

Exploring Similarity Measures for Biometric Databases

Praveer Mansukhani, Venu Govindaraju

Center for Unified Biometrics and Sensors (CUBS)
University at Buffalo
{pdm5, govind}@buffalo.edu

Abstract. Currently biometric system performance is evaluated in terms of its FAR and FRR. The accuracy expressed in such a manner depends on the characteristics of the dataset on which the system has been tested. Using different datasets for system evaluation makes a true comparison of such systems difficult, more so in cases where the systems are designed to work on different biometrics, such as fingerprint and signature. We propose a similarity metric, calculated for a pair of fingerprint templates, which will quantize the “confusion” of a fingerprint matcher in evaluating them. We show how such a metric, can be calculated for an entire biometric database, to give us the amount of difficulty a matcher has, when evaluating fingerprints from that database. This similarity metric can be calculated in linear time, making it suitable for large datasets.

1 Introduction

Biometric-based systems, in their various forms, have become popular over the last few years. Unlike various token or knowledge based systems, they cannot be easily compromised and are used in a variety of low or high security applications. Most systems in place use one biometric (such as facial image, fingerprint data, etc) to correctly identify or authenticate a candidate, though multi-biometric systems have also been developed in which a combination of two or more user characteristics can be used, either in serial or parallel.

The existing biometric based systems can be classified into one of two classes – identification based or authentication based systems. In identification based systems, a biometric template is provided to the system, which must then identify its owner based on the templates of various candidates stored in its database. This can be likened to a 1:N matching. The system outputs the predicted identity of the owner of the template, and in most cases a score, which indicates the confidence level of the match. A simpler form of matching is of 1:1 type, which is carried in authentication based systems, where the claimed identity of the user is given to the system as input, along with the template data. The system must only retrieve the corresponding stored template of the claimed identity, and compare it with the given template. The score outputted by the system can be compared with a set threshold, to authenticate the user.

Biometric based authentication / identification systems are designed to work efficiently and accurately for a large number of users, typically a few thousands or even in some cases, up to a million users. Such systems have to handle large databases,

which may also store more than one template per user. This is especially true in the case of some behavioral based biometrics such as signatures, where two or more templates taken from a single user may significantly vary. Typically such systems require three to five templates per user. Even for a physical biometric such as fingerprint or facial image, where there will be lesser variation across templates (provided enrolled within a short time interval of one another), multiple templates, usually two or three are taken. It becomes imperative that any operations carried out, over the entire database, must be scalable to efficiently run over such large data.

Section 2 describes the existence of measures for estimating the complexity of speech databases and word lexicons. We have also explored currently used methods for estimating the performance of biometric systems and quality of biometric databases. In section 3, we have defined a similarity measure for fingerprint databases, which can serve as an indicator for estimating the difficulty / 'confusion' of a recognizer on a particular database. Section 4 shows how, once the similarity measure has been first defined for a pair of fingerprint templates, and the values computed over pairs of templates in the database, they can be used calculate the similarity metric for the entire database. Section 5 describes the results of various experiments performed to support our theory.

2 Current Work

2.1 Existing Evaluation Measures

Measures to evaluate performance of speech and handwriting recognition systems, on the basis of complexity of the language models, have been researched extensively. Perplexity of a language model can be expressed as the average number of words that appear after a given word in that model. It has been used as a measure of complexity of the model. Nakagawa et al [5] have explored the relation between the size of the vocabulary of a model and the recognition rate of a speech recognition system. Furthermore, they have shown how, given prior knowledge of the performance of the recognizer and information about the perplexity of a language model, they can predict the recognition rate of that recognizer on that language model.

Marti et al [4][8] have used a HMM based handwritten text recognizer on different language models (simple, unigram and bigram) and shown that using a model with a smaller perplexity results in a better recognition rate.

A method to evaluate the difficulty of word recognition can be measured on the basis of the length of the words in the language model. (Grandidier et al [1]). It has been observed that a system is able to correctly recognize larger words better than smaller words. However, this measure cannot be used to predict the performance of the system.

Govindaraju et al [2] have defined 'lexicon density' as a measure of the distance between two handwritten words. It is an average measure of the total cost of edit operations required to convert one word in the lexicon to another. They have shown that a measure based on lexicon density is a better indication of the complexity of a word lexicon than measures such as lexicon size or edit distance.

2.2 Evaluation of Biometric System Performance

Evaluation of biometric based systems is done in terms of their False Accept Ratios (FARs) and False Reject Ratios (FRRs). [7] The FAR is the probability of a system evaluating two biometric templates of different people, to be of the same person. Correspondingly the FRR can be defined as the probability of the system to recognize two templates of the same person to belong to different people. The accuracy of the system could also be expressed in terms of the Equal Error Rate (EER), which is the value of the FRR, at a particular threshold, when it is equal to the FAR. These evaluation measures do not take into account other factors which could affect the performance of the system, such as the similarity of biometric templates or the quality of the template data.

Most biometric systems are tested using databases, of varying sizes, having templates belonging to different users. This makes it impossible to compare the performance of a system, having results on a particular database, with the performance of another, which has been tested on a different database. Moreover, there is no method to compare the performance of systems that use different biometrics to identify / authenticate a user. These two systems will be using totally different types of databases, for example one system may use signature data and the other may use face images.

Bolle et al [3] have attempted to classify fingerprint databases based on the quality of images present. They generate a similarity score using a one-to-one matching between fingerprints of the same user and have used a statistic based on the mean and standard deviation of the scores to separate them into three classes: good, bad and ugly. However, visual inspection of the plot is needed to accurately separate the classes and classify the database. Moreover, this technique is primarily a measure of the quality of templates in the database, and does not give any indication of the difficulty of the recognizer in distinguishing one fingerprint image from another.

3 Calculating Similarity measure for a pair of fingerprints

Currently used fingerprint based biometric systems store a user's fingerprint data not as an image but in the form of a template which is its compact form. The template is considered to be an accurate representation of the users biometric. Matching two fingerprints do not involve the actual comparison of the fingerprint images, but a matching procedure between the two corresponding templates. Typically fingerprint templates contain spatial information about the minutiae detected in the images, and are represented as a list of the extracted minutiae. Fingerprint images typically contain between 30- 50 minutiae, although prints with up to 100 minutiae have also been observed. For a match between two fingerprints a match of 6-8 minutiae points are usually considered sufficient; however up to 12 matching points are needed for some particular law-enforcement applications. [10]

For calculating a similarity metric for a biometric database, we first define a similarity measure between two fingerprint images, and then average it out over all the images in the database. The similarity measure does not indicate whether the two templates belong to the same person, rather it is an indicator of the level of difficulty

of a recognizer in comparing the two templates. Hence, though two templates which are very similar to one another will most likely belong to the same user, such may not always be the case. Similarly, two very dissimilar templates will more often than not, be of different persons, but not every time. A pair of templates having an extreme similarity value ie. either too high or too low should be correctly classified by the biometric system with ease, and pairs with intermediate values are going to take greater computation from a matcher to classify.

Consider a fingerprint template $M = \{m_1, m_2, \dots\}$ where m_i is the i^{th} detected minutiae point. Each point m_i can be expressed as a triple (x_i, y_i, θ_i) denoting the (x,y) pixel location of the point and its orientation respectively. Here we have assumed each template to be extracted from a complete fingerprint image. The similarity measure between two fingerprint templates (M_a and M_b) can be calculated as follows

1. For each m_i in M_a calculate the corresponding m_j in M_b such that for all m_k in M_b , where $j \neq k$, $\text{dist}(m_i, m_j) < \text{dist}(m_i, m_k)$
2. Similarly, for each m_j in M_b calculate the corresponding m_i in M_a such that for all m_k in M_a , where $i \neq k$, $\text{dist}(m_j, m_i) < \text{dist}(m_j, m_k)$
3. Select $\{m^a_1, m^a_2, \dots, m^a_n\}$ from M_a such that for all m^a_x in $\{m^a_1, m^a_2, \dots, m^a_n\}$ exists corresponding m_j in M_b such that $(d^a_x = \text{dist}(m^a_x, m_j)) < \text{dist}(m_i, m_k)$ where m_i in M_a , but m_i not in $\{m^a_1, m^a_2, \dots, m^a_n\}$ and m_k in M_b such that m_k corresponds to m_i from step 1.
4. Similarly, select $\{m^b_1, m^b_2, \dots, m^b_n\}$ from M_b such that for all m^b_x in $\{m^b_1, m^b_2, \dots, m^b_n\}$ exists corresponding m_i in M_a such that $(d^b_x = \text{dist}(m^b_x, m_i)) < \text{dist}(m_j, m_k)$ where m_j in M_b , but m_j not in $\{m^b_1, m^b_2, \dots, m^b_n\}$ and m_k in M_a such that m_k corresponds to m_j from step 2.
5. Similarity $s(M_a, M_b) = (\sigma(d^a_1, d^a_2, \dots, d^a_n) + \sigma(d^b_1, d^b_2, \dots, d^b_n)) / 2$

It could be noted that a large value of the similarity score implies that the two templates are less similar to one another and thus less likely to confuse the recognizer. Conversely, a very small score implies that the two templates are highly similar, most probably belonging to the same person. Intermediate values indicate some possible confusion while matching. A perfectly similar template, which in practice is only going to occur when a template is matched against itself, will give us the similarity score 0.

4 Estimating the similarity metric for a fingerprint database

We can now estimate a similarity measure for a pair of fingerprints. This value indicates the difficulty of a fingerprint matching system in evaluating the two templates. A matcher finds it more difficult to identify the owner of a template from a dataset having a large number of similar prints. Conversely, if the dataset is such that the fingerprint templates stored are not very similar to one another, then the work of the matcher becomes easier. To quantize this level of difficulty, we define a similarity metric for a fingerprint database, calculated from the individual similarity measures for the stored templates, as follows:

$$S(\{M_1, M_2, \dots, M_n\}) = \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n s(M_i, M_j)}{n(n-1)/2}$$

Thus the similarity score for a database of fingerprints is defined as the average of scores calculated for all pairs of templates in the database.

4.1 Similarity measures for large fingerprint databases

We have established a metric for calculating the similarity of the templates in a fingerprint database. However, for a database of size n , calculation of this value will be done in time $O(n^2)$. Though it can easily be calculated for smaller datasets, it does not scale well to databases with a large value of n . Most commercially available matching systems have been designed to run on databases having thousands or even a few million templates. Calculating the similarity measure for such a large database is not feasible, and certainly not efficient.

We propose a randomization method by which a similarity metric can be efficiently estimated for a database in linear time. Consider a dataset $\{M_1, M_2, \dots, M_n\}$ having n fingerprint templates. The similarity measure can be calculated as follows:

1. Randomly select N fingerprint templates, where $N = \lfloor \sqrt{n} \rfloor$
2. For each of the N selected templates do:
 - a. Randomly select N different fingerprint templates
 - b. Calculate similarity measure between the template in (2) and each of the templates selected in (2a)
3. Calculate the mean of all such similarity measures as $S(\{M_1, M_2, \dots, M_n\})$

On comparing the distribution functions of the similarity of all the pairs of fingerprint templates and the distribution function for the similarity calculated as above, using a reduced number of computations, it is observed that they are almost equal, and thus have the same mean. This implies that such a method of calculating similarity of a large biometric database, while it runs in significantly less time than the 'brute-force' method, it gives us the same output.

5 Results

5.1 Description of fingerprint datasets and matcher.

We have used 2 datasets for our experiments, both obtained from CUBS. Dataset 1 (D1) is a smaller dataset consisting of 20 users, each user providing us with 3 templates. Dataset 2 is a larger dataset having 50 users, again 3 templates per user. Each

template stores the spatial information (x, y, θ) of each extracted minutiae point of the user.

The fingerprint matcher used is the standard matcher used at CUBS, and has comparable performance to any fingerprint matching software available commercially today. The output given by the matcher is a score between 0 and 1, where a 1 indicates a match between the input templates.

5.2 Analysis of template similarity metric based on matcher output

For each user, we have computed the similarity score for each pair of his templates (taking number of minutiae points compared $(n) = 10$). (Fig 1.) We have also given those templates to the CUBS fingerprint matcher to obtain a matching score. We have also computed similarity scores for all pairs of templates belonging to different users, and obtained the corresponding matching scores from the fingerprint matcher. Graphs of the similarity score of each template pair vs. the corresponding matcher output are plotted for both datasets.

Points closer to $(0,1)$ represent those pairs of templates which have a high similarity and have been identified by the CUBS matcher as belonging to the same person with a high confidence ratio. This is because fingerprints which are very similar more often than not belong to the same person. Points which have a high x value, are those which have been identified as dissimilar pairs of prints, and have been classified by the matcher as belonging to different people. However most points have a similarity score between 2 to 8, and these pairs are those which would not have been classified by the fingerprint matcher so easily.

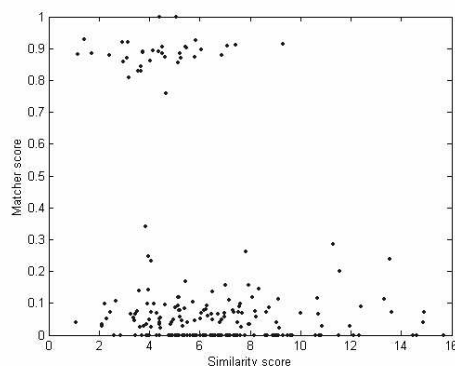


Fig. 1(a). Similarity score vs. Matcher score for individual templates (Dataset 1)

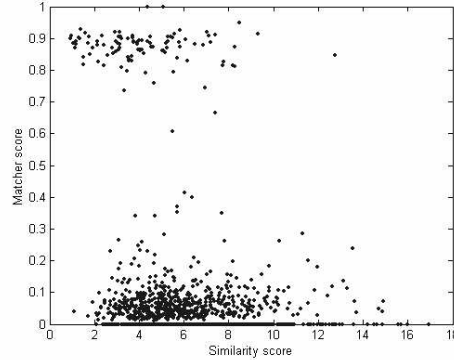


Fig. 1(b). Similarity score vs. Matcher score for individual templates (Dataset 2)

5.3 Applying the fingerprint similarity metric to datasets

We have estimated the similarity metric for both our datasets D1 and D2, as can be seen in Table 1. Recall that a larger value of the similarity measure means that the two templates are more dissimilar to one another. We can see that the database D1 has a greater similarity measure than D2, which implies that, a recognizer will find it more difficult to distinguish between two random templates belonging to D1 than for two templates belonging to D2. However, it may be noted that as D2 has a greater number of fingerprint templates, the size of the dataset will add to the confusion for the matcher. Thus, the similarity score, in conjunction with the size of the dataset (number of users) will be the most accurate way of representing the complexity of the dataset.

Database characteristic	Dataset D1	Dataset D2
Number of users	20	50
Total number of templates	60	150
Similarity measure	6.71	6.32

Table 1. Comparison of similarity metrics calculated for dataset D1 and D2

5.4 Scaling similarity metrics to large databases

We have estimated the similarity metric for two datasets (D1 and D2) having templates from 20 and 50 different users respectively, using a reduced number of computations as described in section 4.1. Table 2 shows that the values of the similarity metric, so estimated, are very close to the actual similarity metric calculated, using a ‘brute force’ method, involving all pairs of templates. However, this estimation requires a significantly less amount of time due to the reduced computations involved.

A running time of $O(n)$ indicates a large speedup over the previously required time of $O(n^2)$, making this method scalable to large databases.

Database characteristic	Database D1	Database D2
Number of users	20	50
Similarity measure calculated using whole database	6.71	6.32
Similarity measure on reduced data	6.58	6.43
Number of similarity comparisons for whole database	190	1225
Number of similarity comparisons using reduced data	~20	~50

Table 2. Comparison of similarity metric calculations for datasets D1 and D2 using whole datasets and reduced data

Figures 2(a) and 2(b) show the similarity score distribution so calculated using the whole set of templates, and the reduced dataset, for datasets D1 and D2 respectively. It is observed that for a particular dataset, the similarity scores, whether calculated for all pairs of templates or a reduced set have the same distribution, and hence the same mean, which is finally used to estimate the similarity metric for the entire database.

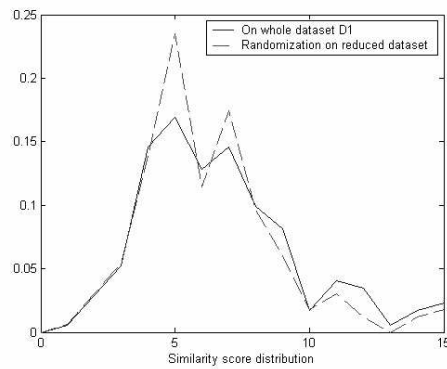


Fig. 2(a). Plot of similarity score distribution on whole dataset and score distribution calculated by using a reduced number of computations – Dataset 1

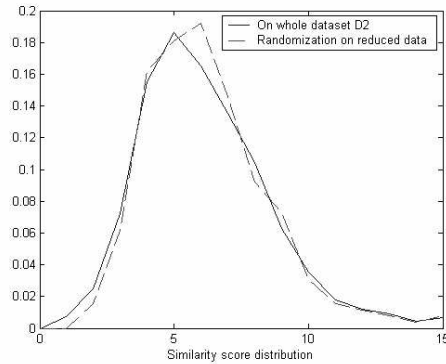


Fig. 2 (b). Plot of similarity score distribution on whole dataset and score distribution calculated by using a reduced number of computations – Dataset 2

6 Conclusions and Future Work

The similarity measure so computed for two templates gives an indication of the difficulty of the recognizer in evaluating them. When computed for the whole database, we can estimate the average “confusion” of the matcher in evaluating two templates from that database. This information, along with the size of the database, used in conjunction with the FAR and FRR values, can give us a better idea of the performance of a fingerprint matcher. Moreover this method can be used to evaluate the complexity of fingerprint datasets generated by artificial means, such as a synthetic fingerprint database generator [9]. An accurate estimation of such a similarity metric can be done in linear time, so as to enable scalability across very large datasets.

We need to develop and evaluate similar metrics for various other biometric systems, such as signature verification systems and hand geometry based biometric systems, among others. Just as we can evaluate the complexity of fingerprint databases, we can do so for other biometrics also. This may also enable us to compare the performance of biometric systems which work on different biometrics. Similarity measures could also be computed for multi-biometric based systems, where a combination of two or more biometrics are used.

Another direction of interest to us, is to further reduce the number of computations required to estimate the similarity score of a fingerprint database. Currently we have proved that we can accurately estimate such a template in linear time. It may be possible to compute the same in a lesser polynomial time, or even in logarithmic time. It remains to be seen how much, if any, loss of accuracy in the computations is caused due to the reduction in running time.

References

1. F. Grandidier, R. Sabourin, A. E. Yacoubi, M. Gilloux, C.Y. Suen: Influence of Word Length on Handwriting Recognition. Proc. Fifth International Conference on Document Analysis and Recognition, Sept 1999
2. H. Xue, V. Govindaraju: On the dependence of Handwritten Word Recognisers on Lexicons. IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol 24, No 12, December 2002
3. Ruud M. Bolle, Sharat Pankanti, Nalini K. Ratha: Evaluation techniques for biometric-based authentication systems(FRR)
4. U. -V. Marti, H. Bunke: On the influence of Vocabulary Size and Language Models in Unconstrained Handwritten Text Recognition. IEEE Proceedings 0-7695-1263-1/01
5. Seiichi Nagakawa, Isao Murase: Relationship between Phoneme Recognition Rate, Perplexity and Sentence Recognition and Comparison of Language Models. IEEE Proceedings 0-7803-0532-9/92
6. Ruud M. Bolle, Nalini K. Ratha, Sharat Pankanti: Evaluating Authentication Systems Using Bootstrap Confidence Intervals
7. Rajiv Khanna, Weicheng Shen: Automated Fingerprint Identification System (AFIS) Benchmarking Using National Institute of Standards and Technology (NIST) Special Database 4. IEEE Proceedings 0-7803- 1479-4/94
8. U. -V. Marti, H. Bunke: Unconstrained Handwriting Recognition: Language Models, Perplexity, and System Performance. L. R. B. Schomaker and L. G. Vuurpijl (Eds.), Proceedings of the Seventh International Workshop on Frontiers in Handwriting Recognition, September 2000, ISBN 90-76942-01-3
9. R. Cappelli, D. Maio, D. Maltoni: Synthetic Fingerprint-Database Generation. IEEE Proceedings 1051-4651/02
10. S. Pankanti, S. Prabhakar, A. K. Jain: On the Individuality of Fingerprints. IEEE Transactions on PAMI, Vol 24, No 8, 2002