

# CSE 474/574: Introduction to Machine Learning (Fall 2018)

Sargur N. Srihari  
University at Buffalo, The State University of New York  
Buffalo, New York 14260  
Contact: 716-645-6162 (O), srihari@buffalo.edu

September 28, 2018

## 1 Overview

You have been hired by the FBI to develop predictive models for detection of crime, your task is to help the Bureau and police departments to solve criminal cases dealing with evidence provided by handwritten documents such as wills and ransom notes. You are assigned to a forensic project by the FBI. The project requires you to apply machine learning 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  $y(x, w)$ .

Your task is to find similarity between the handwritten samples of the known and the questioned writer by using linear regression.

Each instance in the CEDAR “AND” training data consists of set of input features for each handwritten “AND” sample. The features are obtained from two different sources:

1. *Human Observed features*: Features entered by human document examiners manually
2. *GSC features*: Features extracted using Gradient Structural Concavity (GSC) algorithm.

The target values are scalars that can take two values  $\{1:\text{same writer}, 0:\text{different writers}\}$ . Although the training target values are discrete we use linear regression to obtain real values which is more useful for finding similarity (avoids collision into only two possible values).

## 2 Dataset

### 2.1 Source of Dataset

Our dataset uses “AND” images samples extracted from CEDAR Letter dataset. The CEDAR dataset consists of handwritten letter manuscripts written by 1567 writers. Each of the writer has copied a source document thrice. Hence, there are a total of 4701 handwritten letter manuscripts. Each letter has vocabulary size of 156 words (32 duplicate words, 124 unique words) and has one to five occurrences of the word “AND” (cursive and hand printed). Image snippets of the word “AND” were extracted

from each of the manuscript using transcript-mapping function of CEDAR-FOX. The total number of “AND” image fragments available after extraction are 15,518. Figure 1. shows examples of the “AND” image fragments.

Figure 1: Example of Dataset

Sample ID [XXXXy_numZ]	0001a_num1	0001a_num2	0002a_num1	0002a_num2	0005b_num1	0005b_num2	1121a_num1	1121a_num2	1160a_num1	1160a_num2
Writer Number [XXXX]	Writer 0001	Writer 0001	Writer 0002	Writer 0002	Writer 0005	Writer 0005	Writer 1121	Writer 1121	Writer 1160	Writer 1160
Page Number [y]	Page 1	Page 1	Page 1	Page 1	Page 2	Page 2	Page 1	Page 1	Page 1	Page 1
Sample Number [Z]	Sample 1	Sample 2	Sample 1	Sample 2	Sample 1	Sample 2	Sample 1	Sample 2	Sample 1	Sample 2

## 2.2 Types of Datasets:

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

### 2.2.1 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. The description of each of the Human Observed features are given in Table 1. There are total of 18 features for a pair of handwritten “AND” sample (9 features for each sample). The dataset is named as “*HumanObservedDataset.csv*”. Dataset available on UBLearn. The entire dataset consists of 512,345 handwritten sample pairs (rows), each having 2 image ids, 18 features and a target value. Figure 2. shows two sample rows from human observed dataset:

Figure 2: Human Observed Dataset Example

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

Table 1: Feature Description for Human Observed Dataset

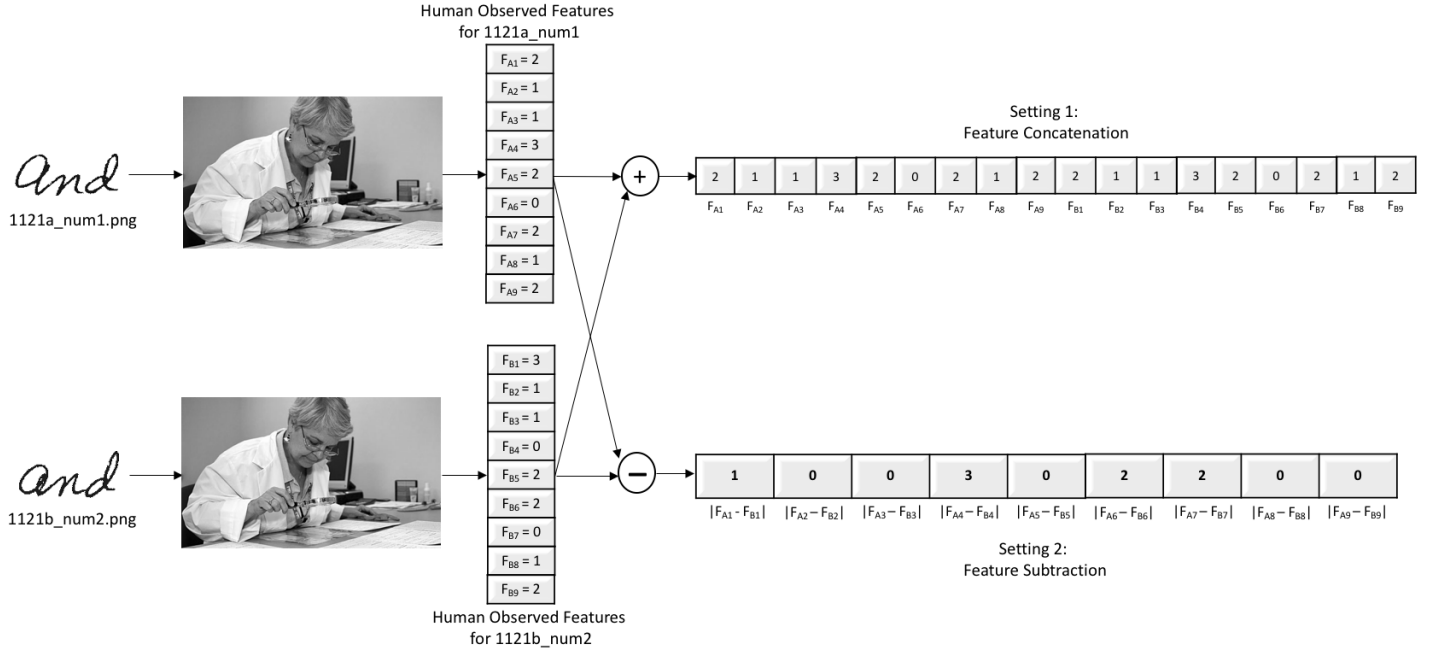
Initial stroke of formation of $a$ ( $x_1$ )	Formation of staff of $a$ ( $x_2$ )	Number of arches of $n$ ( $x_3$ )	Shape of arches of $n$ ( $x_4$ )	Location of mid-point of $n$ ( $x_5$ )	Formation of staff of $d$ ( $x_6$ )	Formation of initial stroke of $d$ ( $x_7$ )	Formation of terminal stroke of $d$ ( $x_8$ )	Symbol in place of the word <i>and</i> ( $x_9$ )
Right of staff (0)	Tented (0)	One (0)	Pointed (0)	Above baseline (0)	Tented (0)	Overhand (0)	Curved up (0)	Formation (0)
Left of staff (1)	Retraced (1)	Two (1)	Rounded (1)	Below baseline (1)	Retraced (1)	Underhand (1)	Straight across (1)	Symbol (1)
Center of staff (2)	Looped (2)	No fixed pattern (2)	Retraced (2)	At baseline (2)	Looped (2)	Straight across (2)	Curved down (2)	None (2)
No fixed pattern (3)	No staff (3)		Combination (3)	No fixed pattern (3)	No fixed pattern (3)	No fixed pattern (3)	No obvious ending stroke (3)	
	No fixed pattern (4)		No fixed pattern (4)				No fixed pattern (4)	

Figure 3. describes the two settings under which you need to perform linear regression

Setting 1: Feature Concatenation [18 features]

Setting 2: Feature subtraction [9 features]

Figure 3: Feature Extraction for Human Observed Dataset



### 2.2.2 GSC Dataset using Feature Engineering

Gradient Structural Concavity algorithm generates 512 sized feature vector for an input handwritten “AND” image. GSC algorithm extracts 192 binary gradient features, 192 binary structural features and 128 concavity features. There are total of 1024 features for a pair of handwritten “AND” sample (512 features for each sample). The dataset is named as “*GSCDataset.csv*”. Dataset is available on UBLearn. The entire dataset consists of 512,345 handwritten sample pairs (rows), each having 2 image ids, 1024 features and a target value. Figure 4. shows two sample rows from human observed dataset:

Figure 4: GSC Dataset Example

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 3. describes the two settings under which you need to perform linear regression

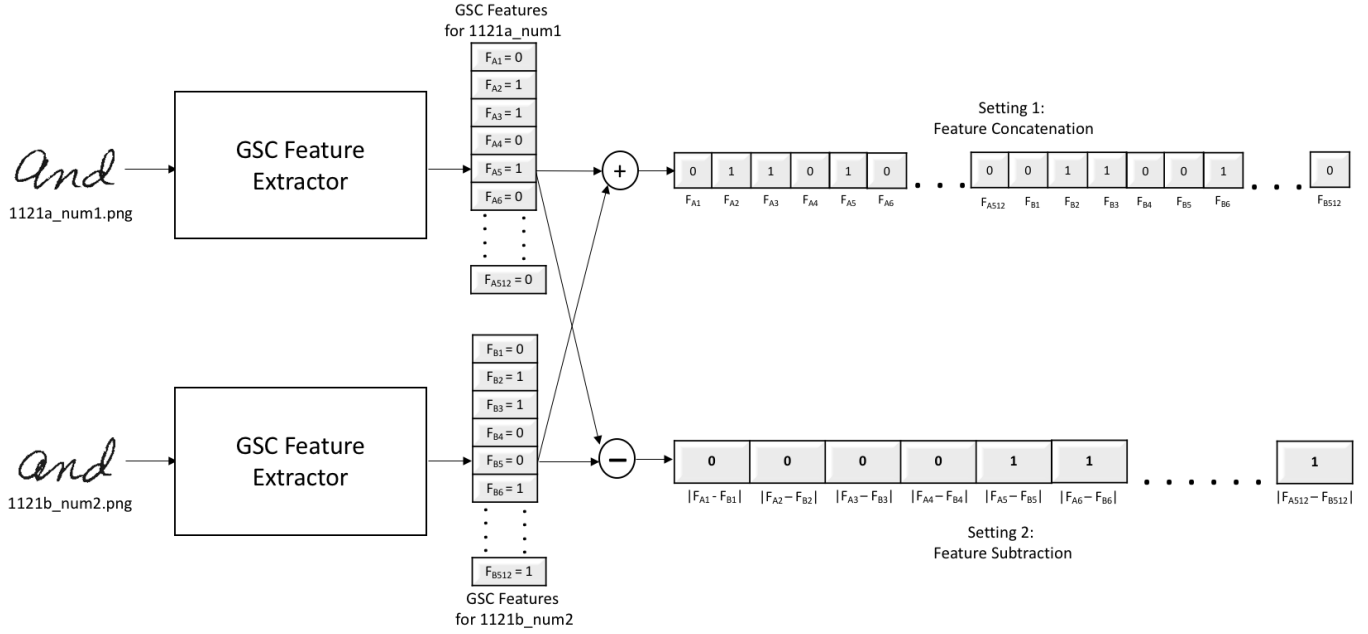
Setting 1: Feature Concatenation [1024 features]

Setting 2: Feature subtraction [512 features]

## 3 Plan of Work

1. **Extract features values and Image Ids from the data:** Process the original CSV data file into a Numpy matrix or Pandas Dataframe.
2. **Data Partitioning:** We use Unseen Writer Partitioning, in this scheme there exists no writer which is present in both the training (Tr) and testing (Ts) writer set simultaneously. Hence, any

Figure 5: Feature Extraction For GSC Dataset



test writer would not be a part of training set and vice-versa.

$$T_r \cap T_s = \emptyset \quad (1)$$

3. **Train using Linear Regression:** Use Gradient Descent to train the regression model using a group of hyperparameters. There are 4 different input dataset for linear regression:
  - (a) Human Observed Dataset with feature concatenation
  - (b) Human Observed Dataset with feature subtraction
  - (c) GSC Dataset with feature concatenation
  - (d) GSC Dataset with feature subtraction
4. **Tune hyper-parameters:** Validate the regression performance of your model on the validation set. Change your hyper-parameters and repeat step 4. Try to find what values those hyper-parameters should take so as to give better performance on the validation set
5. **Test your machine learning scheme on the testing set:** After finishing all the above steps, fix your hyper-parameters and model parameter and test your models performance on the testing set. This shows the ultimate effectiveness of your models generalization power gained by learning

## 4 Evaluation

Evaluate your solution on a test set using Root Mean Square (RMS) error, defined as

$$E_{RMS} = \sqrt{2E(w^*)/N_V} \quad (2)$$

where  $w^*$  is the solution and  $N_V$  is the size of the test dataset.

## 5 Deliverables

There are two deliverables: report and code. After finishing the project, you may be asked to demonstrate it to the TAs, particularly if your results and reasoning in your report are not clear enough.

### 1. Report

The report should describe your results and experimental setup. Submit the PDF on a CSE student server with the following script:

```
submit_cse474 proj2.pdf for undergraduates
```

```
submit_cse574 proj2.pdf for graduates
```

### 2. Code

The code for your implementation should be in Python only. You can submit multiple files, but the name of the entrance file should be `main.py`. Please provide necessary comments in the code. Python code, training and testing files should be packed in a ZIP file named `proj2code.zip`. Submit the Python code on a CSE student server with the following script:

```
submit_cse474 proj2code.zip for undergraduates
```

```
submit_cse574 proj2code.zip for graduates
```

## 6 Scoring Rubric

Conceptual Understanding in Report: 30%

Results in Report - Graphs, Tables etc: 40%

Report Formatting: 5%

Python Code Understanding (Provide comments in your Python Code): 25%