---------------------------------------------------------------------------------------------------------------------------------------------

# Implementing Neural Network in Keystroke Dynamics for a Better Biometric Authentication System

M.Kanimozhi, A.Kanimozhi

*Abstract*— Biometrics is method of measuring the biological features of the human being for the purpose of authorization and authentication. It takes human natural attributes/features as inputs. New method of authenticating right now. Keystroke dynamics is a behavioral biometric that is used to provide authentication during user working with the computer system. Besides the traditional way of providing authentication through password, which is the static authentication, there is a new method called dynamic authentication which recognizes the user past login. This paper is a survey about various techniques and algorithms used for providing dynamic authentication.

*Keywords*— Keystroke dynamics, biometric, Continuous Authentication, mono graph, digraph, neural network , FAR, FRR, EER.

## I. Introduction

Biometrics means "life measurement" but the term is usually associated with the use of unique physiological characteristics to identify an individual. The most important method to prevent from the unauthorized access is user authentication. User authentication means to verify the user id. This authentication is done by matching some indicator of user already that us registered for approve. Behavioral characters are related to the person that means what a person do or how the person uses the body. Now a day's managing multiple passwords, PIN numbers that are not easy to save in mind for a person. So using biometric to overcome these types of problem. Biometrics technologies used to verify and recognize the identity of the living person based on the behavioral characteristics. A number of biometric traits have been developed and are used to authenticate the person's identity. This is biometric authentication in which access is granted to users based on biological signatures such as fingerprint recognition, iris recognition or biometrics based on keystroke behavior. ID Control brings you **Keystroke**ID which is the best way to verify a person while minimizing the crash on confidentiality. This keystroke behavior is used to confirm the identity of a person. Not only single physiological biometrics such as the iris and finger bring us distinctive biometrics. Physiological biometrics defines biological aspects of a person

M.Kanimozhi is with Department of Computer Science and Engineering, Kongunadu College of Engineering & Technology, Trichy, Tamilnadu, India.
A.Kanimozhi is with Department of Computer Science and Engineering, Kongunadu College of Engineering & Technology, Trichy, Tamilnadu, India.

that verify identity. Behavioral biometrics confirms users based on how they conduct a given activity. **Behavioral biometrics** such as the way we sign our name or type in our password are unique as well and have much lower impact on confidentiality and expenses. The way and the manner in which we type on our computer keyboard varies from individual to individual and is considered to be a unique behavioral biometric. Keystroke Recognition is probably one of the easiest biometrics forms to execute and control. This is so because at the current time, Keystroke Recognition is totally a software based explanation. There is no requiring installing any new hardware and software. All that is required is the existing computer and **keyboard** that is already in place and use.

## II. Measuring Factors:

### 2.1 Latency Measures:

In the research field of keystroke dynamics some measurement criteria has to be followed. The modes the key pressed and released are considered as latency of the keystroke data. Some kinds of key pressing modes are press-to-press (PP), press-to-release(PR), release-to-press(RP), release-to-release(RR) [1]. The literature explains the representation of the latency of keystroke by digraph which is the difference between the two presses. In other words it is the time difference between the first key pressed and the second key pressed. The other type of representation is the trigraph. In other words it is the press and the release of the two consecutive keys or it can be called as the time between the press of first key and the release of the second key. [2]
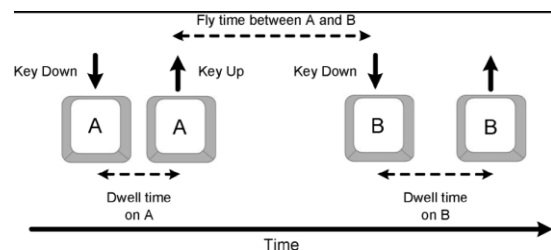


Fig.2.1 Keystroke features and measurements

### 2.2 Timing Information:

The other timing information of the keystroke dynamics are, dwell time or hold time that defines the pressure of a key pressed. It gives the amount of time a particular key is being

-----------------------------------------------------------------------------------------------------------------------------------------

pressed. The kind of information is the flight time that is the pressure of key when it is released. It takes the note of RP latency. The researchers extend the feature up to n-graphs for authentication purpose.

### 2.3 Error rate measurement:

A decision rule, which depends on a static threshold value, decides whether to accept or reject the user into the system. During such matching, some errors may occur. There are two types of such errors namely False Match Rate (FMR), in which the imposters are wrongly accepted to the system. When a system accepts two or more different users as the right person. The other error rate is the False Non-Match Rate (FNMR), in which the authorized or right person is rejected by the system. This happens when the data above the right person is taken from more than one application is not correlated. The system mistakes that the sample doesn't belong to that right person. In a biometric system, there are possibilities for other kind of errors such as Failure to Enroll Error (FER), which arises when captured sample is not properly enrolled into the system. The next is the failure to Capture Rate (FCR), which occurs when the system tales accidentally pressed key during initial collection of sample. [3]

### 2.4 Graphical representation:

ROC (Receiver Operating Characteristics) or DET( Decision Error Tradeoff). These curves can be used to show the performance at the level of threshold. The curve plots true positives (TP), that is (1-FNMR) and false negatives (FN), which is FNMR. To provide good authentication, low FMR is required to reject the imposters at the maximum.

FNR= number of accepted imposters attempts/total no. of imposter attempts

FNMR= no of rejected legitimate users/ total no of legitimate users

The graph represented by the ROC curve gives the clear picture about right attenuation to provide security to the user.
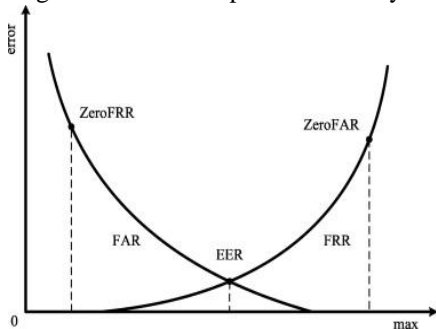


Fig. 2.2 ROC curve for FAR, FRR and EER

### III. PHASES

### 3.1 Data Acquisition:

The main approaches in implementing keystroke dynamics are to collect all the necessary data for evaluation. The data are collected from the number of users in their routine working environment which is termed as data from "uncontrolled environment: that relates to dynamic collection of data. From the collected data such as a) latency time, b) dwell time, c) up-to-up time, the data set is places on the research bed.

The two steps in data acquisition are i) data collection, ii) data analysis [4] based on data collected, data analysis is done. In the first stage, a template has to be creates for the users to work. The template must be of application specific. The interesting features are extracted from each application. These features have to be stored in the database. In the second stage, a new version is created. The user keyboard using rhythm is filtered to collect the features. These features are matched with the database that contains the extracted features done before. A matching is performed to analyze whether the matching is performed to analyze whether the matching is succeeded. During then, we do calculations with FMR and FNMR against certain threshold value. Since we are choosing the interesting features against several applications, the process may be of trial and error method of calculating the timing information over the key pressed and key released.

Analysis of data is done in two stages namely, i) statistical analysis, ii) classification of data. The two consecutive steps are the data mining process done to analyze the data and classify them. It is enough if the statistical accuracy is obtained, but in order to generate machine learning approach and more accurate result algorithms such as neural network and perceptron is preferred.
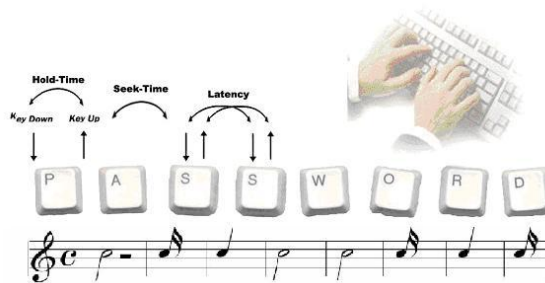


Fig.2.3 Measurements based on key usage

### 3.2 Template constructions:

With the results of the data analysis a template is constructed for each user working on different applications. From the output received we need to generate a template. Before creating a template, the data from the data analysis has to be 'preprocessed'. Data preprocessing is a data mining technique that is used to get better accurate results. The template is developed to finalize that this template is similar to that stored in the database. Concentration on the template designing of the participants is the critical data. Care has to be taken to check, that nothing goes wrong in template construction for a user.

The template creation has been done in two types:
1. General template
2. Personalized template

In general template, the number of features is going to be same for all the user, but not the entries. In personalized data, there would be smaller variations in features from one user to another. The user may have varied keystroke style for every application. Hence personalized template is required

---

additional to generate template. This personalized template is considered as unique template.

### 3.3 Verification:

### 3.3.1 Realizing captured data:

Several attributes are collected regarding the keystroke timing information. The data related to key pressed, key released, pressure of a particular key, time difference among two consecutively pressed( and released) keys. With these data FAR, FRR, EER is determined by working with these data sets in neural network authenticators.

Some of the attributes for realizing the captured data are digraphs, trigraphs, totals username time, total password time total entry time, speed, scan code, edit distance. All related time information data about the keystroke are recorded in nanosecond with 1ns accuracy rate.

### 3.3.2 Comparison with biometrics:

With the available attributes worked on the neural network classification an output is received. This obtained result is matched with the statically stored dynamic keystroke template. If the matches confirms then the right person is continued to work on the system over an application, on the access is timed-out, denied or restricted from the user to use further. This is for what the entire project is concentrating on. Of this verification succeeds,- then the dynamic authentication of users over various applications in a computer system is achieved.

### IV. RELATED WORK

### 4.1 Timing vector based user verification:

When a user types a password on the keyboard, the typing dynamics or timing pattern are measured. Timing vectors is the duration of keystroke time interleaved with keystroke interval time. For a password with 'n' number of characters it has 'n' number of keystroke duration time and 'n-1' keystroke interval time. The sum of these two gives the (n+ (n-1))-dimensional timing vector. The time unit is calculated in millisecond when the next key is pressed before the release of previous key, then the negative time interval is recorded. It is based on the belief that every individual has characteristics and distinctive typing dynamics. A pattern classifier is built to distinguish and identify the right user. to provide a good protection security to the system, the combination of simple password scheme along with pattern classifier is used in spite of negligible increase in cost and processing time.

### (4.1.1) Result:

A password of 7 characters, results in timing vector of 15-dim, since a strike of enter key is also pressed. Example timing vector is 120,60,120,90,120,60,150,-60,120,-30,120,-60,120,120,90,60,150. Where each element was measures in ms. Total of 25 subjects were asked to enter with a new password. For longer password, more number of input and output layers is required in neural network. A 75 vector set of timing information are taken for training set and remaining vectors was taken to train the system. With 10% as threshold

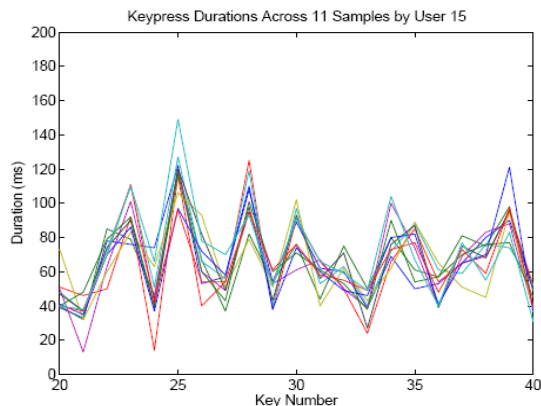acceptance rate, the result was obtained. 4% error rate and 1% of average error rate.


Fig.4.1 keystroke duration across 11 samples user

### 4.2 Identity authorization based on keystroke latencies:

The training set above the keystroke information is collected at different times. Users are allowed to participate under unsupervised conditions. The reference profile collected were representation of n- dimensional feature vector. The data sets were separated into learning and testing sets. These datasets were fed into different classifier techniques such as Non-weighed Probability measures, and Weighed Probability measures,Euclidean distance.

### 4.2.1 Non-weighed probability:

Along with the n- dimensional pattern vector R & U, the additional quadrupts components such as mean, SD, no.of.occurances and data value of $i^{th}$are considered. The score between the reference profiles is calculated by

$$Score(R,U) = \sum_{i=1}^{N} \mathcal{S}_{u_i}$$

where

$$\mathcal{S}_{u_i} = \frac{1}{o_{u_i}} \Big[ \sum_{j=1}^{o_{u_i}} Prob\Big( \frac{X_{ij}^{(u)} - \mu_{r_i}}{\sigma_{r_i}} \Big) \Big]$$

and $X_{ij}^{(u)}$ is the $j^{th}$ occurrence of the $i^{th}$ feature of $U$.

### 4.2.2 Weighed probability measures:

The larger sample set with high frequency in written language are measured, example er, th, sm, et. The score between R&U is calculated as

$$Score(R,U) = \sum_{i=1}^{N} \Big( \mathcal{S}_{u_i} * w_{u_i} \Big)$$

### (4.2.2.1) Result:

The dataset was collected from 63 users. The correct identification rate was 87.18%. The performance of Euclidean distance is 83.22%. The non-weight scoring approach was 85.63%. When examined using Bayesian classifier, it was approximated up to 92.14%, which was almost 5% over the weighed classifier. [6]
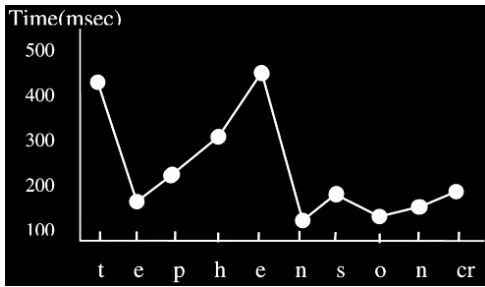
------------------------------------------------------------------------------------------------------------------------


Fig 4.2 Result graph for word Stephenson

### 4.2.3 Euclidean Distance measure:

Similarity can be calculated on pattern vectors using Euclidean distance. Let R=[r1,r2,r3…,rn] and U=[u1,u2,u3,…,un] now the Euclidean distance between the two n- dimensional vector U and R is given by

$$D(R,U) = \left[ \sum_{i=1}^{N} (r_i - u_i)^2 \right]^{1/2}$$

For unknown U, pairwise Euclidean distance is calculated.

### 4.3 Using Ant Colony Optimization for feature subset selection:

It is essential to select the optimized feature from the obtained sample set. Tough there are lot of feature subset selection such as genetic algorithm, artificial intelligence, pattern recognition, neural network, nearest neighbor algorithm, greedy attribute selection, hill climbing algorithm. One such optimization technique used to select the feature sunset is the Ant Colony Optimization. The various steps followed in ACO are,

Step 1: Get the feature values from duration, latency, digraphs of keystroke.

Step2: Calculate the fitness function

Step 3: initialize the no.ofiterations,no.of ants, initial pheromone values and rate of pheromone evaporation.

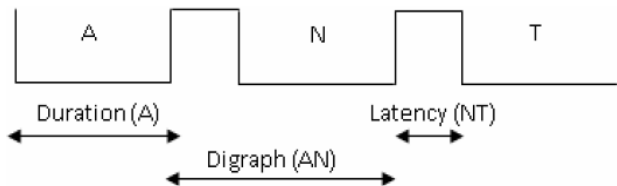Step 4: calculate the local and global optimization value.


Fig.4.3 Duration,Digraph and Latency

### Result:

The fitness value for the calculated duration is 0.425, local minimum for duration is 0.41689, local pheromone update for duration is 0.001. the global duration was calculated as 0.4168 and the global pheromone updating was 0.00225. the remaining ant pheromone update for duration was 0.0001. ant colony optimization can be verified by comparing with BPNN algorithm. The classification error was 0.059% and accuracy was nearly 92.2%. [9]

| nput | I_i | W_ih | I_h | Output of hidden (Hidden) | W_ho | I_o | Sigmoid (Output) o_0 | Target | Difference (Error rate) | Adjusted Weight (W_ih) | Adjusted Weight (W_ih) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Duration** | | | | | | | | | | | |
| Mean | 0.3891 | -0.7 | -0.2723 | 0.46877 | 0.6 | 0.281262 | 0.50073 | 0.1 | 0.160586 | 0.56118 | -0.703482 |
| | | 0.4 | 0.15564 | | | | | | | | 0.402675 |
| SD | 0.7417 | 0.6 | 0.44502 | 0.70889 | -0.5 | -0.354445 | | | | -0.5388 | 0.593362 |
| | | 0.6 | 0.44502 | | | | | | | | 0.605099 |
| **Latency** | | | | | | | | | | | |
| Mean | 0.4169 | -0.7 | -0.29183 | 0.468773 | 0.6 | 0.2812638 | 0.499701 | 0.1 | 0.159761 | 0.561112 | 0.703727 |
| | | 0.4 | 0.16676 | | | | | | | | 0.402848 |
| SD | 0.7437 | 0.6 | 0.44622 | 0.709393 | -0.5 | -0.3546965 | | | | -0.538888 | 0.593351 |
| | | 0.6 | 0.44622 | | | | | | | | 0.605081 |
| **Digraph** | | | | | | | | | | | |
| Mean | 0.4243 | -0.7 | -0.2970 | 0.468222 | 0.6 | 0.2809332 | 0.499448 | 0.1 | 0.159558 | 0.561059 | -0.703791 |
| | | 0.4 | 0.16972 | | | | | | | | 0.402892 |
| SD | 0.7483 | 0.6 | 0.44898 | 0.710530 | -0.5 | -0.355265 | | | | -0.538941 | 0.593313 |
| | | 0.6 | 0.44898 | | | | | | | | 0.605101 |

Table 4.2. Results for duration, digraph and latency

### 4.4 Trigraph features used for identification of keystroke dynamics:

The three conseqitive keys typed are called trigraphs. Trigraph duration is the time between the 1[st] key pressed and the 3[rd] key released. For example if the user types indie, then the sequence of trigraphs and duration (msec).

S1: Ind: 277, ndi: 255, dia: 297, ia + enter key: 326. Now the feature vector is sorted in ascending order. The vector is measured for various vectors S2, S3, etc. with varied timing information. Then the distance is calculated and the value is normalized.

### (4.4.1) Results:

The genuine users and 110 imposters are made to enter the text. 5 samples of genuine users are taken. The experiment was made up to 350 different trigraphs. Let the template for user A is [A1, A2, A3]. The interclass variability of user A is determined using the vectors. The distance is found between the two vectors and the normalized value is obtained, which is between 0 and 1.

| Raw Sample Log | | |
|---|---|---|
| Down | Up | Time |
| A | | 0 |
| | A | 6386 |
| P | | 14824 |
| | P | 11512 |
| P | | 5752 |
| L | | 6594 |
| | P | 4921 |
| | L | 9056 |
| E | | 4943 |
| | E | 5752 |
| S | | 5761 |
| | S | 7393 |

| Dwell Subsample | |
|---|---|
| Graph | Time |
| A | 6386 |
| E | 5752 |
| L | 13977 |
| P | 11512,11515 |
| S | 7393 |

| Trigraph Subsample | |
|---|---|
| Graph | Time |
| A+P+P | 38474 |
| L+E+S | 30433 |
| P+L+E | 25514 |
| P+P+L | 23858 |

| Flight Subsample | |
|---|---|
| Graph | Time |
| A-P | 14824 |
| E-S | 5761 |
| L-E | 4943 |
| P-L | -4921 |
| P-P | 5752 |

| Fourgraph Subsample | |
|---|---|
| Graph | Time |
| A+P+P+L | 45068 |
| P+L+E+S | 37027 |
| P+P+L+E | 42778 |

| Digraph Subsample | |
|---|---|
| Graph | Time |
| A+P | 21210 |
| E+S | 11513 |
| L+E | 18920 |
| P+L | 6594 |

Table 4.1.Trigraph features

## V.ALGORITHM THAT WORKS/ APPROACHES

Once the feature extraction process is done, then the templates are created. The users are classified based on the

---------------------------------------------------------------------------------------------------------------------------------------------------

similarity measures. There are some statistical algorithms that are used to classify the users. Sometimes the combinatorial algorithm can also be used.

### 5.1 Statistical algorithm:

To compute mean and SD of the features. Then the values are compared against the threshold. The comparisons can be done by using hypothesis test, T-test, distance measures such as absolute distance, Euclidean distanced, Manhattan distance, etc. since keystroke dynamics is continuous authentication and non-linear in nature. It is not appreciable to use the linear, statistical methods to compute the features. Moreover, training the datasets is not encouraged by the statistical method to great extent. Hence there is a necessity for more appropriate approaches.

### 5.2 Back propagation neural network:

Back propagation neural network has a forward pass and a backward pass. The features extracted are fed as input to the input layers. It propagates to a value to the hidden layer. The values generated by the hidden layer are fed as input to the output layer. The output layer in turn calculates the output value for the given inputs since the weight are random values sometimes the output vector is not related. Hence there requires a backward pass. This is achieved by back propagation. The steps involved here are 1. To compute error in the output layer, 2.To compute error in hidden layer, 3. Adjust the weight values to improve the performance, 4. Sum up the total error [13].

### 5.3 Neural networks:

Neural network is also called Artificial Neural Network (ANN). It is a non-linear statistical data modeling tool. There are basically two different ways of learning the training data sets. They are supervised learning and unsupervised learning. The most popular supervised learning is back propagation [10]. The other supervised learning algorithm examples are train a decision tree, cross validation, neural networks, transduction, ensembles. The popular unsupervised learning are Hopfield Neural Network (HNN). The other unsupervised learning algorithm examples are clustering; dimensionality reduction using PCA, independent component analysis, etc. neural networks is suggested by many researches to give best results. Neural network can handle many parameters. Due to the black box feature of neural network, it is considered as a problem during continuous keystroke authentication [11]

### 5.4 Support vector Machine Algorithm:

Keystroke dynamics concentrates on identifying the correct users. On the other side imposters are also identified. Support vector machine (SVM) is one such algorithm used to detect the imposters. It is considered as consistent and low complexity algorithm. The approach is carried out in two ways, 1. One class svm (OC-SVM), 2. Two class svm (TC-SVM). Ocscm is used to capture data with probable values. Tcsvm provides training data with overall coverage of objects. (i.e. imposters). The FAR and FRR calculated using two

approaches are compared and best performance is evaluated. [12].

### 5.5 Genetic algorithm:

Genetic algorithm is a class of probability optimization algorithm, inspired by the biological evaluation process. It uses the concept of natural selection and genetic inheritance (Darwin 1859). It was originally developed by John Holland (1975). The feature extractions are considered as population. There are several steps involved sequentially after the population is selected.



$$c_1(\phi) = \min_{z \in 1}(w \cdot z + b)$$

$$c_2(\phi) = \max_{z \in -1}(w \cdot z + b)$$

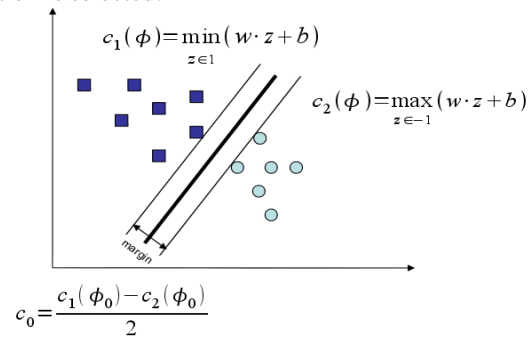$$c_0 = \frac{c_1(\phi_0) - c_2(\phi_0)}{2}$$

Fig.5.1 SVM vector used for classification of users

1. The populations are ranked according to their fitness.
2. The population is made to reproduce by two steps such as crossover and mutation. It is based on the concept of 'a pair of parents produces two children'.
3. The steps are repeated until the desired fitness level is reached

### 5.5.1 PSEUDOCODE:

*Simple genetic algorithm*
*Produce initial population of intervals*
*Evaluate the fitness of all individuals*
*While (termination not met)*
*{*
*Do*
*{*
*Select fittest individual for reproduction;*
*Recombine (i.e. crossover individuals;*
*Mutation individuals;*
*Evaluated the fittest modified individual;*

*Generate new population;*
*}*
*}*
*End while* [14]

Genetic algorithm can be preceded in travelling salesman problem (TSP). The travelling salesman must visit every city at least once and return to the starting point at the minimum total cost of the entire travel. The TSP can be approximately a genetic algorithm. The advantage of GA is
1. It can solve every optimization problem
2. It is easy to understand and simulate
The disadvantage is if the fitness value is poorly functioned, then the optimization will be at risk.
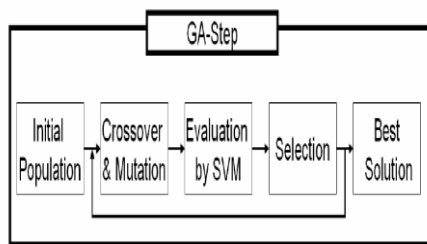
--------------------------------------------------------------------------------------------------------------------------------



Fig.5.2 Crossover and mutation process

### 5.6 Ant colony optimization used to solve TSP:

Ant colony optimization was introduced by Marco Dorigo in Italy in the year 1992 in his doctoral thesis. It is used to solve TSP ants go through the food by laying down the pheromone traits. The shortest path is found via pheromone traits.

1. The ant move in random
2. After some time, ants follow the traits which have more amount of pheromone.
3. Meanwhile, all the ants will follow the pheromone traits
4. The previous path is evaporated.

Each ant located in city I has to move to city j. d (I,j) is the attractiveness, which is the function that gives the inverse of cost. T(I,j) is the trait level, detecting the amount of pheromone trait. The set of cities not visited by the ant k in city I is $T_k(i)$. the probability that ant k $P_k(I,j)$ in city i will go to city j, is calculated.

### 5.6.1 PSEUDOCODE:for general ant colony:

*Initialize the base attractiveness $\tau$ for each edge*
*For (each ant) do*
*P( choose the edge)*
*Add and move to table list of each ant*
*Repeat until each ant complete solution*
*End;*
*For (each ant that completes a soln)*
*Update $\tau$ for each edge the ant traversed*
*End;*
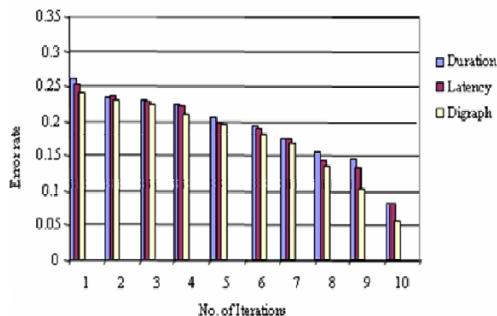*If (local better than global) save local*
*End;*
*End;*
[15]



Fig.5.3 Results based on ACO and BPNN

The benefit of ACO is that it can solve NP-Hard problem in short time. It balances the previous solution and new exploring solution. Optimal solution is obtained.The limitations of ACO are the coding differs for different applications. Ineffective utilization of previous solution affects global solution.

### CONCLUSION

The upcoming of biometric technologies is promising. Biometric device and application still grow worldwide. There are many factors that may move forward the growth of biometric technologies. An important matter of the growth of biometrics has been the price to apply them. Moreover, improved accuracy rates can play a big division within the receiving of biometric technologies.

The evolution and research into biometric error testing false refuse and false stay for has been of enthusiastic attention to biometric developers. In this paper represent a free text analysis with monograph and diagraph analysis and using the neural network proposed approaches that help to solve the security challenges in the existing system.

### REFERENCES

[1]    F. Bergadano, D. Gunetti, and C. Picardi, "User authentication through keystroke dynamics," ACM Trans. Inform. Syst. Security, vol. 5, no. 4, pp. 367–397, Nov. 2002.
[2]    M. Brown and S. J. Rogers, "User identification via keystroke characteristics of typed names using neural networks," Int. J. Man-Mach. Stud., vol. 39, no. 6, pp. 999–1014, Dec. 1993.
[3]    P. Dowland, S. Furnell, and M. Papadaki, "Keystroke analysis as a method of advanced user authentication and response," in Proc. IFIP TC11 17th Int. Conf. Inform. Security: Visions Persp., May 7–9, 2002, pp. 215–226.
[4]    D. Gunetti and C. Picardi, "Keystroke analysis of free text," ACM Trans. Inform. Syst. Security, vol. 8, no. 3, pp. 312–347, Aug. 2005.
[5]    F. Monrose and A. Rubin, "Authentication via keystroke dynamics," in Proc. Fourth ACM Conf. Comput. Commun. Security, pp. 48–56, Apr. 1997.
[6]    D. Polemi. "Biometric techniques: Review and evaluation of biometric techniques for identification and authentication, including an appraisal of the areas where they are most applicable," [Online].
[7]    S. Ross, "Peirce's criterion for the elimination of suspect experimental data," J. Eng. Technol., vol. 20, no. 2, pp. 38–41, Oct. 2003.
[8]    M. Villani, C. Tappert, N. Giang, J. Simone, St. H. Fort, and S.-H. Cha, "Keystroke biometric recognition studies on long-text input under ideal and application-oriented conditions," in Proc. 2006 Conf. Comput. Vis. Pattern Recognit. Workshop (CVPRW'06), June 2006, p. 39.
[9]    L. Ballard, D. Lopresti, and F. Monrose, "Forgery quality and its implication for behavioral biometrics security," IEEE Trans. Syst. Man Cybernet., Part B, vol. 37, no. 5, pp. 1107–1118, Oct. 2007.
[10] M. S. Obaidat and B. Sadoun, "Verification of computer users using keystroke dynamics," IEEE Trans. Syst., Man Cybernet., Part B, vol. 27, no. 2, pp. 261–269, 1997.