

# Classification on Colab using MNIST dataset

Submitted by Manish Bafna

Student Id: 19655

Instructor: Dr. Henry Chang

GIT: [Classification using MNIST dataset](#)

# Table of Content

Introduction

Design

Implementation

Test

Enhancement Ideas

Conclusion

References

# Introduction

The most common supervised learning tasks are

- Regression (predicting values)
  - Linear Regression
  - Non-linear Regression
    - Decision Trees
      - Random Forests
- Classification (predicting classes)

# Introduction

MNIST dataset is a set of 70,000 small images of handwritten digits. Scikit-Learn is a Machine Learning Library which provides many helper functions to download popular datasets, MNIST is one of them.

Instead of recognize 9 digits

- Let's simplify the problem for now  
and only try to identify one digit
  - For example, the number 5.



# Introduction

## Classification Model Evaluation Metric: Confusion Matrix



Low Accuracy  
High Precision



High Accuracy  
Low Precision



High Accuracy  
High Precision

Suppose you're working as a **security guard** at a **concert** and your **job** is to

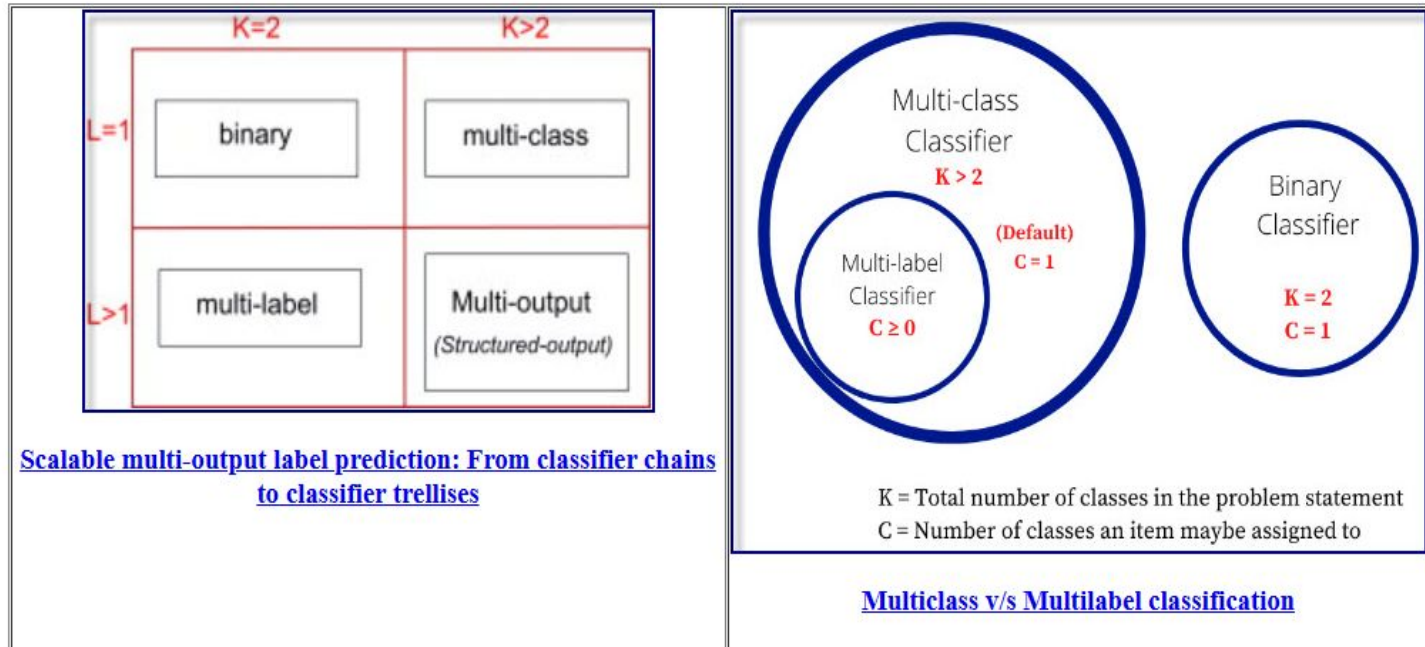
- identify any **individuals** who were **troublemakers** at **previous events** (represented as "**positive**" class in this **problem**).

Recall	F1 Score
<p>Recall measures the <b>proportion</b> of <b>true positive</b> instances (<b>troublemakers</b>) that were <b>correctly identified</b> by the <b>classifier</b>.</p> <ul style="list-style-type: none"><li>• The <b>fraction</b> of <b>actual troublemakers</b> that were <b>correctly identified</b> by you as <b>troublemakers</b> and allowed into the <b>concert</b>.</li></ul>	<p>F1 Score is a balance between</p> <ul style="list-style-type: none"><li>• <b>Precision</b> and<ul style="list-style-type: none"><li>◦ <b>Precision</b> is the <b>fraction</b> of <b>instances</b> classified as <b>positive</b> that are actually <b>positive</b>,</li></ul></li><li>• <b>Recall</b><ul style="list-style-type: none"><li>◦ <b>Recall</b> is the <b>fraction</b> of <b>positive instances</b> that were correctly classified as <b>positive</b>.</li></ul></li></ul>
<p>A <b>high recall score</b> means that you are able to catch most of the <b>troublemakers</b>, but you may also let in some <b>non-troublemakers</b> (<b>false positives</b>).</p>	<p>In this analogy, <b>F1 Score</b> would represent the <b>balance</b> between catching all the <b>troublemakers</b> (<b>high recall</b>) and only letting in the <b>non-troublemakers</b> (<b>high precision</b>).</p> <ul style="list-style-type: none"><li>◦ A <b>high F1 Score</b> means that you have found a <b>good balance</b> between these <b>two metrics</b>.</li></ul>

# Introduction

## Types of Classification

### Summary



# Implementation

- This “5-detector” will be an example of a binary classifier, capable of distinguishing between just two classes, 5 and not-5.
- Let’s create the target vectors for this classification task:  
`y_train_5 = (y_train == 5) # True for all 5s, False for all other digits`
- `y_test_5 = (y_test == 5)`
- Now let’s pick a classifier and train it.
  - A good place to start is with a Stochastic Gradient Descent (SGD) classifier, using Scikit-Learn’s `SGDClassifier` class.
  - This classifier has the advantage of being capable of handling very large datasets efficiently.
  - This is in part because SGD deals with training instances independently, one at a time (which also makes SGD well suited for online learning), as we will see later.

# Implementation

Let's create an SGDClassifier and train it on the whole training set:

```
from sklearn.linear_model import SGDClassifier
```

```
sgd_clf = SGDClassifier(random_state=42)
```

```
sgd_clf.fit(X_train, y_train_5)
```

- The SGDClassifier relies on randomness during training (hence the name “stochastic”).
  - If you want reproducible results, you should set the `random_state` parameter.
- Now we can use it to detect images of the number 5:

```
>>> sgd_clf.predict([some_digit])
```

```
array([ True])
```

- The classifier guesses that this image represents a 5 (True).
  - Looks like it guessed right in this particular case!



# Implementation

```
▶ %matplotlib inline
import matplotlib as mpl
import matplotlib.pyplot as plt

some_digit = X[0]
some_digit_image = some_digit.reshape(28, 28)
plt.imshow(some_digit_image, cmap=mpl.cm.binary)
plt.axis("off")

save_fig("some_digit_plot")
plt.show()
```

👤 Saving figure some\_digit\_plot



# Implementation

```
plt.figure(figsize=(9,9))  
example_images = X[:100]  
plot_digits(example_images, images_per_row=10)  
save_fig("more_digits_plot")  
plt.show()
```

Saving figure more\_digits\_plot



# Implementation

Use Google Colab

Understanding MNIST database

Run the [03\\_classification.ipynb](#) on google colab

# Test

## Measuring Accuracy Using Cross-Validation

```
[18] from sklearn.model_selection import StratifiedKFold
      from sklearn.base import clone

      skfolds = StratifiedKFold(n_splits=3, shuffle=True, random_state=42)

      for train_index, test_index in skfolds.split(X_train, y_train_5):
          clone_clf = clone(sgd_clf)
          X_train_folds = X_train[train_index]
          y_train_folds = y_train_5[train_index]
          X_test_fold = X_train[test_index]
          y_test_fold = y_train_5[test_index]

          clone_clf.fit(X_train_folds, y_train_folds)
          y_pred = clone_clf.predict(X_test_fold)
          n_correct = sum(y_pred == y_test_fold)
          print(n_correct / len(y_pred))
```

```
0.9669
0.91625
0.96785
```

# Test

## Confusion Matrix

```
✓ [23] y_train_perfect_predictions = y_train_5 # pretend we reached perfection  
In confusion_matrix(y_train_5, y_train_perfect_predictions)
```

```
array([[54579,    0],  
       [    0,  5421]])
```

# Test

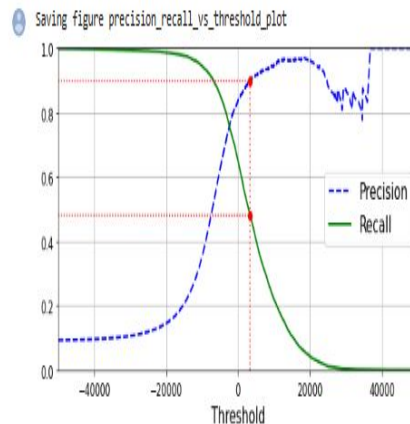
## Confusion Matrix

## Precision and Recall:

```
[36] def plot_precision_recall_vs_threshold(precisions, recalls, thresholds):  
    plt.plot(thresholds, precisions[:-1], "b--", label="Precision", linewidth=2)  
    plt.plot(thresholds, recalls[:-1], "g-", label="Recall", linewidth=2)  
    plt.legend(loc="center right", fontsize=16) # Not shown in the book  
    plt.xlabel("Threshold", fontsize=16)      # Not shown  
    plt.grid(True)                          # Not shown  
    plt.axis([-50000, 50000, 0, 1])         # Not shown
```

```
recall_90_precision = recalls[np.argmax(precisions >= 0.90)]  
threshold_90_precision = thresholds[np.argmax(precisions >= 0.90)]
```

```
plt.figure(figsize=(8, 4))  
plot_precision_recall_vs_threshold(precisions, recalls, thresholds) # Not shown  
plt.plot([threshold_90_precision, threshold_90_precision], [0., 0.9], "r:") # Not shown  
plt.plot([-50000, threshold_90_precision], [0.9, 0.9], "r:") # Not shown  
plt.plot([-50000, threshold_90_precision], [recall_90_precision, recall_90_precision], "r:") # Not shown  
plt.plot([threshold_90_precision], [0.9], "ro") # Not shown  
plt.plot([threshold_90_precision], [recall_90_precision], "ro") # Not shown  
save_fig("precision_recall_vs_threshold_plot") # Not shown  
plt.show()
```



# Test

```
✓ [86] ambiguous_digit = X_test[2589]  
    knn_clf.predict_proba([ambiguous_digit])  
  
array([[0.24579675, 0.        , 0.        , 0.        ,  
        0.        , 0.        , 0.        , 0.75420325]])
```

```
✓ [87] plot_digit(ambiguous_digit)
```



# Enhancement

**Hyperparameter Tuning:** The performance of KNN can be greatly improved by fine-tuning the hyperparameters such as the number of neighbors, the distance metric, etc. We can use techniques such as grid search or random search to find the optimal hyperparameters.

**Anomaly Detection:** We can use KNN for anomaly detection on MNIST, where the goal is to identify images that are significantly different from the rest of the data.

**Feature Engineering:** Feature engineering is an important step in enhancing the performance of KNN. We can try different features such as edge detection, shape detection, histograms, etc.



# Conclusion

In conclusion, KNN is a simple yet powerful classification algorithm that can be applied to the MNIST dataset.

By following best practices such as data preprocessing, feature engineering, dimensionality reduction, hyperparameter tuning, and ensemble methods, we can significantly improve the performance of KNN on MNIST.

However, it's important to remember that there is no one-size-fits-all solution, and the best approach may vary based on the specific use case and requirements. It's always a good idea to experiment with different approaches and evaluate their performance to find the best solution for the problem.

# References

[3 Classification.in \(sfbu.edu\)](#)

[03\\_classification.ipynb - Colaboratory \(google.com\)](#)