

---

# Project 2: Learning to Rank using Linear Regression

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

The goal of this project is to solve the handwriting comparison task in forensics. We formulate this as a problem of linear regression where we map a set of input features  $x$  to a real-valued scalar target  $y(x, w)$ .

Our task is to find similarity between the handwritten samples of the known and the questioned writer by using linear regression models and Neural Network.

Each instance in the CEDAR "AND" training data consists of set of input features for each hand-written "AND" sample. The features are obtained from two different sources:

(a) Human Observed features: Features entered by human document examiners manually.

(b) GSC features: Features extracted using Gradient Structural Concavity (GSC) algorithm.

The target values are scalars that can take two values {1:same writer, 0:different writers}. Although the training target values are discrete we use regression models to obtain real values which is more useful for finding similarity.

We have two objectives Three this project:

1. Train a Linear Regression model on both Human read feature data-set and GSC features data-set using Concatenation and Subtraction technique.

2. Train a Logistic Regression model on both Human read feature data-set and GSC features data-set using Concatenation and Subtraction technique.

3. Train a Neural Network model (Keras) on both Human read feature data-set and GSC features data-set using Concatenation and Subtraction technique.

## 1 Types of Datasets:

Based on feature extraction process, we have provided two datasets:

### 1.1 Human Observed Dataset:

The Human Observed dataset shows only the cursive samples in the data set, where for each image the features are entered by the human document examiner. There are 9 distinct features for each image in human observed dataset.

Both concatenation and subtraction has been performed based on the image id. After performing concatenation based on the image id, each sample will have 18 feature values and the dataset will look like the following image.

img_id_A	img_id_B	f <sub>A1</sub>	f <sub>A2</sub>	f <sub>A3</sub>	f <sub>A4</sub>	f <sub>A5</sub>	f <sub>A6</sub>	f <sub>A7</sub>	f <sub>A8</sub>	f <sub>A9</sub>	f <sub>B1</sub>	f <sub>B2</sub>	f <sub>B3</sub>	f <sub>B4</sub>	f <sub>B5</sub>	f <sub>B6</sub>	f <sub>B7</sub>	f <sub>B8</sub>	f <sub>B9</sub>	t
1121a_num1	1121b_num2	2	1	1	3	2	2	0	1	2	2	1	1	0	2	2	0	3	2	1
1121a_num1	1386b_num1	2	1	1	3	2	2	0	1	2	3	1	1	0	2	2	0	1	2	0

Figure 1: Neural Network

## 1.2 GSC Dataset:

Gradient Structural Concavity algorithm generates 512 features for an input handwritten "AND" image. The dataset is named as "GSC-Features-Data". Similar to the Human observed dataset.

Both concatenation and subtraction has been performed based on the image id. After performing concatenation based on the image id, each sample will have 1024 feature values and the dataset will look like the following image.

img_id_A	img_id_B	f <sub>A1</sub>	f <sub>A2</sub>	f <sub>A3</sub>	f <sub>A4</sub>	f <sub>A5</sub>	f <sub>A6</sub>	...	f <sub>A512</sub>	f <sub>B1</sub>	f <sub>B2</sub>	f <sub>B3</sub>	f <sub>B4</sub>	f <sub>B5</sub>	f <sub>B6</sub>	...	f <sub>B512</sub>	t
1121a_num1	1121b_num2	0	1	1	0	1	0	...	0	0	1	1	0	0	1	...	1	1
1121a_num1	1386b_num1	0	1	1	0	1	0	...	0	1	1	1	0	1	0	...	0	0

Figure 2: Neural Network

## 1.3 Data Cleansing:

The GSC dataset contains several such features which contains all zero values (or all same values) which in terms generate a determinant of 0. These makes the dataset "non-inversable".

For this reason, after performing concatenation/subtraction and removing the img\_id columns and target column, all the features containing all 0 values are removed from the GSC dataset.

## 2 E\_RMS Graphs:

The following diagrams describes the E\_RMS values for each iteration plotted in a line graph:

### 2.1 Linear Regression on Human observed dataset:

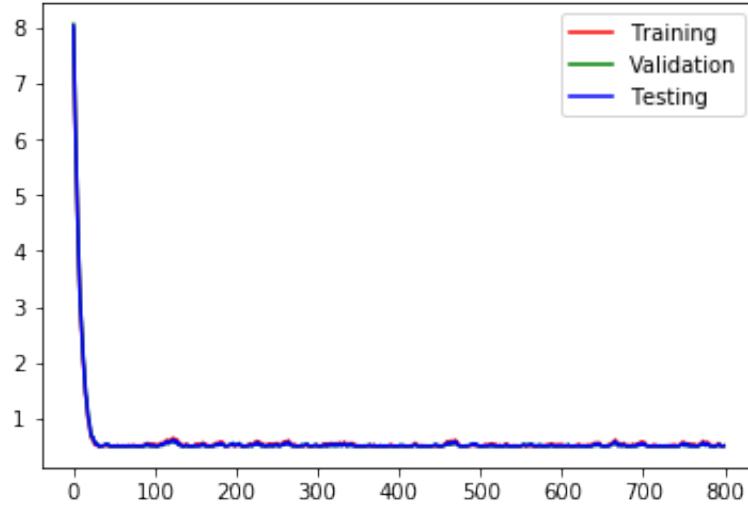


Figure 3: Concatenation Technique

E\_rms Training = 0.49908  
E\_rms Validation = 0.49628  
E\_rms Testing = 0.49789

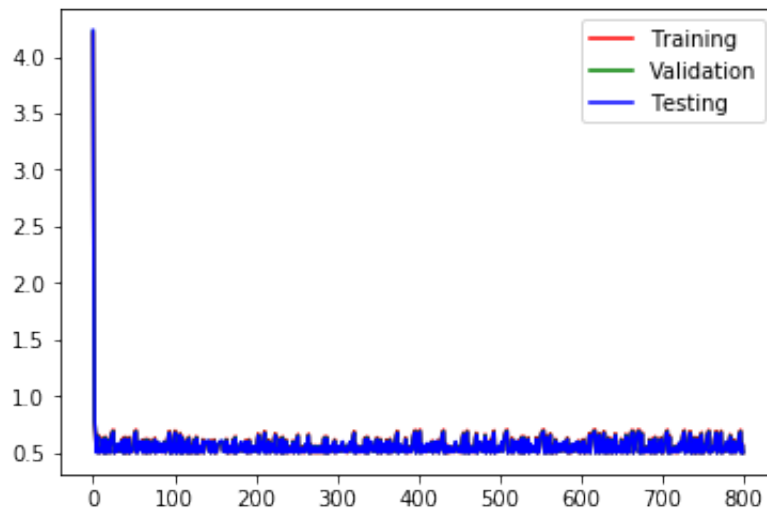


Figure 4: Subtraction Technique

E\_rms Training = 0.49993  
E\_rms Validation = 0.49704  
E\_rms Testing = 0.49898

## 2.2 Logistic Regression on GSC dataset:

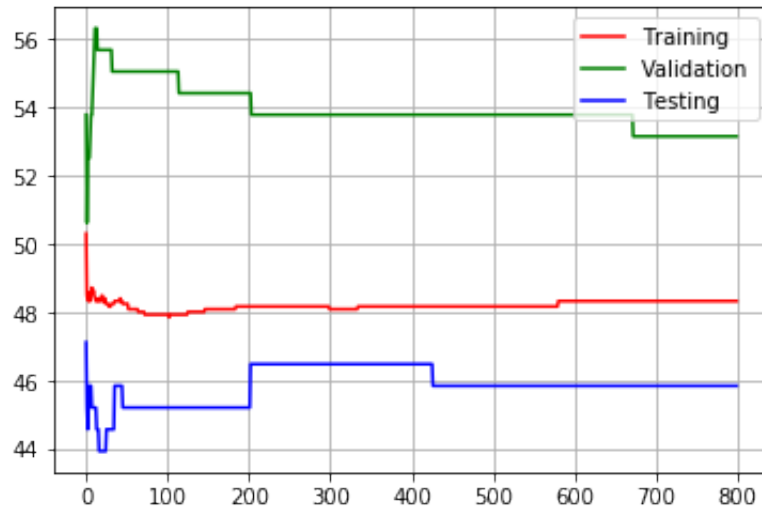


Figure 5: Concatenation Technique

Training\_accuracy: 50.00  
 Validation\_accuracy: 51.89873  
 Testing\_accuracy: 47.7707

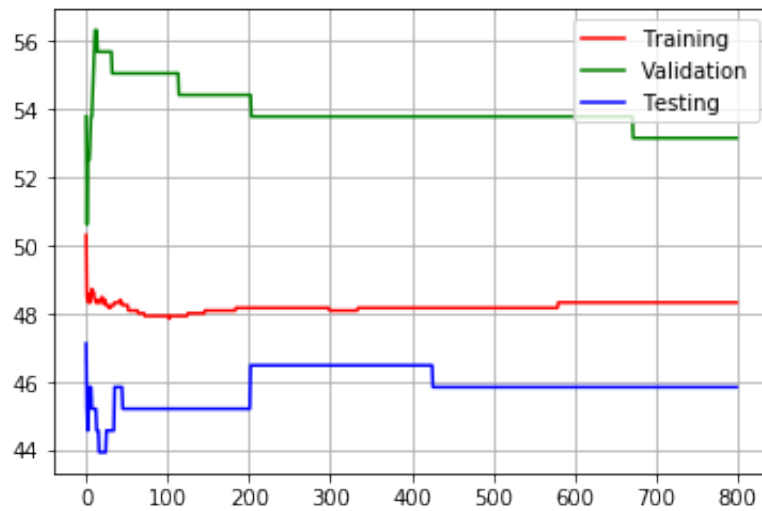


Figure 6: Subtraction Technique

Training\_accuracy: 50.31596  
 Validation\_accuracy: 56.32911  
 Testing\_accuracy: 47.13376

### 64 2.3 Neural Network Model on Human observed dataset:

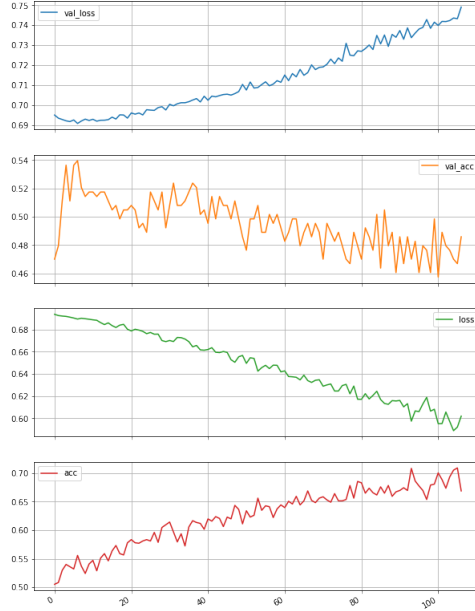


Figure 7: Concatenation Technique

65 Training\_accuracy: 0.5238  
66 Validation\_accuracy: 0.4858  
67 Testing\_accuracy: 0.4902  
68

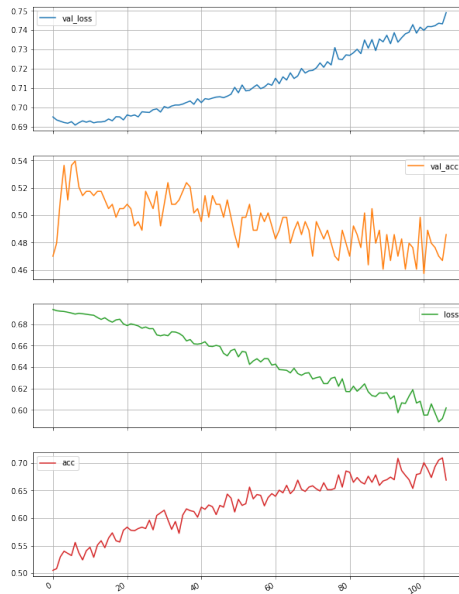


Figure 8: Subtraction Technique

69 Training\_accuracy: 0.5127  
70 Validation\_accuracy: 0.4700  
71 Testing\_accuracy: 0.4895  
72

73 **2.4 Linear Regression on GSC dataset:**

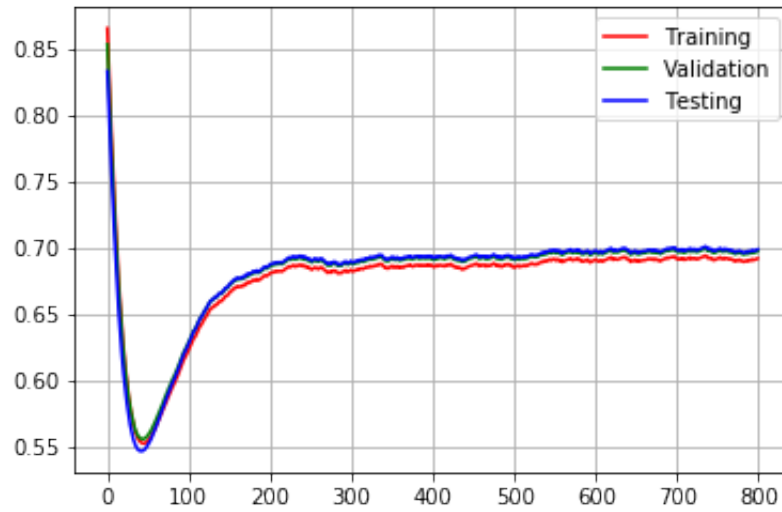


Figure 9: Concatenation Technique

74 E\_rms Training = 0.55253  
 75 E\_rms Validation = 0.55582  
 76 E\_rms Testing = 0.54671  
 77

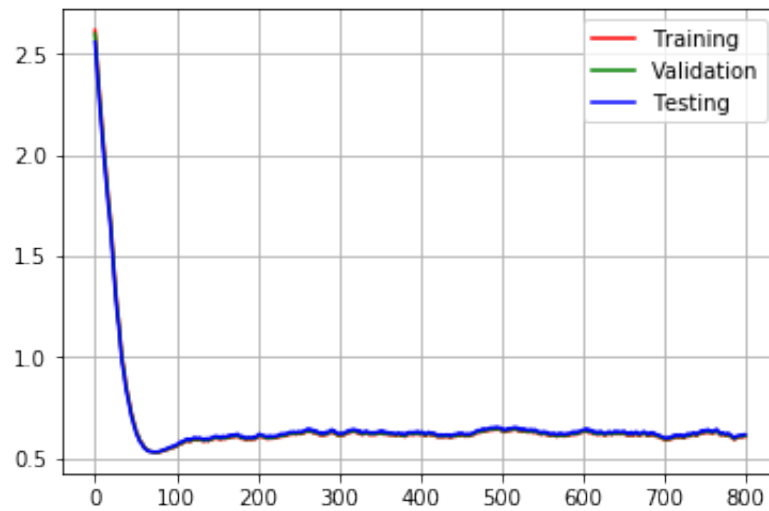


Figure 10: Subtraction Technique

78 E\_rms Training = 0.5284  
 79 E\_rms Validation = 0.52922  
 80 E\_rms Testing = 0.52997  
 81

82 **2.5 Logistic Regression on GSC dataset:**

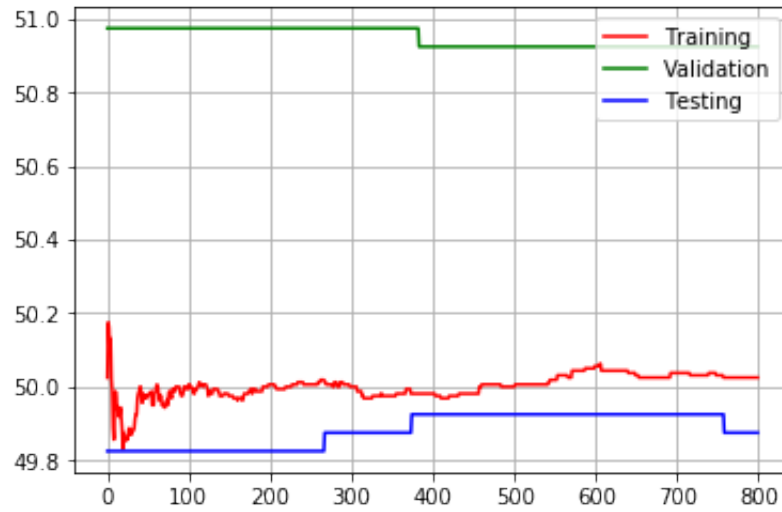


Figure 11: Concatenation Technique

83 Training\_accuracy: 50.14375  
84 Validation\_accuracy: 50.07504  
85 Testing\_accuracy: 48.77439  
86

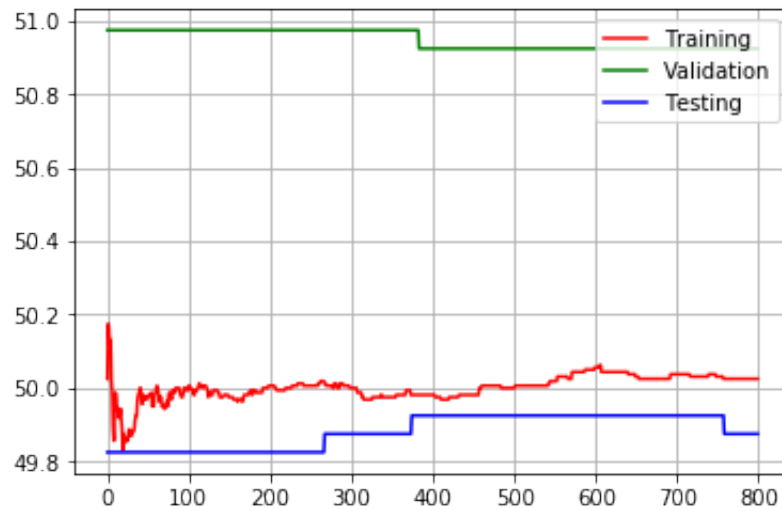


Figure 12: Subtraction Technique

87 Training\_accuracy: 50.175  
88 Validation\_accuracy: 50.97549  
89 Testing\_accuracy: 49.92496  
90

91 **2.6 Neural Network on GSC dataset:**

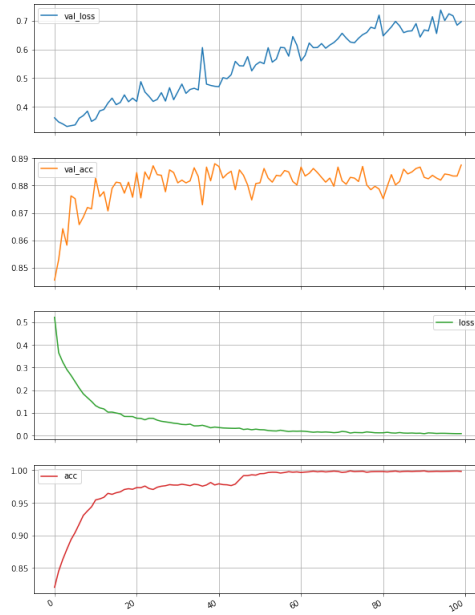


Figure 13: Concatenation Technique

92 Training\_accuracy: 0.9982  
 93 Validation\_accuracy: 0.8875  
 94 Testing\_accuracy: 0.8528  
 95

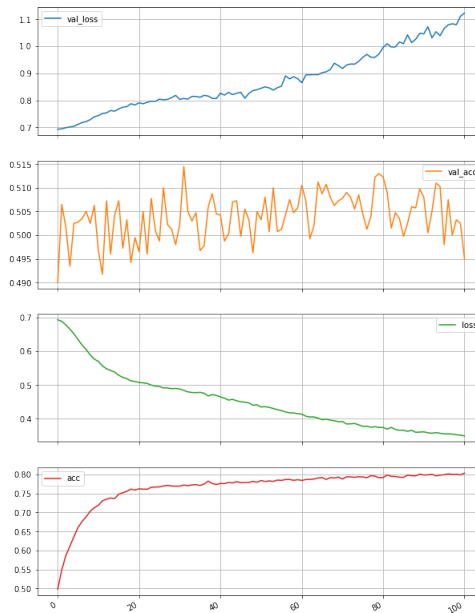


Figure 14: Subtraction Technique

96 Training\_accuracy: 0.8875  
 97 Validation\_accuracy: 0.4950  
 98 Testing\_accuracy: 0.5249  
 99



### 100 **3 references**

101 <https://www.coursera.org/learn/machine-learning>

102 Bishop - Pattern Recognition And Machine Learning - Springer 2006

[https://en.wikipedia.org/wiki/Machine<sub>l</sub>earning](https://en.wikipedia.org/wiki/Machine_learning)