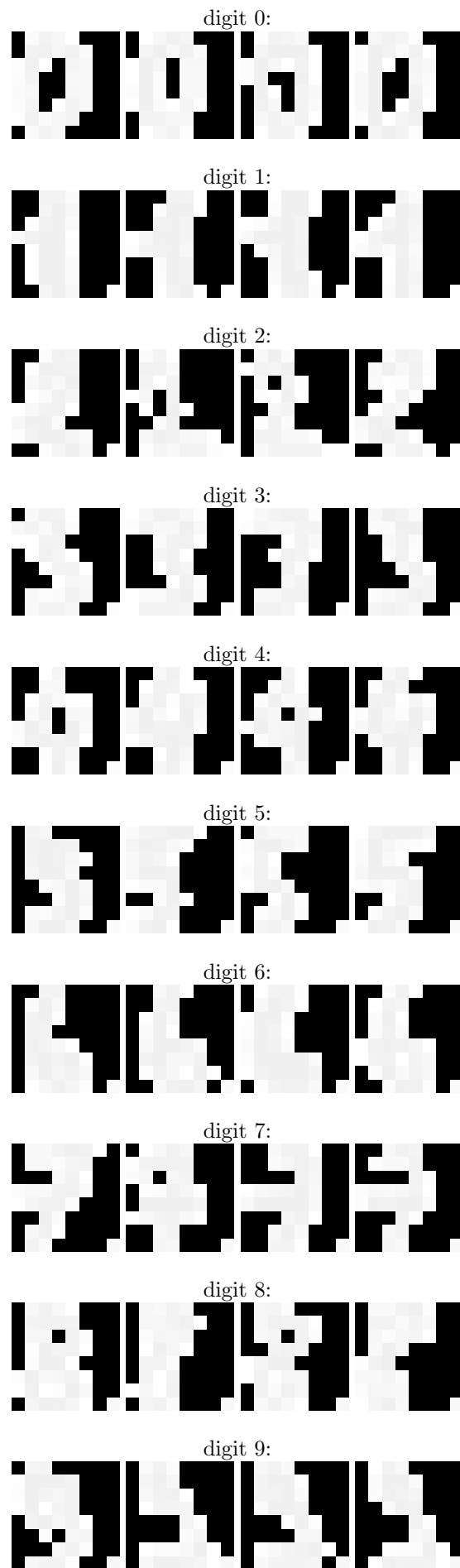


Figure 3: Four sample images per digit (0-9); shows dataset balance and handwriting variation.

We confirm dataset balance with this class distribution plot:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 df = pd.read_csv('digits_dataset.csv')
5 label_counts = df['label'].value_counts().sort_index()
6
7 ax = label_counts.plot(kind='bar', color='black', edgecolor='gray', title='Class Distribution
8 ')
9 plt.xlabel('Digit')
10 plt.ylabel('Count')
11
12 for idx, val in enumerate(label_counts):
13     ax.text(idx, val + 0.5, str(val), ha='center', va='bottom')
14
15 plt.show()
```

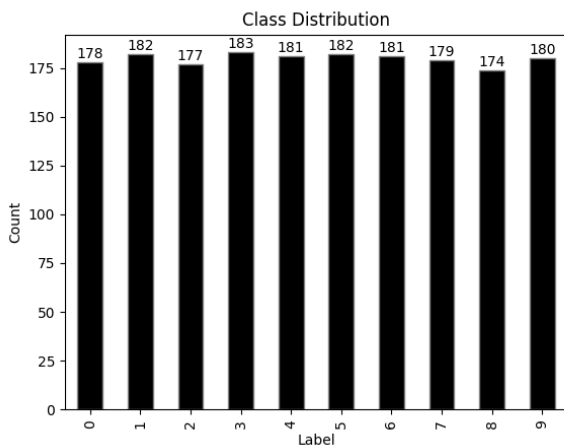


Figure 4: Digit distribution for scikit-learn dataset

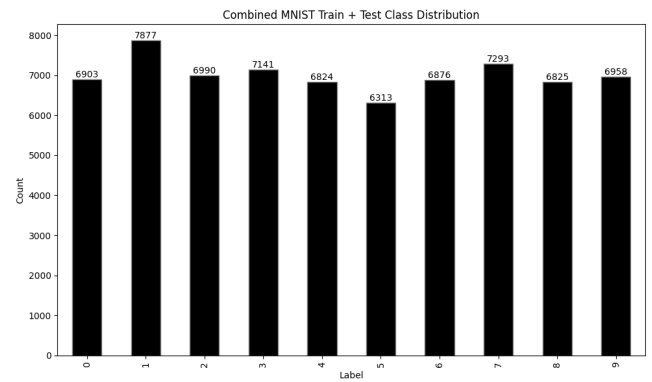


Figure 5: Digit distribution for MNIST dataset

## 1.1 Training Setup and Evaluation Approach

We apply an 80/20 train/test split and use logistic regression with tuned parameters to train a handwritten digit classifier. The setup is designed to ensure efficient convergence and balanced performance.

The model's performance will be evaluated using key metrics:

- Overall accuracy
- Precision, recall, and F1 scores (via classification report)
- Error patterns (via confusion matrix)

The implementation relies on established machine learning libraries and standard evaluation workflows to ensure reliable and reproducible results.

## 1.2 What is Logistic Regression?

Logistic regression is a simple but powerful model for classification. In stats, it is often used to tell if something belongs to a set or not — like yes (1) or no (0).

In our project, we are not doing just yes or no, we have 10 classes (digits 0 to 9). Because we use the 'lbfgs' solver, scikit-learn handles multiclass directly using softmax [3]. This means the model looks at all the classes together and picks the one with the highest score.

After we train the model on the train data, we use the left over data to check how good it works.

## 2 Parameter Tuning & Testing

For parameter tuning, you can explore various options available in the [LogisticRegression documentation](#) [3].

In this project, I iteratively test multiple values for each key hyperparameter separately, aiming to find the value that yields the best model accuracy. This approach is a form of manual hyperparameter tuning where each parameter is varied individually rather than exhaustively searching all combinations.

The parameters tuned include:

- **max\_iter**: Maximum iterations allowed for the solver to converge.
- **solver**: The optimization algorithm used (e.g., `lbfgs`, `saga`, `newton-cg`, `liblinear`).
- **class\_weight**: Strategy to handle imbalanced classes (e.g., `None` or `balanced`).
- **C**: Inverse of regularization strength; smaller values specify stronger regularization.
- **normalization value**: Pixel value scaling factor used to normalize input data (e.g., dividing by 16).

This method is simpler than full grid search or random search, focusing on tuning each parameter individually to identify its impact on model performance. For more details on hyperparameter tuning techniques including grid and random search, see [5] and the scikit-learn documentation [6, 7].

### Normalization Value vs Accuracy

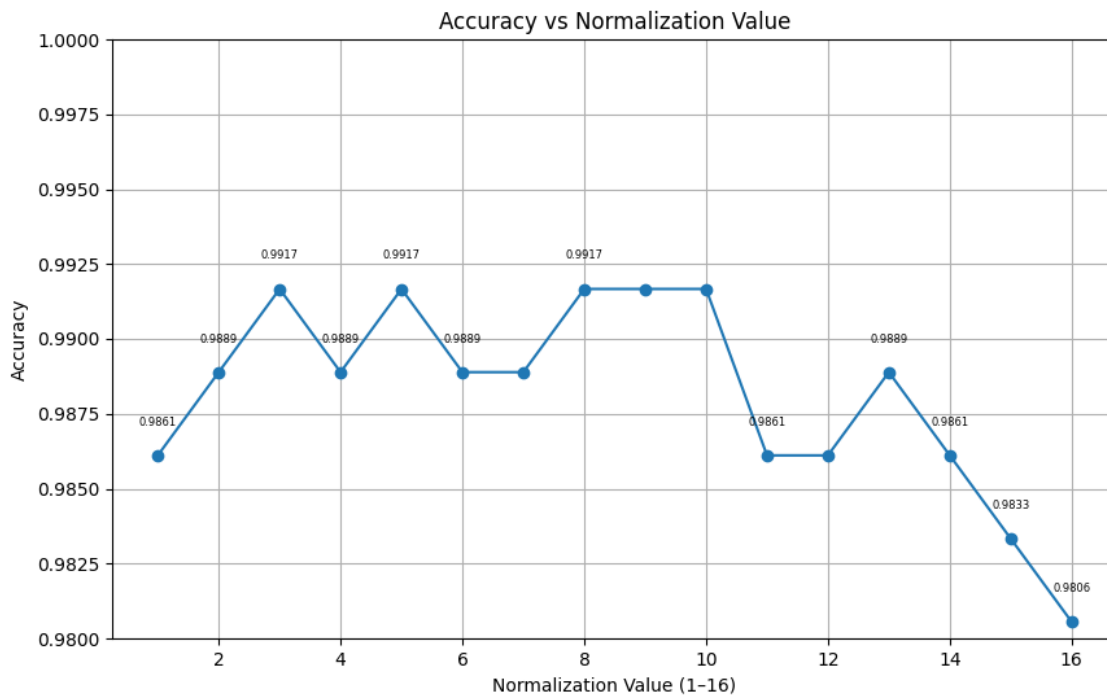


Figure 6: Accuracy stays high (98.0%–99.2%) across scaling factors 1 to 16. The best accuracy (99.2%) occurs at 3,5,8,9,10. Note: Only dividing by 16 produces true normalization (0–1 range); other divisors apply partial scaling, not full normalization.

## Test Size vs Accuracy

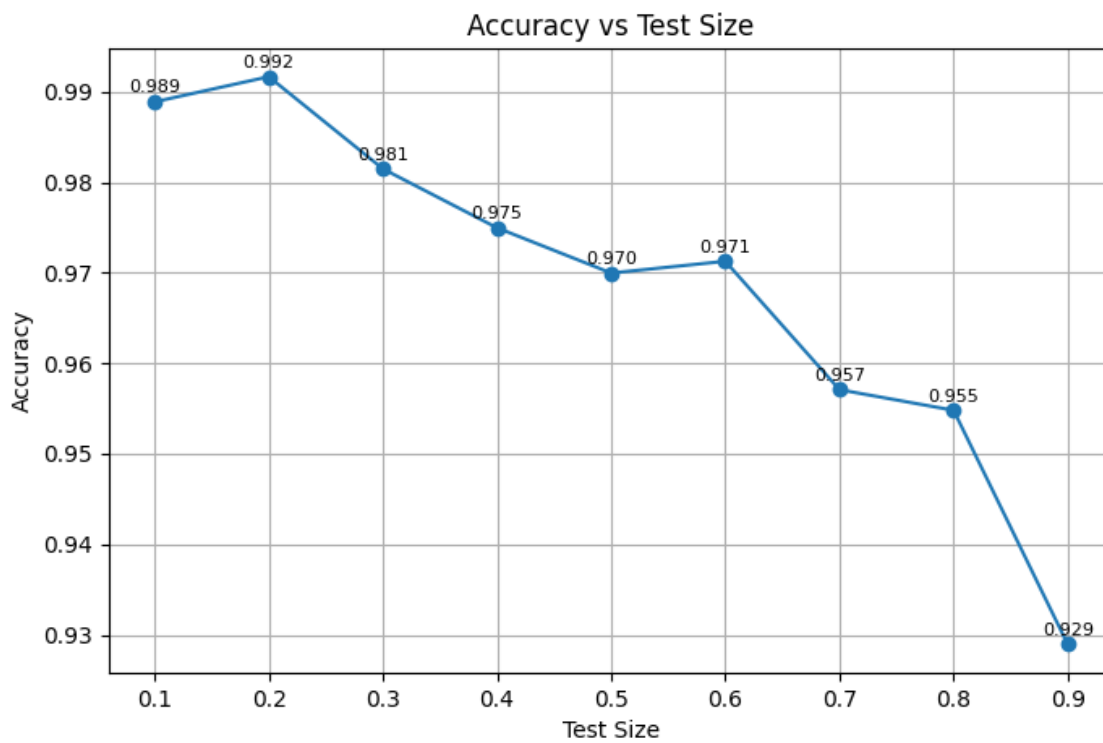


Figure 7: Accuracy vs Test Size; accuracy drops sharply below 80% training share

## Accuracy vs Random Seed

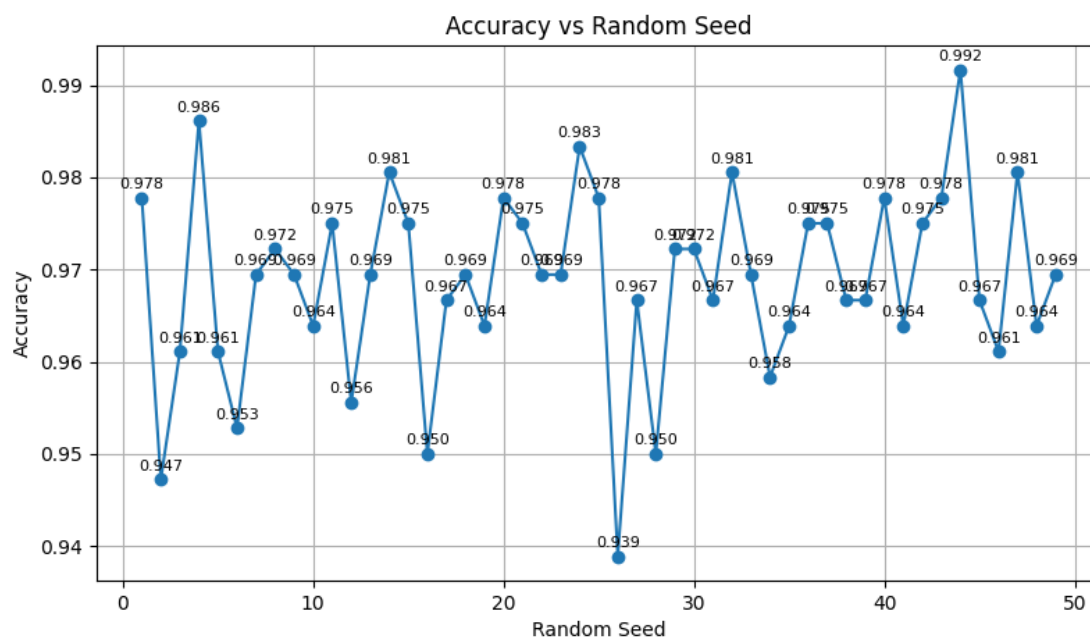


Figure 8: Accuracy vs Random Seed; accuracy remains stable regardless of seed choice max at 44

## Accuracy vs Max Iteration

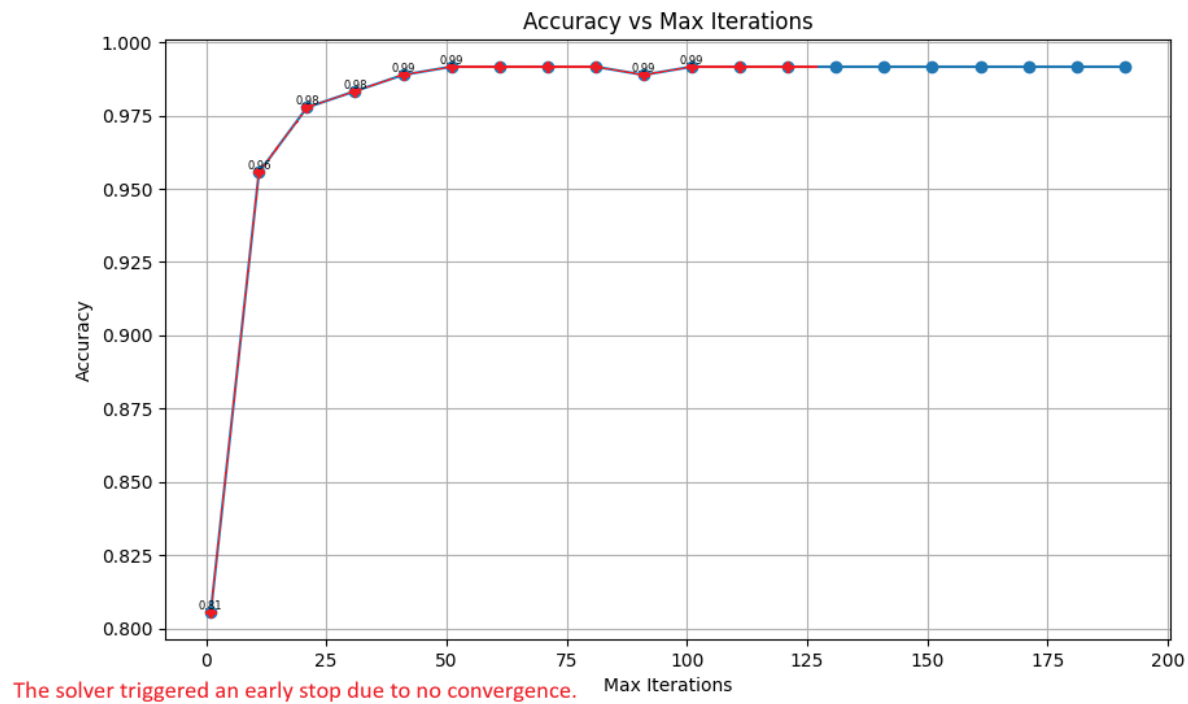


Figure 9: Accuracy vs. Max Iterations: The model’s accuracy levels off after about 126 iterations. Setting `max_iter` to 128 gives good results without wasting time. Raising `max_iter` higher does not add much benefit and only makes training slower. “However, if the number is set too high, the training process might become unnecessarily long without significant gains in performance.” [4]



## Accuracy Comparison Across Solvers

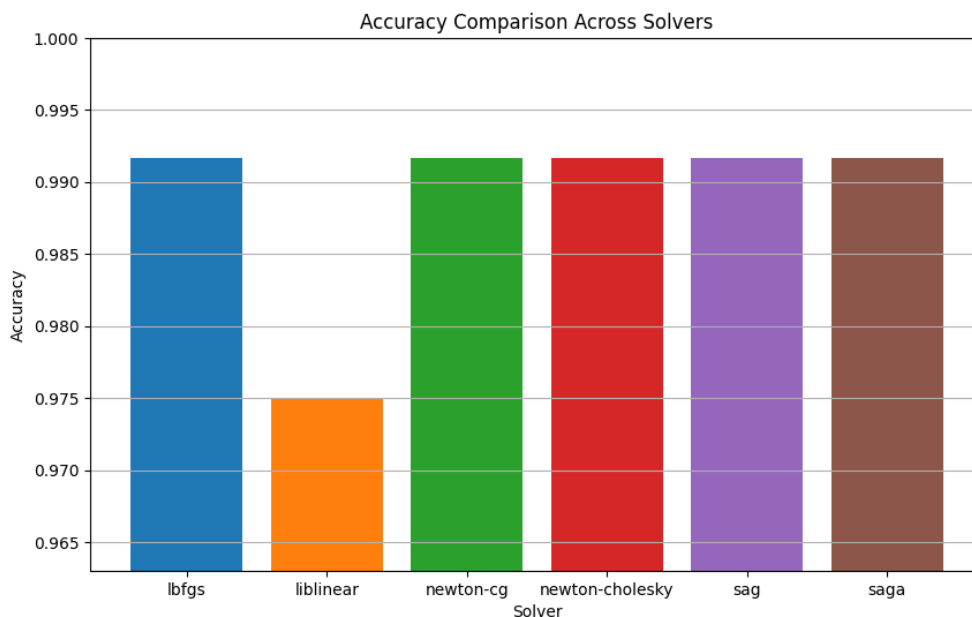


Figure 10: Accuracy comparison across solvers. Multiclass solvers perform similarly; only `liblinear` differs. We use the default `lbfgs`.

## Accuracy vs Regularization

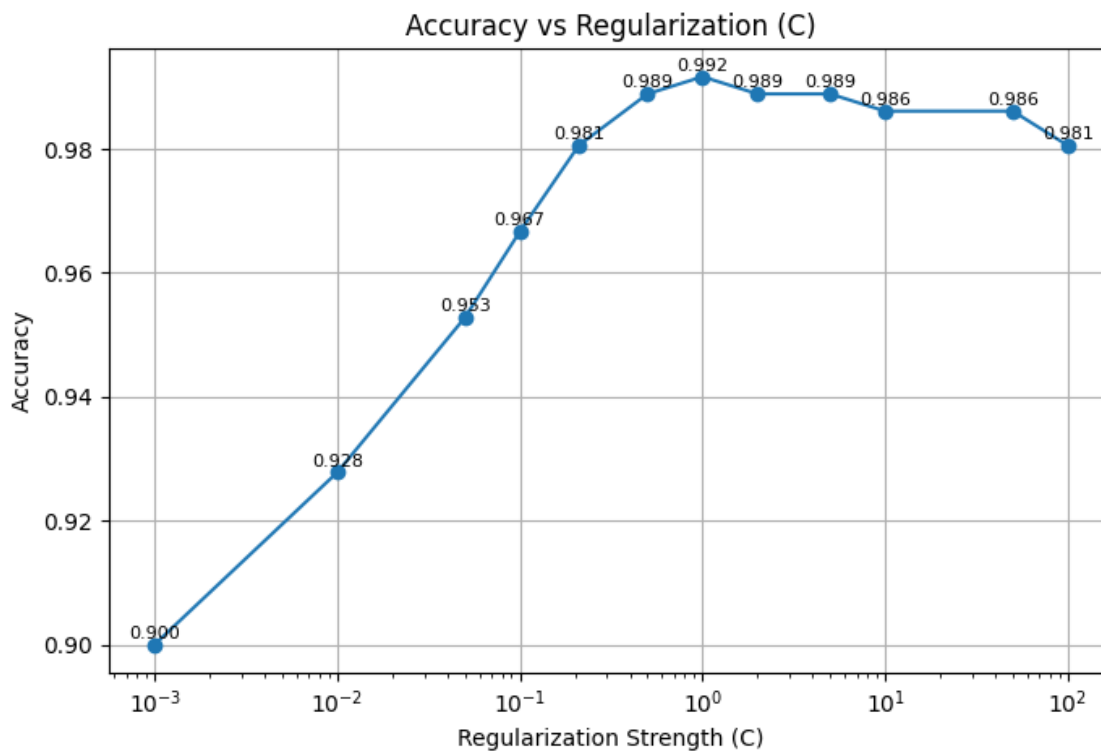


Figure 11: Accuracy vs Regularization; best performance at medium  $C$  values with  $C=1$

This means  $C=1$  works best because it avoids being too strict or too loose, keeping the model balanced.

## Accuracy vs Class Weight

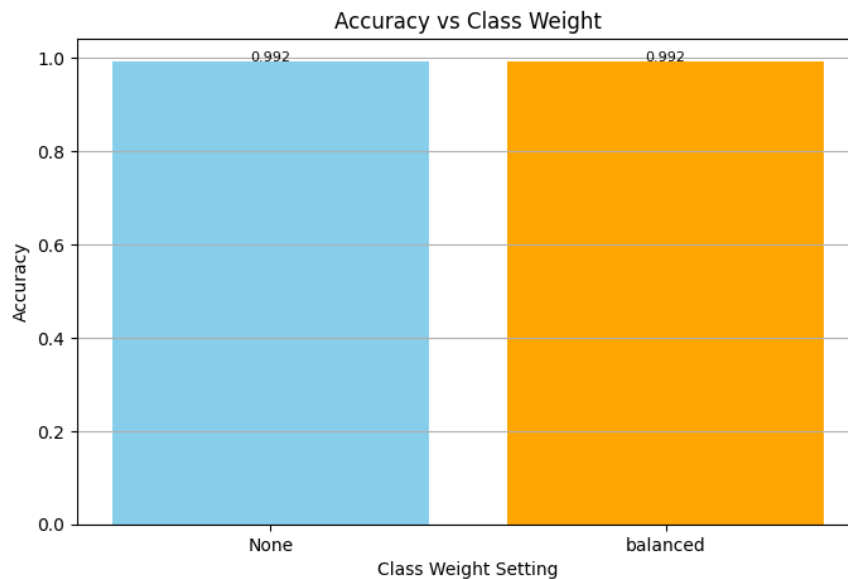


Figure 12: Accuracy vs Class Weight: This plot compares model accuracy using `class_weight=None` and `class_weight='balanced'`. Since the dataset is highly balanced (each digit averages 176 samples, range 174–183), applying class weights provides no meaningful performance gain. The classifier already handles slight class size differences, resulting in nearly identical accuracy for both settings.

## Final Recommendations

We set `max_iter = 128` since accuracy levels off beyond this, avoiding long training [3, 4]. We scaled pixels by 9.0 to keep feature values manageable. We used `lbfgs` as it works well for multiclass tasks [3].

- Solver: `lbfgs`
- Regularization: `C = 1.0`
- Pixel scaling: divide by 9.0
- Iterations: `max_iter = 128`
- Random seed: 44
- Split: 80/20 (`test_size = 0.2`)
- Class weight: None

## Best Parameter Configuration

Listing 1: Best Parameters Applied

```
1 # Train logistic regression
2 model = LogisticRegression(
3     max_iter=128,          # Ensure enough iterations for convergence
4     solver='lbfgs',        # Best multiclass solver
5     class_weight=None,     # Balanced dataset, no weighting needed
6     C=1.0                  # Regularization strength
7 ) # Set up logistic regression
8 model.fit(X_train, y_train) # Train the model
```

Achieved accuracy:

Accuracy: 0.9916666666666667

**Note:** Increasing `max_iter` to 128 improved accuracy compared to lower values, as 126 appears to be the minimum needed for convergence without triggering early stopping.

## 3 Results

### 3.1 Model and Evaluation

We used an 80/20 train/test split and logistic regression with `max_iter=128` for convergence.

*Note:* Setting `random_state=44` ensures consistent train/test splits for reproducible results. Initially, `max_iter=10000` guaranteed convergence, but later experiments show that `max_iter=128` suffices, improving efficiency without sacrificing accuracy.

Model performance was evaluated using accuracy, the classification report (precision, recall, F1), and the confusion matrix.

Imports required to run the project code:

```
1 from sklearn.datasets import load_digits
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
5 import numpy as np
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8 import pandas as pd
```

Set up the parameters:

```
1 # parameter setup for logistic regression
2 iterations_max = 128 # maximum amount of iterations
3 number_to_divide_pixels_for_normal = 9.0 # divisor for partial normalisation
4 size_of_test_data = 0.2 # 20% testing, 80% training
5 random_seed_number = 44 # random seed for reproducibility
6 solver = 'lbfgs' # solver type
7 regularization_c = 1.0 # regularization strength
```

Load and split the data:

```
1 digits = load_digits() # load built-in digit data from sklearn
2 X = digits.data / number_to_divide_pixels_for_normal # image data
3 y = digits.target # label data
4 all_the_indexes = np.arange(len(X))
5
6 # split into training and testing sets (80% train, 20% test)
7 X_train, X_test, y_train, y_test, idx_train, idx_test = train_test_split(
8     X, y, all_the_indexes, test_size=size_of_test_data, random_state=random_seed_number
9 )
```

Set up and train the logistic regression model:

```
1 model = LogisticRegression(
2     max_iter=max(iterations_max, 128), # ensure enough iterations
3     solver=solver, # best multiclass solver
4     class_weight=None, # balanced dataset, no weighting needed
5     C=regularization_c # regularization strength
6 )
7 model.fit(X_train, y_train) # train the model
```

Evaluate and print accuracy:

```
1 y_pred = model.predict(X_test) # make predictions
2 acc = accuracy_score(y_test, y_pred) # compute accuracy
3 print("Accuracy:", acc)
```

Output:

```
Accuracy: 0.9916666666666667
```

full code : [Section 6.1](#)

## 3.2 Classification Report

We refer to the `plot_classification_report` function (Listing 3) inside `plot_utils.py`, which visualizes the classification report as a heatmap.

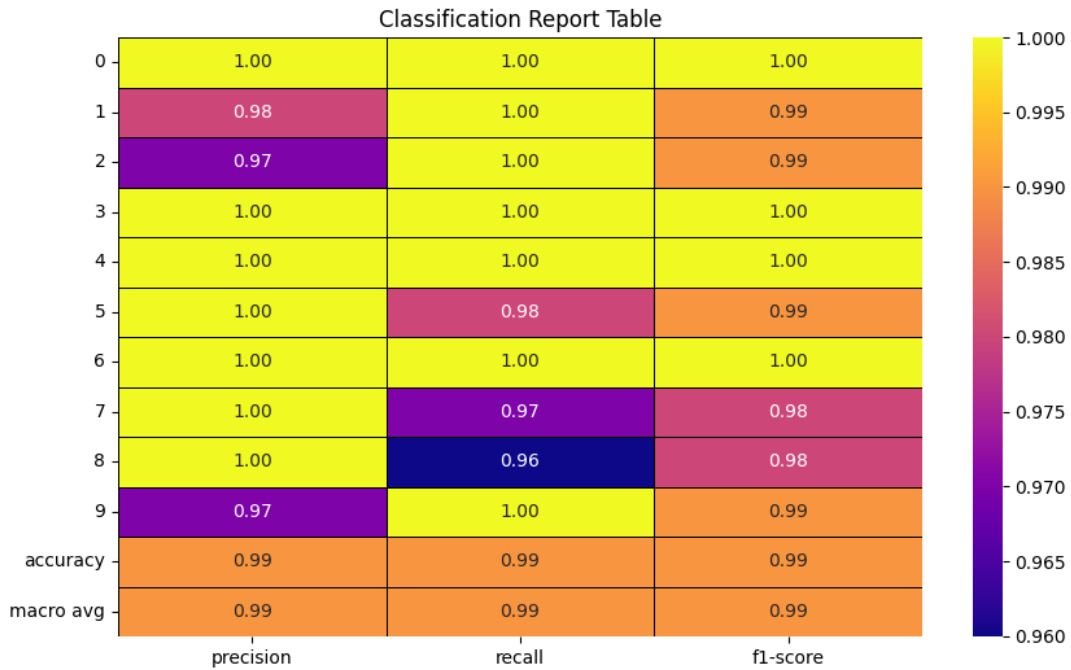


Figure 13: Classification Report Table; overall ~99% accuracy, minor challenges on digits 5 and 9

The model shows very strong performance with precision, recall, and F1 scores between 0.99 and 1.00. The macro and weighted averages are approximately 0.99, indicating consistent, near-perfect accuracy across all digits.

### 3.3 Confusion Matrix

We refer to the `plot_confusion_matrix` function (Listing 4) inside `plot_utils.py`, which visualizes the confusion matrix as a heatmap.

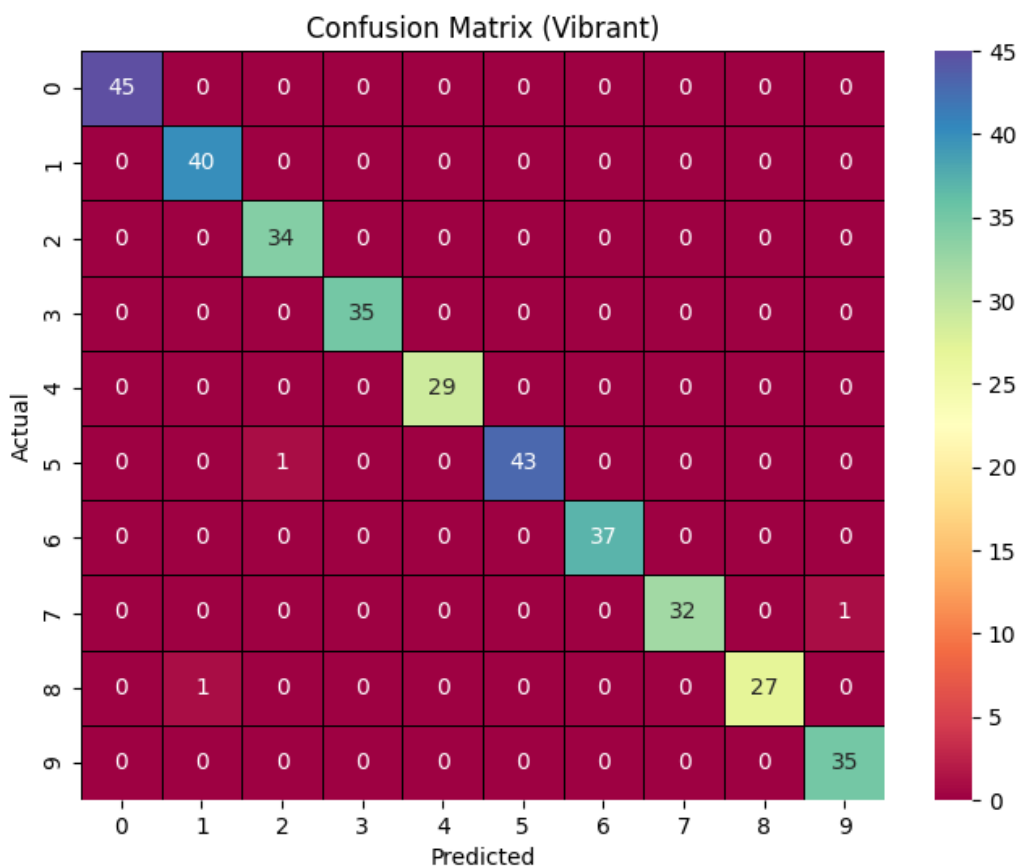


Figure 14: Confusion Matrix showing near-perfect predictions with only minor errors on digits 5, 7, and 8.

The confusion matrix shows almost all predictions are correct, with just a few small errors mainly on digits 5, 7, and 8 — reflecting the overall high accuracy of 99.2

## 4 Sample Prediction

Below is an example of a test image with its predicted and actual label (Listing 2) inside `plot_utils.py`, including the accuracy status:

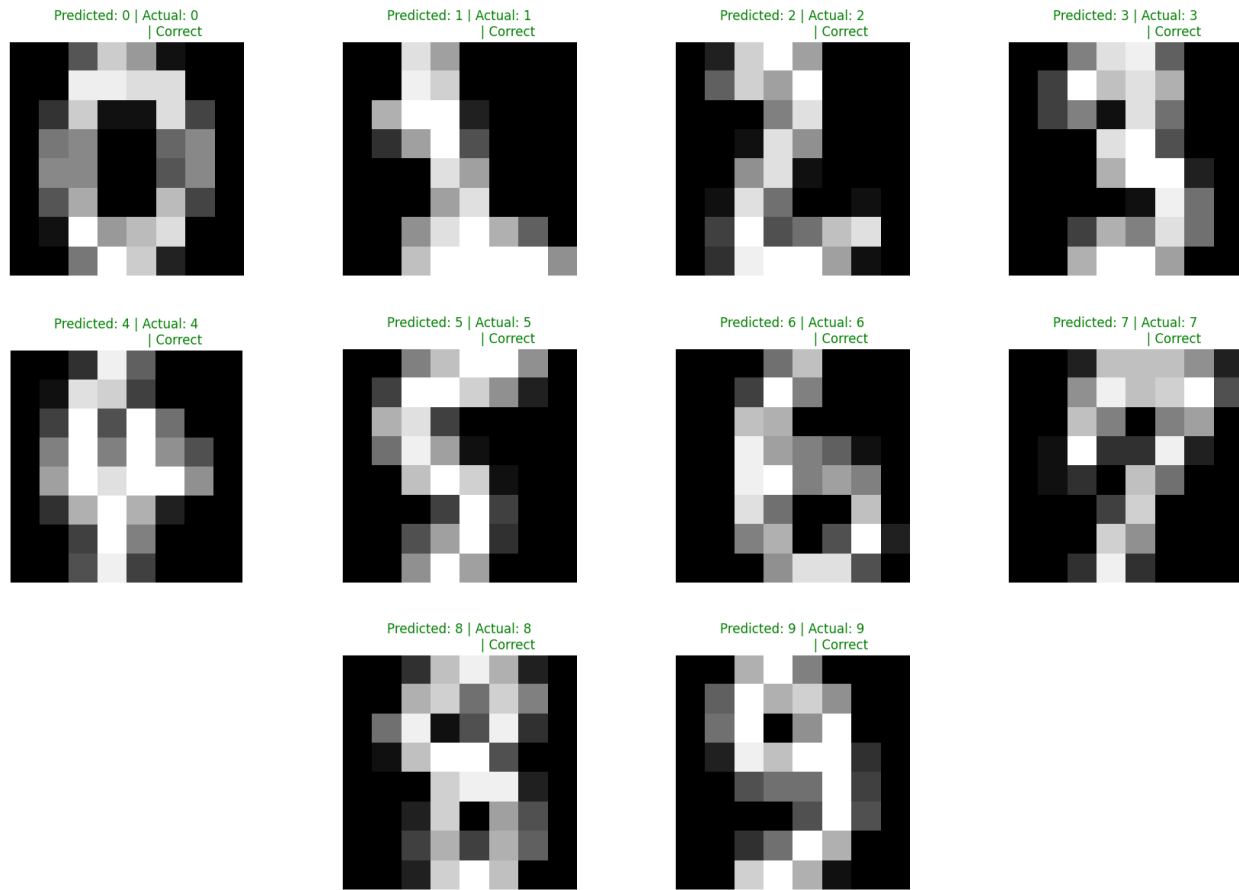


Figure 15: Sample predictions and key visual results

## 5 MNIST Logistic Regression: Small vs Large Dataset

[MNIST Dataset on Kaggle \[2\]](#)

### Parameters

Same as Section 2, except for:

- **max\_iter = 1000**: Maximum iterations allowed for solver convergence (since its a bigger dataset).
- **Normalization divisor = 255.0**: Pixel values scaled to the  $[0, 1]$  range.

The code for this is in Section 6.2.

### Overview

We compare logistic regression performance on:

- Small dataset (scikit-learn dataset): 1,438 train / 359 test (80/20 split)
- Full MNIST: 60,000 train / 10,000 test (80/20 split)

This isolates dataset size effects on accuracy and generalization.

### Accuracy Comparison

- **Small dataset**: 99.2%
- **Full dataset**: 92.2%

The 4.5% drop shows larger data size challenges model generalization despite proportional splits.

### Confusion Matrices

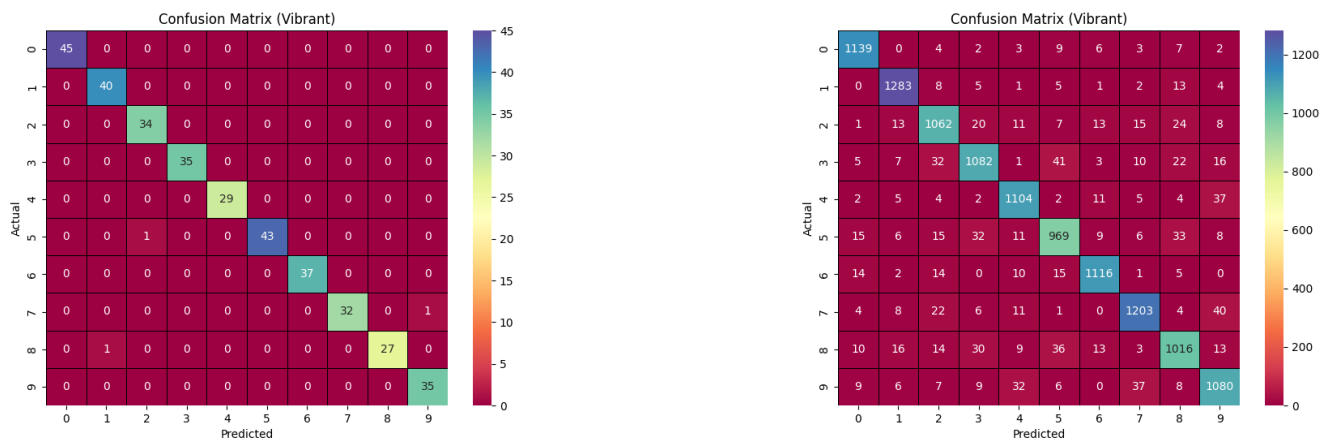


Figure 16: Small dataset: near-perfect classification; Full dataset: more misclassifications due to complexity.

Small data shows near-perfect results; full data reveals errors mainly among digits 2, 3, 5, 8, and 9.

### Classification Reports

#### Test Set Size Impact

- Both use 80/20 splits.
- Small test set: 360 samples.
- Full test set: 10,000 samples.

Larger test sets provide more reliable, robust evaluation, uncovering weaknesses hidden by small samples.

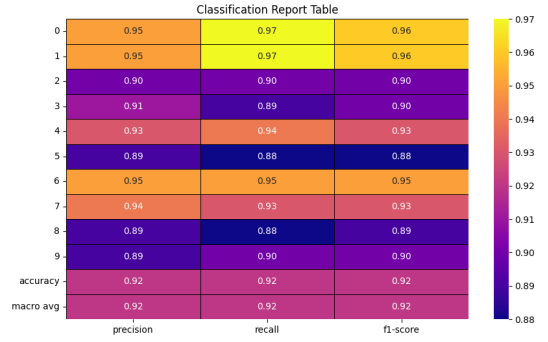
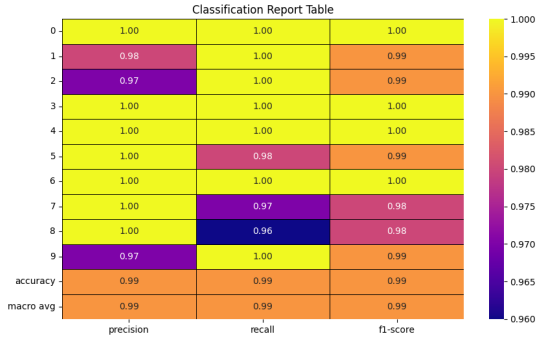


Figure 17: Small dataset: near-perfect F1 scores; Full dataset: average F1 0.92 reflecting greater diversity.

## Scaling Effects and Recommendations

The full dataset is  $33\times$  larger in train and test size. Scaling improves generalization but requires:

- Stronger regularization.
- Careful hyperparameter tuning.
- Possibly advanced models (e.g., neural networks) for top accuracy.

## Conclusion

Small datasets can mask true model weaknesses. Large datasets reveal realistic performance, ensuring better real-world readiness despite minor accuracy drops.

## Final Recommendations

- Set `max_iter = 128` to ensure convergence without unnecessary computation.
- Use `solver = 'lbfgs'` for efficient multiclass logistic regression.
- Keep regularization parameter `C = 1.0` for balanced fitting.
- Normalize pixel values by dividing by 9 for best observed accuracy.
- Use `random_state = 44` for consistent train/test splits.
- Maintain an 80/20 train/test split.

This setup achieves a top accuracy of **99.2%** on the `sklearn` dataset and **92.2%** on the MNIST dataset, reflecting excellent model performance.



## 6 Project Code

This section will contain all the relevant code used threw out the project

### 6.1 Sklearn Dataset Code

File: sklearn\_digits\_logreg\_tuned.py code for the sklearn data set

```
1 # tuned_mnist_logistic_regression_parameter.py
2
3 from sklearn.datasets import load_digits
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import LogisticRegression
6 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
7 import numpy as np
8
9 # very simple version
10 def run_the_logistic_regression_model():
11     # parameter setup for logistic regression
12     iterations_max = 128
13     number_to_divide_pixels_for_normal = 9.0
14     size_of_test_data = 0.2 # 20% testing, 80% training
15     random_seed_number = 44 # random seed
16     solver = 'lbfgs' # solver type
17     regularization_c = 1.0 # regularization strength
18
19     digits = load_digits() # load built-in digit data from sklearn
20     X = digits.data / number_to_divide_pixels_for_normal # image data
21     y = digits.target # label data
22     all_the_indexes = np.arange(len(X))
23
24     # split into training and testing sets
25     X_train, X_test, y_train, y_test, idx_train, idx_test = train_test_split(
26         X, y, all_the_indexes, test_size=size_of_test_data, random_state=random_seed_number
27     )
28
29     model = LogisticRegression(
30         max_iter=max(iterations_max, 128), # ensure enough iterations
31         solver=solver,
32         class_weight=None,
33         C=regularization_c
34     )
35     model.fit(X_train, y_train) # train the model
36
37     y_pred = model.predict(X_test) # make predictions
38     acc = accuracy_score(y_test, y_pred)
39     print("Accuracy:", acc)
40     print(classification_report(y_test, y_pred))
41     print(confusion_matrix(y_test, y_pred))
42
43     return acc
44
45 if __name__ == "__main__":
46     final_accuracy = run_the_logistic_regression_model()
47     print(f"Returned accuracy: {final_accuracy}")
```

## 6.2 MNIST Dataset Code

code for the mnsit dataset File: mnist\_logistic\_regression.py

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import plot_utils
5 import seaborn as sns
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
8 from sklearn.model_selection import train_test_split
9
10 # Load both MNIST train and test CSVs
11 train_data = pd.read_csv("mnist/mnist_train.csv")
12 test_data = pd.read_csv("mnist/mnist_test.csv")
13
14 # Combine datasets
15 full_data = pd.concat([train_data, test_data], ignore_index=True)
16
17
18 # Separate features and labels, normalize
19 X = full_data.drop('label', axis=1) / 255.0
20 y = full_data['label']
21
22 # Train-test split 80/20
23 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
24
25 # Train logistic regression
26 model = LogisticRegression(
27     max_iter=1000,
28     solver='lbfgs',
29     class_weight=None,
30     C=1.0
31 )
32 model.fit(X_train, y_train)
33
34 # Predict and evaluate
35 y_pred = model.predict(X_test)
36 acc = accuracy_score(y_test, y_pred)
37 print("Accuracy:", acc)
38 print(classification_report(y_test, y_pred))
39 print(confusion_matrix(y_test, y_pred))
40
41 # Plot confusion matrix and classification report
42 plot_utils.plot_confusion_matrix(y_test, y_pred, labels=list(range(10)))
43 plot_utils.plot_classification_report(y_test, y_pred)
44
45 # Display one test image with prediction
46 index = 1
47 y_test_array = np.array(y_test)
48 is_correct = y_pred[index] == y_test_array[index]
49
50 image_data = X_test.iloc[index].values.reshape(28, 28)
51
52 plot_utils.display_prediction_image(
53     image=image_data,
54     predicted=y_pred[index],
55     actual=y_test_array[index],
56     accuracy=acc,
57     is_correct=is_correct,
58     y_test=y_test_array,
59     y_pred=y_pred,
60     reshape_size=28
61 )
```

## 6.3 Utility Functions

### File: plot\_utils.py

This code allows the main script to display graphical interpretations, including the confusion matrix, the classification report, or a single predicted image.

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 import numpy as np
4 import pandas as pd
5 from sklearn.metrics import confusion_matrix, classification_report
```

### 6.3.1 display\_prediction\_image Function

Listing 2: display\_prediction\_image function

```
1 # display a single prediction image with details
2 def display_prediction_image(image, predicted, actual, accuracy, is_correct, y_test, y_pred,
   reshape_size):
3     indices = np.where(y_test == actual)[0]
4     total = len(indices)
5     correct = np.sum(y_pred[indices] == y_test[indices])
6     digit_acc = correct / total if total > 0 else 0
7
8     plt.figure(figsize=(4, 4))
9     plt.imshow(image.reshape(reshape_size, reshape_size), cmap='gray')
10    result_text = "Correct" if is_correct else "Wrong"
11    color = 'green' if is_correct else 'red'
12    plt.title(f"\nPredicted: {predicted} | Actual: {actual}\n"
13             f"Overall Acc: {accuracy:.2f} | Digit Acc: {digit_acc:.2f} | {result_text}",
14             fontsize=10, color=color)
15    plt.axis('off')
16    plt.tight_layout()
17    plt.show()
```

### 6.3.2 plot\_classification\_report Function

Listing 3: plot\_classification\_report function

```
1 def plot_classification_report(y_true, y_pred):
2     report_dict = classification_report(y_true, y_pred, output_dict=True)
3     df = pd.DataFrame(report_dict).transpose().round(2)
4     plt.figure(figsize=(10, 6))
5     sns.heatmap(df.iloc[:-1, :-1], annot=True, fmt='.2f', cmap='plasma',
6                cbar=True, linewidths=0.5, linecolor='black')
7     plt.title('Classification Report Table')
8     plt.yticks(rotation=0)
9     plt.show()
```

### 6.3.3 plot\_confusion\_matrix Function

Listing 4: plot\_confusion\_matrix function

```
1 # plot confusion matrix as heatmap
2 def plot_confusion_matrix(y_true, y_pred, labels):
3     cm = confusion_matrix(y_true, y_pred)
4     plt.figure(figsize=(8, 6))
5     sns.heatmap(cm, annot=True, fmt='d', cmap='Spectral',
6                cbar=True, linewidths=0.5, linecolor='black',
7                xticklabels=labels, yticklabels=labels)
8     plt.xlabel('Predicted')
9     plt.ylabel('Actual')
10    plt.title('Confusion Matrix (Vibrant)')
11    plt.show()
```

## References

- [1] Scikit-learn documentation,  
[https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load\\_digits.html](https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_digits.html)
- [2] MNIST dataset on Kaggle,  
<https://www.kaggle.com/datasets/oddrationalle/mnist-in-csv>
- [3] Scikit-learn LogisticRegression documentation,  
[https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LogisticRegression.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html)
- [4] GeeksforGeeks, *How to Optimize Logistic Regression Performance*,  
[https://www.geeksforgeeks.org/how-to-optimize-logistic-regression-performance/?utm\\_source=chatgpt.com](https://www.geeksforgeeks.org/how-to-optimize-logistic-regression-performance/?utm_source=chatgpt.com)
- [5] Jason Brownlee, *Hyperparameter Optimization With Random Search and Grid Search*, MachineLearningMastery.com, 2020.  
<https://machinelearningmastery.com/hyperparameter-optimization-with-random-search-and-grid-search/>
- [6] Scikit-learn GridSearchCV documentation,  
[https://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.GridSearchCV.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html)
- [7] Scikit-learn RandomizedSearchCV documentation,  
[https://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.RandomizedSearchCV.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.RandomizedSearchCV.html)