

Development of Online Signature Recognition System

Outline

1. Introduction
2. Data
3. Feature Extraction
4. Performance Evaluation

1. Introduction

The objective of this session is to DEVELOP and EVALUATE an online signature recognition algorithm. According to the theory sessions, signature recognition systems can be divided into two categories:

- **Off-line:** the input is a static image of the signature.
- **On-line:** the signature is acquired using a specific digital sensor which includes the static images and dynamic signals related with the way the signature was done: x,y coordinates and pressure as a function of time.

Figure 1 shows a block diagram of a typical online signature recognition algorithm where $[x,y,p]$ are the captured signals by the sensor (Cartesian coordinates and pressure), f_t is the feature vector of the query signature to be compared with the f_c feature vector of the signature stored in the database (claimed identity).

In this session we will assume that the data is available (previously acquired) and we will focus on the development of two modules:

- Feature Extraction Module.
- Matcher.

You must complete the tasks proposed in this document and answer the questions included.

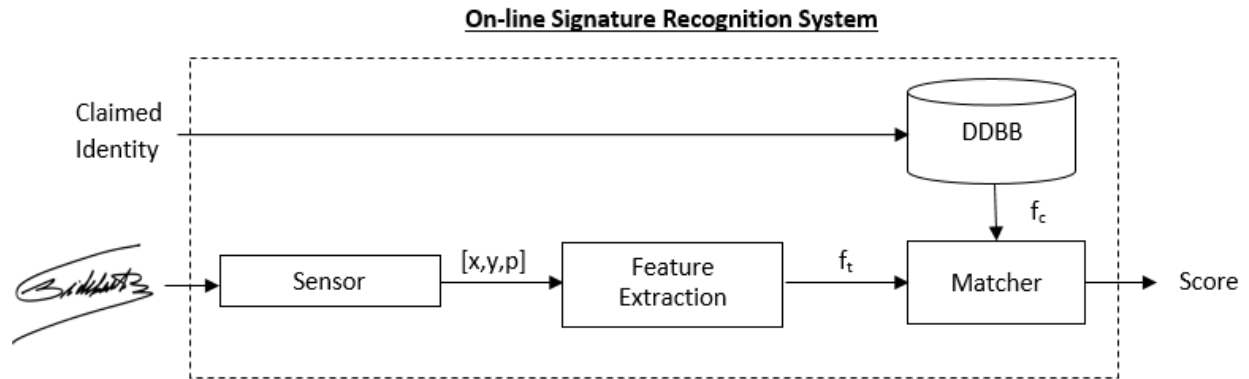


Figure 1. Block Diagram of a typical online signature recognition system

2. Data

For the practice we will use 50 users from the BiosecurID database. Each of the users have 28 signatures acquired in 4 sessions with a time lapse of 2 months. From the 28 signatures, 16 are genuine (4 per session) and 12 are forgers (3 per session). In this practice we will only consider the genuine signatures.

Each of the signatures is stored in .mat file which contains three vectors of same length with the x, y coordinates and the pressure as functions of time.

The formatting of the files is uXXXXsYYYY_sgZZZZ.mat:

- XXXX: user number
- YYYY: session number
- ZZZZ: signature number

The GENUINE signatures of each session are those with ZZZZ=[0001,0002,0006,0007].

The signatures with ZZZZ=[0003,0004,0005] are the FORGERS and they will NOT be used in this practice.

QUESTION. Choose a signature (from a random user) and show (assuming that the sensor has a 200 samples/second acquisition rate):

- Signal x as a function of signal y.
- Signal x as a function of time.
- Signal y as a function of time.
- Signal p as a function of time.

x as a function of y

x as a function of t

y as a function of t

p as a function of t

Repeat the task with another signature of the same user.

QUESTION: are the different signals reasonable? Are they the same length? Why?

3. Feature Extraction

The comparison of signals with different lengths is not trivial. Therefore, we will extract 5 global parameters of each of the signatures. So, all signature will be represented by a feature vector with fixed size equal to 5. These parameters are:

- Total duration of the signature: T
- Number of *pen-up* (number of times the pen was lifted). It means the number of times (not the number of samples) that p is equal to 0.
- Duration of *pen-down* (signal p is different to 0) T_d divided by the total duration T : T_d/T
- Number of maximums and minimums of signal x .
- Number of maximums and minimums of signal y .

You have to develop 4 functions to extract each of the parameters:

- $T=T_{total}(x)$
- $N_{pu}=N_{penups}(p)$
- $T_{pd}=T_{pendown}(p)$
- $N_{maxx}=N_{maxima}(x)$ (analogous, $N_{maxy}=N_{maxima}(y)$).

According to those functions, we will develop a new function with input data (x,y,p) of a given signature and output data the feature vector containing the 5 parameters ($FeatVect=featureExtractor(x,y,p)$).

Based on your function `featureExtractor` you have to develop a program (`ProcessBiosecurID.m`) to extract all the feature vectors from the database and store it in a matrix with 3 dimensions:

- Dimension 1: number of user (1:50)
- Dimension 2: number of signature (1:16)
- Dimension 3: number of parameter (1:5)

You have to save this matrix into the file `BiosecurIDparameters.mat`

Once you have the file `BiosecurIDparameters.mat`, you have to plot the distributions normalized between 0 and 1 (dividing by the total number of points of the distribution) for each of the 5 parameters.

You can use the Matlab functions `hist` and `histc`.

QUESTION: Plot the 5 distributions.

Total duration

N pen-ups

T pen-down / T

N maximums and minimums x

N maximum and minimum y

4. Performance Evaluation

We will evaluate the performance of our system according the number of signatures N in the enrollment set (N=1, N=4 and N=12).

The similarity score between a query/test signature and the enrollment signatures (signatures in the database) will be the Euclidean distance between feature vectors (vectors with 5 parameters). The final score will be the average score of the N comparisons (comparison between the query/test sample and the N enrollment samples).

You have to develop the function $\text{Score} = \text{Matcher}(\text{test}, \text{Model})$ where:

- Score: is the final score of the comparison.
- test: is the feature vector of the query/test signature (1x5)
- Model: is a matrix containing the feature vectors of the signatures enrolled in the database. Therefore, this matrix contains Nx5 values in which N is the number of signatures enrolled for the claimed identity.

There are two cases to be analyzed:

Genuine Scores: scores obtained when you compare a signature with his real enrolled identity (claimed identity = enrolled identity). So these users should be accepted by the system. For each user you will use N signatures as enrolled samples and the rest for testing:

- For N=1 we will have SG=15 genuine scores.
- For N=4 we will have SG=12 genuine scores.
- For N=12 we will have SG=4 genuine scores.
-

For each of the scenarios (N=1,4,12) you have to save all the genuine scores into a matrix (with dimension 50xSG). Each of the three matrixes will be stored into a .mat file with name: GenuineScores_N.mat.

Impostor Scores: scores obtained when you compare a signature with the enrolled samples of other users (claimed identity \neq enrolled identity). So these users should be rejected by the system. In this case, we will compare one signature of each user (the first one) with the models of the rest of the users (excluding the genuine case). Therefore, we will obtain SI=49 impostor scores for each user and each scenario (N=1,4,12).

For each scenario (N=1,4,12) these impostor scores will be saved into a matrix with dimensions 50xSI (50x49). Each of the three matrixes will be stored into a .mat file with name: ImpostorScores_N.mat.

Once we obtain the genuine and impostor scores, we will evaluate the performance of our system for each of the three scenarios (N=1,4,12) as a function of: FAR/FRR, EER and DET curves.

To obtain these performance metric you will have available the next functions:

[EER]=Eval_Det(GenuineScores, ImpostorScores, 'b')

- EER: value of the Equal Error Rate (error when FAR and FRR are equal)
- GenuineScores: the scores from target or genuine comparisons
- ImpostorScores: the scores from non target or impostors comparisons

QUESTION. Plot the performance graphics (DET curves) using the genuine and impostors score stored in their respective matrixes (for each of the scenarios N=1,4,12). Indicate the EER value.

N=1

N=4

N=12

QUESTION. According to the results, are they reasonable? What metrics are more illustrative? When do you obtain the best performance?