
Writer Identification from Handwriting

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A decade ago, Handwriting recognition was a very difficult task but now due
2 advancements in Machine Learning and access powerful computing power , it
3 has become a simple toy problem. Identifying character in a handwritten text is
4 quite easily doable but the problem that we are tackling is slightly different, we try
5 to identify the writer given two handwritten texts, by comparing them through a
6 machine learning model whether they belong to the same writer or not. We can
7 test this for multiple known writers. We chose three domains of machine learning
8 to compare and solve this problem. We first use a Multinet to get log likelihood
9 ratio of the features being from the same writer. Secondly we used a SIFT feature
10 extractor and Siamese Network. Finally we used a modern approach to solve this
11 problem which is Deep Siamese Network.

12 1 Introduction

13 In this project we use the AND DATASET which has two forms , one with human-engineered features
14 and other with images of hand-written text 'and'. We use the human engineered features to build our
15 two Bayesian Networks where one would represent the two features are from the same writer and the
16 other would represent they are from different writer. Using this knowledge of probability obtained
17 from Bayesian Networks we can find out whether the two features belong to handwriting of the same
18 writer or not.

19 2 PGM

20 In this approach, we try to predict whether a given sample pair belongs to the same writer or different
21 writer using Probabilistic Graphical Model, specifically a Bayesian Network Model also known as
22 Belief Network. Bayesian Networks are DAG's used to represent knowledge of a uncertain domain
23 using parameters of the domain. Each node in the network is a random variable and edges represent
24 conditional dependencies and in essence capture the causal relations between each node. For a given
25 set of random variables their joint probability can be calculated using a Bayesian Network using a
26 factorization of the Bayesian Network.

27 2.1 Dataset

28 The dataset used for this approach is in the form of features extracted by the human experts. For each
29 handwriting sample there are 9 features that define the characteristics of that sample such as initial
30 stroke, number and shape of arches. There are total 1026 writers and total 3524 handwritten samples
31 whose features are provided. From this dataset we are creating 1581 positive pair samples (pairs in
32 which both samples belong to the same writer) and 5000 negative pair samples (pairs in which each
33 sample belongs to different writer).

34 2.1.1 Positive Pair Generation

35 For a writer we take two feature vectors in the dataset where the writer id are the same. Once first
36 feature vector is chosen the other one is chosen at random to avoid choosing the exactly same feature
37 vector twice. This is done for all writers.

38 2.1.2 Negative Pair Generation

39 We chose a writer and then we choose another different (ID must be different) writer, We take any
40 two feature vectors from these two writer as negative pair. We do this for all writers.

41 2.2 Approach

42 Multinet:

43 Bayesian networks offer a compromise between the two extremes by encoding independence
44 when possible and dependence when necessary. They allow a wide spectrum of independence
45 assertions to be considered by the model builder so that a practical balance can be established
46 between computational needs and adequacy of conclusions.[5] When the independence relations
47 change wrt certain parameters then just single bayesian network fail to represent both independencies
48 properly i.e single bayesian network can not handle asymmetric independencies.

49
50 So we use two bayesian networks to capture the variation in dependencies wrt a single
51 parameter(hypothesis) In our case we will have one Bayesian network of 18 nodes where each node
52 is a variable of the feature vector of a writer since there are two writers we will have 18 random
53 variables. This Bayesian Network will be formed under the hypothesis that the 18 features belong
54 to the same writer. The other Bayesian network will be formed using the hypothesis that the 18
55 features(9 from writer1 and 9 from writer2) are from different writer.

56
57 Joint Probability can be calculated using the product of the conditional probability distribu-
58 tion of each node. The joint probability of 18 variables from both Bayesian networks would give the
59 probability of the 18 values being from the same writer and different writer. The log of the ratio of
60 these probabilities will give us the log likelihood ratio.

61 2.3 Implementation

62 For implementation of PGM approach, we have used pandas and pgmpy. We are using *pandas.DataFrame*
63 to group all positive pairs in one dataframe and all negative pairs in another. These
64 dataframes have 18 columns such that first 9 columns represent the features of first image whereas
65 the remaining represent the features of second image.

66 There are 2 broad techniques for structure learning, constraint based structure learning and score
67 based structure learning. Content based structure learning is more useful when using some hypothesis
68 tests, one can identify the dependencies in the structure and use these dependencies to construct
69 Bayesian network. Since there are no hypothesis tests that can identify dependencies in gives dataset,
70 we will use score based estimation techniques.

71 To get the Bayesian network, we have used 3 different types of score based estimators, Bdeu Score,
72 BIC score, K2 score.

73 **Bdeu Score:** Bdeu score stands for Bayesian Dirichlet equivalent uniform score. The scores take on
74 the most adversary and the most beneficial priors among those within a contamination set around the
75 symmetric one. [1]

76 **BIC score:** Bayesian information criterion (BIC) or Schwarz criterion is a criterion for model
77 selection among a finite set of models; the model with the lowest BIC is preferred. It is based, in
78 part, on the likelihood function. When fitting models, it is possible to increase the likelihood by
79 adding parameters, but doing so may result in overfitting. BIC attempts to resolve this problem by
80 introducing a penalty term for the number of parameters in the model. [2]

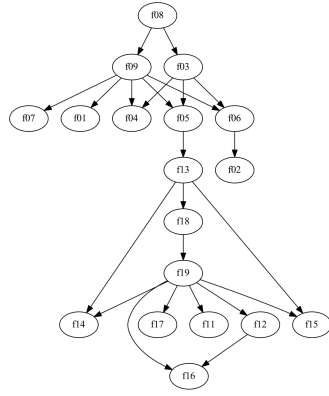


Figure 1: Bayesian Network for Same Writer (BDeu Score)

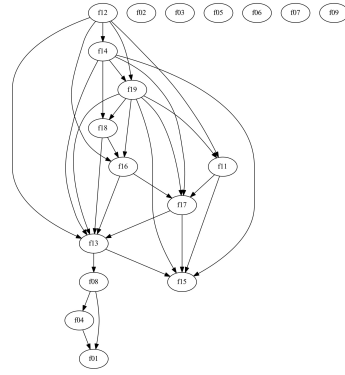


Figure 2: Bayesian Network for Different Writer (BDeu Score)

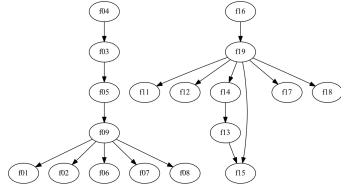


Figure 3: Bayesian Network for Same Writer (BIC Score)

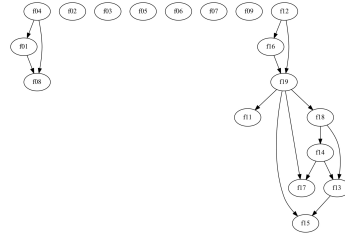


Figure 4: Bayesian Network for Different Writer (BIC Score)

81 **K2 score:** The K2-score is a greedy search method, which incrementally adds a node to a parent
 82 set and finds the best parent set to maximize the joint probability of the structure and the database. [3]

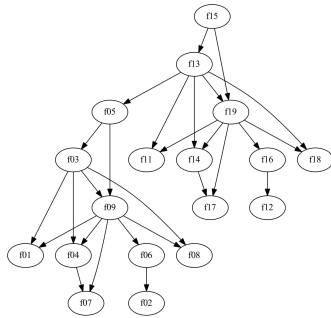


Figure 5: Bayesian Network for Same Writer (K2 Score)

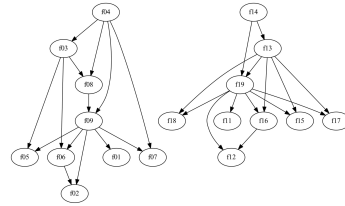


Figure 6: Bayesian Network for Different Writer (K2 Score)

83 On each of these score estimators we performed local hill climbing search to find the Bayesian model
 84 that has optimal score with respect to the scoring method supplied in the constructor. [4]
 85 Out of all these models, it is found that K2 score estimator gives the most optimum model. Hence the
 86 final model used is the one found using K2 estimator.

87 Using this model we get the joint probabilities from the product of CPDs for observed values of
88 features. If log ratio of a given pair is positive, then the samples in the pair are considered to belong
89 to the same writer else the samples belong to different writers.

90 3 Simple Machine Learning

91 3.1 Dataset

92 The dataset given for this approach consists of images of handwritten word 'and'. There are total
93 1568 writers and 15518 images in total written by the writers.

94 3.1.1 PreProcessing

95 We arrange the given feature dataset in a hash-table like fashion where key is the writer number and
96 value is a list of features of all those images where writer is the key(writer)

97 3.1.2 Positive Pair Generation

98 For a writer we take two feature vectors in the value-space in the hashtable of the writer. Once first
99 feature vector is chosen the other one is chosen at random to avoid choosing the exactly same feature
100 vector twice. This is done for all writers.

101 3.1.3 Negative Pair Generation

102 We chose a writer and then we choosing a non-zero random increment, using this increment we
103 choose another writer, and since the increment is non-zero we can not choose the same writer. We
104 take any two feature vectors from these two writer as negative pair. We do this for all writers.

105 3.2 Approach:

106 Feature Extraction:

107 In this approach we have extracted SIFT features from the images and these features won't be of
108 same length due to the nature of SIFT features and this will cause trouble from implementing neural
109 network because it would require input of the same size to train on.

110 To tackle this problem we form an array of SIFT descriptors from all the images. Then we find the
111 cluster centers using K-means [Note : The number clusters specified will be the dimensionality of the
112 feature vector] using the array of descriptors as data points and we store the centres in a file for later
113 use. We can now cluster the data points according to the centers obtained. The population of descriptors
114 for an image in a cluster for every cluster forms the feature vector.

115 In simple words the feature vector is the histogram of population of descriptors for an image in each
116 cluster.

117 Neural Network:

118 To determine whether given pair of features belongs to same writer or not we calculate the eu-
119 clidean distance between the two feature vector and the network updates weights to decrease/increase
120 this euclidean distance depending upon the known labels. If the features are from the same writer
121 we try to decrease the distance otherwise we try to increase the distance. For this updation we are
122 making use of contrastive loss function which uses the following formula:

$$124 \quad L(W, Y, X_1, X_2) = (1 - Y) \frac{1}{2} (D_w)^2 + (Y) \frac{1}{2} \{ \max(0, m - D_w) \}^2 \quad [6]$$

125 where D_w is the Euclidean distance

126
127 Since the labels are binary, in the above function only one part can be non-zero because
128 other one has to be zero, this then changes weights in order to favour/hinder distance. The objective of
129 this architecture is not to classify input images, but to differentiate between them. So, a classification
130 loss function (such as cross entropy) would not be the best fit. Instead, this architecture is better
131 suited to use a contrastive function.[7]

132 3.3 Implementation:

133 We implement a Siamese network which learns the similarity of the two feature vectors. We use
 134 Keras library to handle all the gradient operations and model architecture. We concatenate the two
 135 inputs in a single vector one after the other and while calculating the distance we calculate the
 136 distance between first half and the second half this helps in keeping the model simple and to avoid
 two identical branches. The model architecture is as follows:

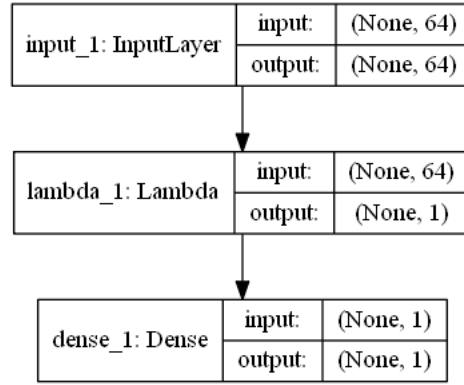


Figure 7: Simple Machine Learning Model Architecture

137

138 4 Deep Siamese Network

139 4.1 Dataset:

140 The dataset given is same as that used from simple machine learning which consists of 15518 images
 141 of handwritten word 'and', contributed by 1568 writers.

142 4.1.1 PreProcessing

143 We arrange the given image dataset in a hash-table like fashion where key is the writer number and
 144 value is a list of features of all those images where writer is the key(writer)

145 4.1.2 Positive Pair Generation

146 For a writer we take two images in the value-space in the hashtable of the writer. Once first feature
 147 vector is chosen the other one is chosen at random to avoid choosing the exactly same feature vector
 148 twice. This is done for all writers.

149 4.1.3 Negative Pair Generation

150 We chose a writer and then we choosing a non-zero random increment, using this increment we
 151 choose another writer, and since the increment is non-zero we can not choose the same writer. We
 152 take any two images from these two writer as negative pair. We do this for all writers.

153 4.2 Approach:

154 For deep learning approach, we are going to make use of Siamese Architecture. Since the problem
 155 requires to compare 2 samples of images and determine whether they belong to the same writer, we
 156 are going to use Siamese Neural Network. We have two branches of Convolutional Neural Network
 157 which will work as automatic feature extractors and the outputs of these sister branches are then
 158 merged in a lambda layer of Euclidean Distance and to determine whether the two samples belong to
 159 same writer or not. A contrastive loss function is used to update weights in order to favour/hinder the
 160 distance observed.

161

$$L(W, Y, X_1, X_2) = (1 - Y) \frac{1}{2} (D_w)^2 + (Y) \frac{1}{2} \{ \max(0, m - D_w) \}^2 \quad [6]$$

4.3 Implementation:

Siamese Neural Network: A Siamese networks consists of two identical convolutional neural networks, each taking one of the two input images. The output of the two networks are then fed to a euclidean distance function , which calculates the similarity between the two images. The loss function to update weights is the contrastive loss function There are two sister networks, which are identical neural networks, with the exact same weights. Each image in the image pair is fed to one of these networks.[7]

Contrastive Loss Function: The objective of the siamese architecture is not to classify input images, but to differentiate between them. So, a classification loss function (such as cross entropy) would not be the best fit. Instead, this architecture is better suited to use a contrastive function. Intuitively, this function just evaluates how well the network is distinguishing a given pair of images.[7]

In this approach for feature extraction, we are making use of Convolutional Neural Network

Convolutional Neural Network: Convolutional Neural Networks are very similar to ordinary Neural Networks they are made up of neurons that have learnable weights and biases. Each neuron receives some inputs, performs a dot product and optionally follows it with a non-linearity. The whole network still expresses a single differentiable score function: from the raw image pixels on one end to class scores at the other. A simple CNN is a sequence of layers, and every layer of a CNN transforms one volume of activations to another through a differentiable function. There are three main types of layers used to build CNN architectures: Convolutional Layer, Pooling Layer, and Fully-Connected Layer. Convolutional layer will compute the output of neurons that are connected to local regions in the input, each computing a dot product between their weights and a small region they are connected to in the input volume. Pooling layer will perform a downsampling operation along the spatial dimensions (width, height). Finally the fully connected layer gives the class score, thereby deciding whether given sample belongs to the same writer or not. [8]

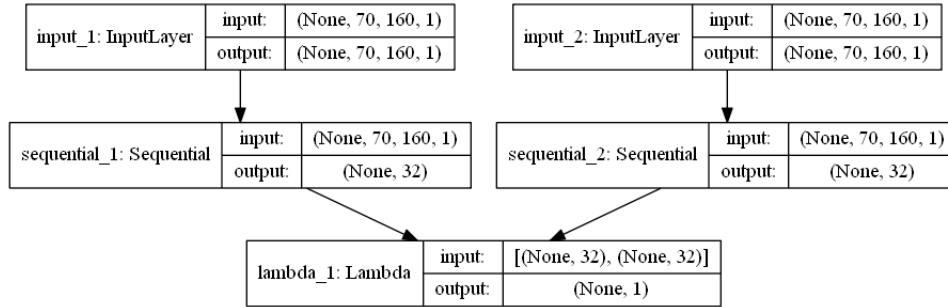


Figure 8: Deep Siamese Model Architecture

5 Conclusion

We can see from the graph [Figure 8] and the table below that the approach involving Deep Siamese Network yields the best results. But it should be also considered that the Deep Learning model required a lot more of computation power than the other approaches. The PGM approach required the least computation power (excluding Bayesian Structure Learning as that can be done manually). Overall the simple machine learning model's performance was in the Goldilocks zone as its accuracy is also reasonable and it didn't require much computing power.

Table 1: Accuracies of each method

Approach	Accuracy
PGM [Mulinet]	54.15%
Simple Machine Learning	63.23%
Deep Siamese Network	82.53%

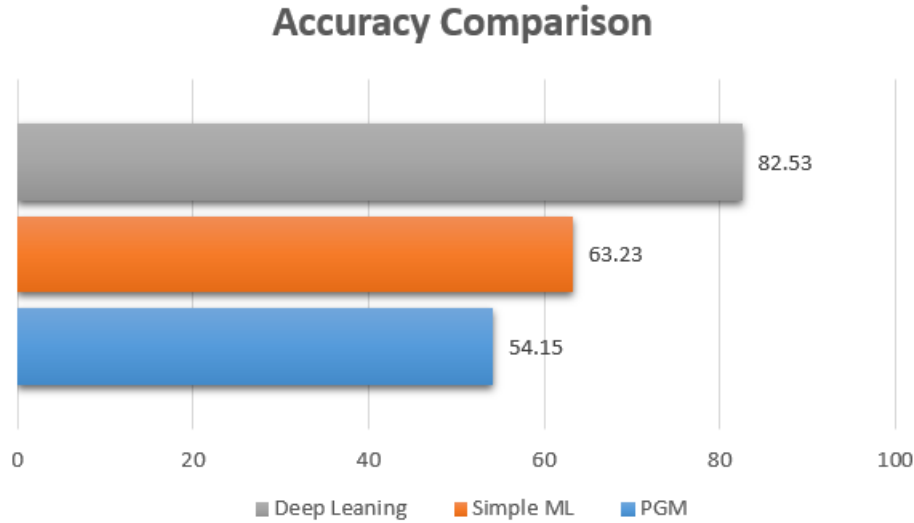


Figure 9: Accuracy

Acknowledgments:

We are thankful to Professor Sargur N. Srihari and teaching assistants, Mihir Chauhan and Mohammad Shaikh for giving us this opportunity and guiding us throughout the course of the project.

References

- [1] Mauro Scanagatta, Cassio P. de Campos, and Marco Zaffalon *Min-BDeu and Max-BDeu Scores for Learning Bayesian Networks*
- [2] https://en.wikipedia.org/wiki/Bayesian_information_criterion
- [3] Shulin Yang and Kuo-Chu Chang *Comparison of Score Metrics for Bayesian Network learning*
- [4] <http://pgmpy.org/estimators.html#hill-climb-search>
- [5] <https://arxiv.org/ftp/arxiv/papers/1303/1303.5718.pdf>
- [6] Raia Hadsell, Sumit Chopra, Yann LeCun *Dimensionality Reduction by Learning an Invariant Mapping*, The Courant Institute of Mathematical Sciences, New York University, 719 Broadway, New York, NY 1003, USA.
- [7] <https://hackernoon.com/one-shot-learning-with-siamese-networks-in-pytorch-8ddaab10340e>
- [8] <http://cs231n.github.io/convolutional-networks/>