

**NISTIR 7273**

**Using Chebyshev's Inequality to  
Determine Sample Size in Biometric  
Evaluation of Fingerprint Data**

Jin Chu Wu  
Charles L. Wilson

**NIST**

**National Institute of Standards and Technology**  
Technology Administration, U.S. Department of Commerce

NISTIR 7273

# Using Chebyshev's Inequality to Determine Sample Size in Biometric Evaluation of Fingerprint Data

Jin Chu Wu

Charles L. Wilson

*Image Group, Information Access Division  
Information Technology Laboratory*

November 2005



**U.S. DEPARTMENT OF COMMERCE**

*Carlos M. Gutierrez, Secretary*

**TECHNOLOGY ADMINISTRATION**

*Michelle O'Neill, Acting Under Secretary of Commerce for Technology*

**NATIONAL INSTITUTE OF STANDARDS AND TECHNOLOGY**

*William Jeffrey, Director*

# *Using Chebyshev's Inequality to Determine Sample Size* **in Biometric Evaluation of Fingerprint Data**

Jin Chu Wu\* and Charles L. Wilson

Image Group, Information Access Division, Information Technology Laboratory

National Institute of Standards and Technology, Gaithersburg, MD 20899

## **Abstract**

The fingerprint datasets in some cases may exceed a million of samples. The underlying distribution functions of the similarity scores are unknown. Therefore, the needed size of a biometric evaluation test set is an important question in terms of both efficiency and accuracy. In this article, Chebyshev's inequality, in combination with simple random sampling, is used to determine the sample size for biometric applications. The performance of fingerprint-image matcher is measured by both the area under a Receiver Operating Characteristic (ROC) curve and the True Accept Rate (TAR) at an operational False Accept Rate (FAR). The Chebyshev's greater-than-95% intervals of these two criteria based on 500 Monte Carlo iterations are computed for different sample sizes as well as for both high- and low-quality fingerprint-image matchers. The stability of such Monte Carlo calculations with respect to the number of iterations is also presented. The choice of sample size is dependent on the qualities of fingerprint-image matchers as well as on which performance criterion is invoked. However, in general, for 6000 match similarity scores, 50000 to 70000 scores randomly selected from 35994000 non-match similarity scores can ensure reasonable accuracy with greater-than-95% probability.

*Keywords:* Fingerprint Matching; Chebyshev's Inequality; Sample Size; Simple Random Sampling; Receiver Operating Characteristic (ROC) Curve; Biometrics

### **1. Introduction**

The fingerprint datasets in many cases may exceed a million of samples. Therefore, the size of biometric evaluation test sample is an important question in terms of both efficiency and accuracy. Since two years ago, the National Institute of Standards and Technology (NIST) has used large samples of fingerprint data from a wide range of government sources to evaluate the fingerprint-image matchers from different vendors\* [1,2]. In the SDK tests [2], 6000 subjects' fingerprint images

---

\* Corresponding author. Tel: + 301-975-6996; fax: + 301-975-5287. E-mail address: [jinchu.wu@nist.gov](mailto:jinchu.wu@nist.gov) (J.C. Wu).

\* These tests were performed for the Department of Homeland Security in accordance with section 303 of the Border Security Act, codified at 8 U.S.C. 1732. Specific hardware and software products identified in this report were used in order to adequately

are used as a probe, and 6000 second fingerprint images of the same subjects are used as a gallery. The probe is matched against the gallery. This creates 6000 match similarity scores from the same subjects' different fingerprint-image comparisons, and 35994000 non-match similarity scores from different subjects' fingerprint-image comparisons.

For such fingerprint data, there is usually no underlying parametric distribution function for match and non-match similarity scores, respectively. Thus, the nonparametric approach must be employed to analyze the data and evaluate matchers [3]. Nonetheless the fingerprint-image matcher is designed in such a way that the higher (lower) values of similarity scores tend to indicate that two fingerprint images are more (less) similar. Hence, the distribution of the match similarity scores is always centered at higher scores than the distribution of the non-match similarity scores. The True Accept Rate (TAR) is defined as the cumulative probability of the match similarity scores at a specified similarity score (i.e., threshold) from the highest match similarity score. And the False Accept Rate (FAR) is determined as the cumulative probability of the non-match similarity scores at a threshold from the highest non-match similarity score [3].

The fingerprint-image matcher can be evaluated by a Receiver Operating Characteristic (ROC) curve. An ROC curve is constructed based on TAR and FAR by moving the threshold, one similarity score at a time, from the highest similarity score to the lowest similarity score. Thus, any ROC curve has two fixed endpoints, i.e., starting from (0, 0) and ending at (1, 1), in the FAR-and-TAR coordinate system. Usually, it is above the straight line that connects these two points [3]. An ROC curve can be measured by using either the area under the ROC curve [3] or the TAR value at an operational FAR value [1,2]. In this article, both of these two criteria will be employed. The size of our fingerprint datasets is very large in comparison to the applications of ROC curves in other areas [4,5,6,7, and references therein]. However, the principles remain the same.

How much fingerprint data should be selected from a large dataset to obtain both efficiency and accuracy in biometric evaluation? Different sizes of samples generate different ROC curves. Hence, the sample size can be determined by the accepted deviations of ROC curves for samples with reduced sizes from the ROC curve in the baseline, i.e.,  $\Delta(ROC\ curve)$ , in terms of both or either of the above two criteria, at a specified probability (e.g., 95%). The baseline can be generated from the largest dataset that the available computer power can handle from the largest consolidated dataset.

In the SDK tests [2], all performances of fingerprint-image matchers were evaluated based upon comparing the distribution of 6000 match similarity scores with the distribution of 35994000 non-match similarity scores. If the current SDK evaluation is set to be a baseline, then with respect to 6000 match similarity scores, out of 35994000 non-match similarity scores, how many non-match similarity scores are needed to achieve the same performance? In other words, the issue of determining the sample size for SDK turns out to be: 1) reduce the number of non-match similarity scores, 2) take *one* trial, 3) the result must be close to the baseline result within an accepted tolerance at a specified probability. In this article, for simplicity, we restricted ourselves to the scenario in which the number of match similarity scores is fixed as 6000, but the number of non-match similarity scores can be varied. As a matter of fact, the same methodology can be applied to other scenarios.

---

support the development of technology to conduct the performance evaluations described in this document. In no case does such identification imply recommendation or endorsement by the National Institute of Standards and Technology, nor does it imply that the products and equipment identified are necessarily the best available for the purpose.

As specified above, one of our requirements is that the test be performed only once. To satisfy this objective, Chebyshev's inequality is invoked. Using Chebyshev's inequality, an interval in which a percentage of population resides can be determined, provided that the lower bound on the probability is specified. If an interval can contain the baseline result as well as, for example, greater than 95% of the test results in a specified circumstance, then the one-trial test result will have greater-than-95% probability to fall in that interval and its deviation from the baseline result will not exceed the length of the interval. In other words, only matters the absolute error, not the relative error. This is consistent with the above statement of  $\Delta$  (*ROC curve*).

Further, as stated above, the number of match similarity scores is fixed as 6000. Thus, to ensure that the ROC curves from the test results are close to the ROC curve in the baseline in terms of the above criteria, the distributions of non-match similarity scores with reduced sizes must be "very similar" to the distribution of 35994000 non-match similarity scores in the baseline. To serve this purpose, the simple random sampling without replacement is applied. A simple random sample selected from 35994000 non-match similarity scores constitutes a new set of non-match similarity scores, and its distribution is used with the distribution of 6000 match similarity scores to generate an ROC curve.

A Chebyshev's greater-than-95% interval can be obtained using a Monte Carlo calculation. Different sizes of simple random samples are selected from 35994000 non-match similarity scores in the baseline. 500 Monte Carlo iterations are carried out for different sample sizes. Thereafter, the sample size of non-match similarity scores can be determined according to whether the Chebyshev's interval is within an accepted tolerance. In addition, the stability of the Monte Carlo calculation with respect to the number of iterations is also dealt with in this article. It is quantified by the worst deviation of the test result from the baseline result within the Chebyshev's interval.

The methods, i.e., the Chebyshev's inequality and Chebyshev's greater-than-95% interval, the simple random sampling, and the stability metric, are presented in Section 2. The results of their applications to determining sample sizes of non-match similarity scores in biometric evaluation of high-quality and low-quality fingerprint-image matchers are provided in Section 3. Discussion of the sampling error of the sample mean, the scope of the application of this methodology, the matcher-quality dependence of the results, and other issues can be found in Section 4. Finally, the conclusion is stated in Section 5.

## **2. Methods**

Chebyshev's inequality, in combination with simple random sampling, is used to determine the sample size for biometric applications. The stability of the calculation with respect to the number of Monte Carlo iterations will be addressed as well.

### **2.1 Chebyshev's Inequality [8] and Chebyshev's Greater-Than-95% Interval**

If  $\xi$  is a random variable and its mean and variance exist, i.e.,  $M(\xi) = \mu < \infty$  and  $V(\xi) = \sigma^2 < \infty$ , then Chebyshev's inequality

$$P\{|\xi - \mu| \geq k\sigma\} \leq \frac{1}{k^2} \quad (1)$$

is valid for any  $k > 1$ . A variation of Chebyshev's inequality can be expressed as

$$P\{|\xi - \mu| < k\sigma\} > 1 - \frac{1}{k^2} \quad (2)$$

It states that greater than  $(1 - 1/k^2)$  of population falls within  $k$  ( $k > 1$ ) standard deviations, i.e.,  $k\sigma$ , from the population mean  $\mu$ .

The proof of Chebyshev's inequality is trivial. However, its concept is profound. First of all, it is important that Chebyshev's inequality holds good without any assumption regarding the shape of the distribution of population as long as the mean and variance exist. This nonparametric characteristic is just the one that was encountered and dealt with for fingerprint data distributed with respect to similarity scores generated by fingerprint-image matchers [3]. Second, Chebyshev's inequality provides a way to compute a quantitative relationship between an interval, which is greater than one standard deviation from the population mean, and the lower bound on the probability at which observation values of a random variable fall into that interval. In fact, there are many other implications of Chebyshev's inequality, which are out of the scope of this article.

On the other hand, Chebyshev's inequality cannot offer the lower bound of the proportion of the population that lies within one standard deviation or less from the population mean. Furthermore, for distributions that have a special shape, such as normal distribution, etc., the probability at which the population falls into an interval that is greater than one standard deviation from the population mean is much larger than the lower bound on the probability calculated using Chebyshev's inequality for the same size of interval. In other words, for instance, for the normal distribution, 95% of the population is within  $1.96\sigma$  from the population mean. However, if the lower bound on the probability for any type of distribution is also set to be 95%, then the required interval provided by Chebyshev's inequality is  $4.48\sigma$ , that is about 2.29 times larger.

Such an interval that is  $4.48\sigma$  from the population mean is defined as Chebyshev's greater-than-95% interval in this article. Chebyshev's greater-than-95% interval is different from a 95% confidence interval for an estimate of the population mean, which is a consequence of the Central Limit Theorem. Chebyshev's interval only describes the fact that the probability at which the population falls into an interval that is within  $4.48\sigma$  from the population mean is greater than 95%, in spite of the shape of the population distribution. Therefore, no inferential statistics, such as hypothesis tests, etc., can be carried out based on Chebyshev's interval.

Chebyshev's greater-than-95% interval serves the objective that the result of taking *one* trial must be close to the baseline result within a desired tolerance at a specified probability. Since greater than 95% of population lies in Chebyshev's interval, assuming the interval contains the baseline result, the probability at which the one-trial test result falls in that interval and satisfies the requirement is greater than 95%. In addition, 95% for Chebyshev's interval is the lower bound on the probability. Thus, it provides conservative estimates in its applications.

In Chebyshev's greater-than-95% interval, the population mean  $\mu$  and the population standard deviation  $\sigma$  are used.

The sample mean  $\hat{\mu} = \sum_{i=1}^n x_i / n$  and the sample standard deviation  $\hat{\sigma} = \sqrt{\sum_{i=1}^n (x_i - \hat{\mu})^2 / (n - 1)}$ , where  $x_i$ 's are independent observation values of a random variable  $\xi$  and  $n$  is the number of observations, are unbiased point estimators of  $\mu$  and  $\sigma$ , respectively. However, according to the Law of Large Numbers,  $\hat{\mu}$  and  $\hat{\sigma}$  converge to  $\mu$  and  $\sigma$ , respectively, as the number of observations increases [9]. While comparing the sample mean with the baseline result, the sampling error of the sample mean, i.e., the absolute value of the difference between  $\hat{\mu}$  and  $\mu$ , might not be needed to be taken into account in our applications. This will be discussed in Section 4.

$4.48 \hat{\sigma}$  is not a small quantity in many applications, and it could happen that  $4.48 \hat{\sigma}$  went beyond the allowed range of random variables. However, in our applications, thanks to the simple random sampling and the large size of fingerprint data, the sample standard deviation is very small (see Section 3.2). Therefore,  $4.48 \hat{\sigma}$ , i.e., Chebyshev's greater-than-95% interval, can be used as a criterion to determine the sample size in biometric evaluation of fingerprint data.

## 2.2 Simple Random Sampling

Both match and non-match similarity scores will be referred to as similarity scores in this section. In order to test how far the number of similarity scores can be reduced with respect to the baseline, a simple random sample of similarity scores is selected from the finite set of similarity scores in the baseline. The simple random sampling (SRS) applied in this article is carried out under three assumptions: 1) the population is finite, 2) each member in the population has the same probability of being selected, and 3) it is a sampling without replacement (WOR) for members in the population.

The similarity scores can be represented as integers within different ranges, depending on different fingerprint-image matchers [3]. Let the integral score set be  $\{s\} = \{s_{\min}, s_{\min}+1, \dots, s_{\max}\}$ , consecutively from  $s_{\min}$  to  $s_{\max}$ , where  $s_{\min}$  and  $s_{\max}$  are the minimum and maximum similarity scores, respectively. Thus, the similarity score set is a set of integral scores,

$$\mathbf{S} = \{s_i \mid \forall i \in \{1, \dots, N\}\} \quad (3)$$

where  $s_i \in \{s\}$  and  $N$  is the total number of similarity scores. The similarity scores  $s_i$  may not exhaust all members in the integral score set  $\{s\}$ . Moreover, some of the fingerprint-image comparisons may very well share the same integral value. Therefore, the similarity score set  $\mathbf{S}$  can be partitioned into pairwise-disjoint subsets  $\{\mathbf{S}_s\}$ . In each of the subsets,  $\mathbf{S}_s$ , the members have the same integer  $s \in \{s\}$ . The similarity score set  $\mathbf{S}$  is the union of all these subsets  $\{\mathbf{S}_s\}$ .

The frequency  $f(s)$  of the similarity score  $s$ , which appears in the similarity score set  $\mathbf{S}$ , is the size of the subset  $\mathbf{S}_s$ . The corresponding probability  $p(s)$  equals the frequency  $f(s)$  divided by the total number of similarity scores, i.e.,  $p(s) = f(s) / N$ . Thus, in the baseline, the discrete probability distribution function of the similarity scores, by including zero probability caused by some similarity scores that appear in the integral score set  $\{s\}$  but not in the similarity score set  $\mathbf{S}$ , can be represented in terms of the probability  $p(s)$  as

$$\mathbf{P} = \{ p(s) \mid \forall s \in \{s\} \text{ and } \sum_{\tau=s \min}^{s \max} p(\tau) = 1 \} \quad (4)$$

According to the SRS as stated above, each member in the similarity score set  $\mathbf{S}$  with size  $N$  in the baseline has the same probability of being selected. Hence, the probability of being chosen for such a member is  $1/N$ . Furthermore, the SRS is assumed to be WOR for scores in the similarity score set  $\mathbf{S}$ . Thus, for any similarity score  $s$  in the integral score set  $\{s\}$ , whose frequency of appearing in the similarity score set  $\mathbf{S}$  is  $f(s)$ , the probability of being selected is  $f(s) / N$ . As a result, after sufficiently large amount of such selections, the discrete probability distribution function of the selected similarity scores will approach to the discrete probability distribution function of the similarity scores in the baseline as expressed in Equation (4).

The size of the similarity scores is relatively large. Therefore, for a large amount of simple random samples, the variance of area under an ROC curve and even the variance of the TAR value at an operational FAR value, caused by the discrepancy between the distribution of the selected similarity scores and the distribution in the baseline, are quite small. It follows that  $4.48 \hat{\sigma}$ , i.e., Chebyshev's greater-than-95% interval, can be applied as a criterion to determine the sample size and is suitable to serve our objectives. As for using the Kolmogorov-Smirnov Test to see the difference between such two distributions, it depends on how to deal with the ties of similarity scores in these two discrete probability distribution functions (see Section 3.1).

### 2.3 The Stability Metric

Chebyshev's greater-than-95% interval can be obtained using Monte Carlo calculation, i.e., by running a number of Monte Carlo iterations. How many iterations of the Monte Carlo calculation based upon SRS are needed to determine the sample size of similarity scores in the biometric evaluation of fingerprint data for a fingerprint-image matcher? In other words, how stable is the Monte Carlo calculation with respect to the number of iterations? The Monte Carlo stability is related to how much the sample size of similarity scores is, which fingerprint-image matcher is dealt with, and which criterion of evaluation of ROC curve is involved.

The discrete probability distribution functions of the selected similarity scores for the amount of sample sizes discussed in this article do not deviate very much from the discrete probability distribution function in the baseline. Thus, Chebyshev's intervals always contain the baseline result as observed in our tests (see Section 3.2). Thereafter, the maximum of two distances between the baseline result and two end points of Chebyshev's greater-than-95% interval, respectively, can be chosen as a metric to measure the stability of our Monte Carlo calculation.

Such a stability metric can be expressed as

$$M_n = \max[b - (\hat{\mu}_n - 4.48\hat{\sigma}_n), (\hat{\mu}_n + 4.48\hat{\sigma}_n) - b] \quad (5)$$

where  $M_n$  is the stability metric for  $n$  Monte Carlo iterations,  $b$  is the baseline result (either the area under an ROC curve or the TAR value at an operational FAR value in the baseline), and  $\hat{\mu}_n$  and  $\hat{\sigma}_n$  are the unbiased point estimators of the population mean  $\mu$  and the population standard deviation  $\sigma$  after  $n$  iterations, respectively. And the population is determined by the sample size of similarity scores and a chosen fingerprint-image matcher. This stability metric describes the worst

deviation of the one-trial test result from the baseline result inside Chebyshev's greater-than-95% interval. That is, with greater-than-95% probability, a one-trial test result will not deviate from the baseline result by more than the stability metric in a specified circumstance.

Outside Chebyshev's greater-than-95% interval, some points can have deviations from the baseline result less than the stability metric, if the baseline result is not in the middle of Chebyshev's interval. If these points happen to be one-trial test results, the real probability at which the one-trial test result deviates from the baseline result less than the stability metric will be greater than the exact probability at which the population lies inside the Chebyshev's interval, which is subsequently greater than the lower bound on the probability computed using Chebyshev's inequality. Therefore, the stability metric is a more conservative measurement than Chebyshev's interval. Once the variation of the stability metric over the number of iterations is within an accepted tolerance, the stability of the Monte Carlo calculation is achieved.

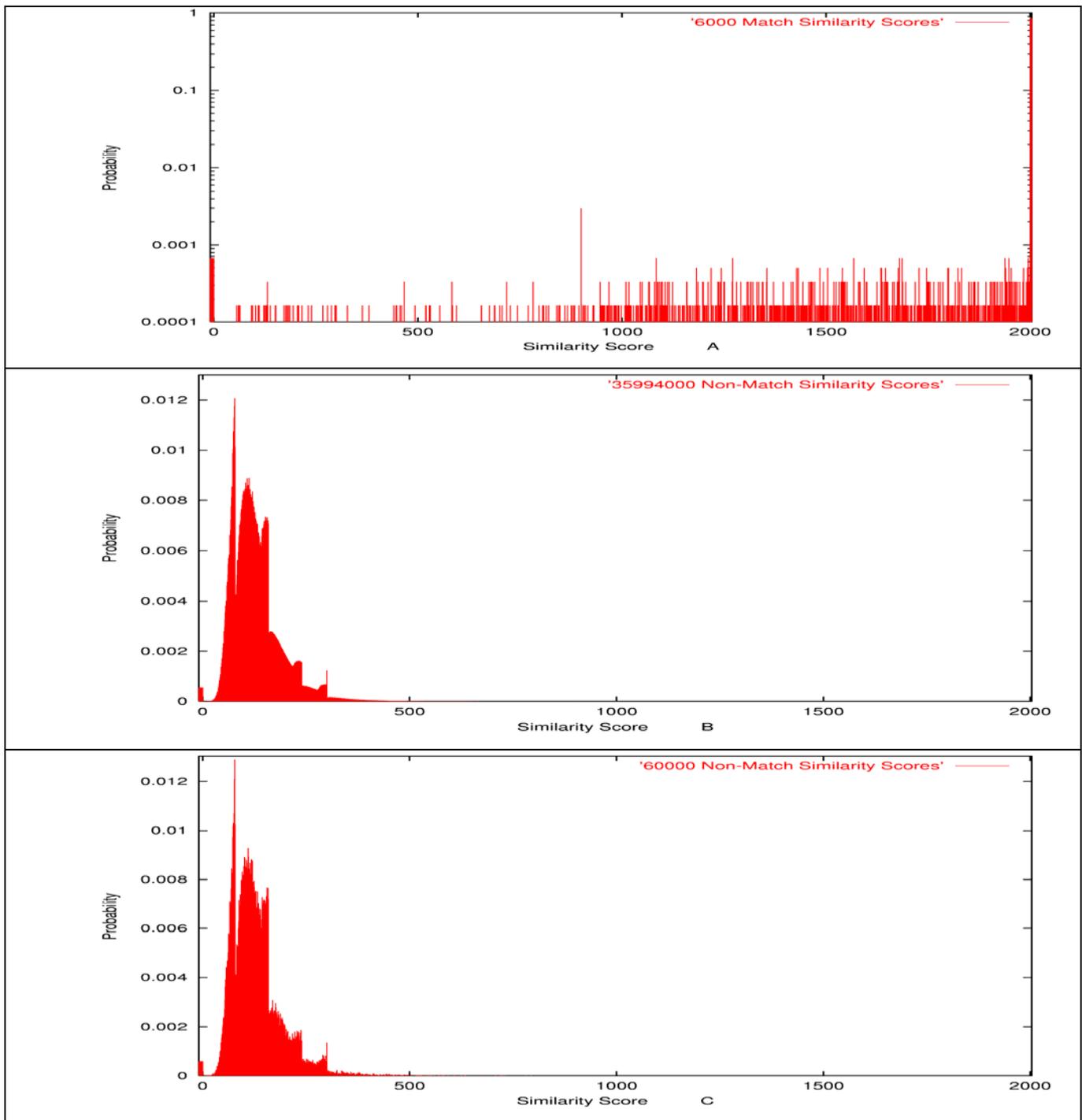
### **3. Results**

Chebyshev's greater-than-95% intervals vary depending upon 1) the quality of the fingerprint-image matcher, 2) the criterion to evaluate ROC curves, 3) the sample size of SRS, and 4) the number of Monte Carlo iterations. Two fingerprint-image matchers were taken as examples, among which Matcher 1 was high-quality matcher and Matcher 2 was low-quality matcher. Both matchers were executed on the same fingerprint dataset. And two criteria were employed, namely, the area under an ROC curve (i.e., AUROC) as well as the TAR value at an operational FAR value that is set to be 0.001 (abbreviated as TVAFV in the following).

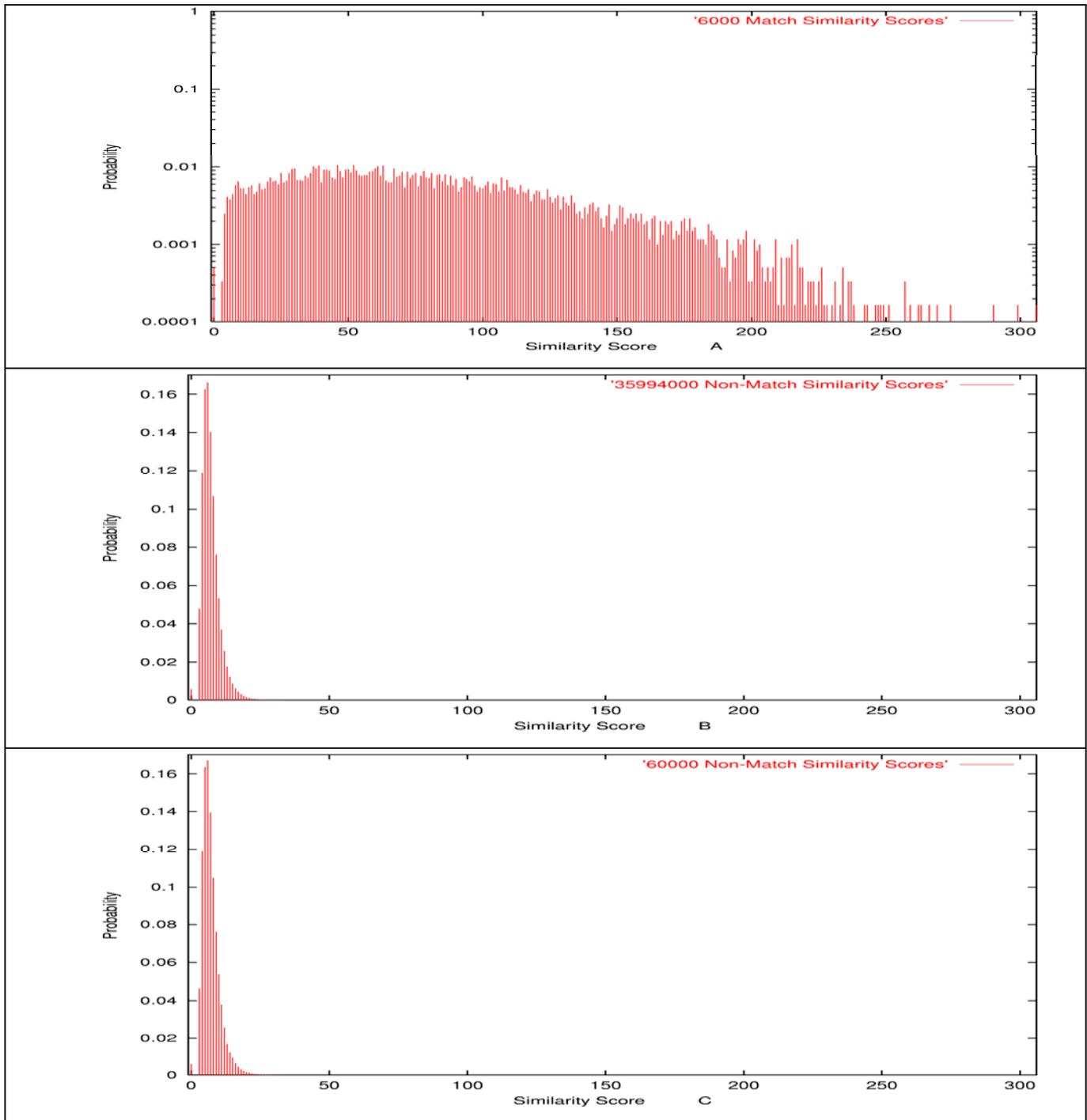
The baseline result of a matcher was obtained using 6000 match similarity scores and 35994000 non-match similarity scores, as performed in the current SDK tests. For the results presented here, the number of match similarity scores was fixed as 6000, however the number of non-match similarity scores was reduced from 35994000. So the issue turns out to be how much low the number of non-match similarity scores can go to keep the one-trial test result as close as to the baseline result within accepted tolerance at a specified probability. In addition, to test the stability, different numbers of Monte Carlo iterations were carried out.

#### **3.1 The Discrete Probability Distribution Functions [3]**

It is always a good thing to take a look at the distribution function first. For Matcher 1, the discrete probability distribution functions of the 6000 match similarity scores, the 35994000 non-match similarity scores, and the 60000 non-match similarity scores that were a simple random sample selected from 35994000 non-match similarity scores, are shown in Figure 1 A, B, and C, respectively. The integral similarity scores of Matcher 1 run from 0 through 2000. And for Matcher 2, the corresponding discrete probability distribution functions are presented in Figure 2 A, B, and C, respectively. Its integral similarity scores run from 0 to 306.



**Figure 1** The discrete probability distribution functions of the 6000 match similarity scores (A), the 35994000 non-match similarity scores (B), and the 60000 non-match similarity scores selected simply randomly from 35994000 non-match similarity scores (C), respectively. All distributions were generated by using the fingerprint-image Matcher 1. The integral similarity scores run from 0 to 2000. The widths of peaks at the score 2000 and zero are enlarged to show the characteristics of the distribution.



**Figure 2** The discrete probability distribution functions of the 6000 match similarity scores (A), the 35994000 non-match similarity scores (B), and the 60000 non-match similarity scores selected simply randomly from 35994000 non-match similarity scores (C), respectively. All distributions were generated by using the fingerprint-image Matcher 2. The integral similarity scores run from 0 to 306.

For the discrete probability distribution function of 6000 match similarity scores, Matcher 1 has a sharp peak at the highest similarity score as depicted in Figure 1A. However, Matcher 2 does not have such kind of peak as shown in Figure 2A. To show the contrast, the probability is depicted in logarithmic scale in these two figures. For the distribution function of

35994000 non-match similarity scores, Matcher 1 and 2 are completely different as presented in Figure 1B and 2B, respectively. All these demonstrate that distributions of similarity scores vary very much from matcher to matcher and in many cases there is no parametric model to fit [3].

To explore the relationship between the distribution of simple random samples of non-match similarity scores and the distribution of 35994000 non-match similarity scores in the baseline, the discrete probability distribution functions of 60000 non-match similarity scores of simple random samples were examined for Matcher 1 and 2, as shown in Figure 1C and 2C, respectively. By using Kolmogorov-Smirnov Test, if a tie of non-match similarity scores in a discrete distribution was treated as separated individual scores that shared the same value while comparing two cumulative distribution functions [10], it was found that the two-tailed p-values of two distribution functions (i.e., 35994000 against 60000 non-match similarity scores) for Matcher 1 and 2, respectively, were much less than 1%. This indicates that these two distributions are likely to be different.

However, if a tie of non-match similarity scores was dealt with as a single bar at the shared value of these scores, then the two-tailed p-values of the Kolmogorov-Smirnov Test were much larger than 5%. In this sense, these two distributions are unlikely to be different. This is why it is hard to see the difference between these two distribution functions visually. Indeed, such a treatment of ties matches the way of formation of ROC curve, which is generated by moving the threshold, one similarity score at a time, from the highest similarity score to the lowest, for two distributions of match and non-match similarity scores [3]. Therefore, the SRS has little impact on ROC curves, even when the sample size of non-match similarity scores is as small as 60000. In other words, the variances of AUROC and TVAFV, caused by using SRS, i.e., by the discrepancy between the distribution of the selected similarity scores and the distribution in the baseline, are so small that Chebyshev's greater-than-95% interval can be invoked.

### 3.2 Chebyshev's Greater-Than-95% Interval

The results of fingerprint-image Matcher 1 and 2 are presented. The quality of Matcher 1 is higher than that of Matcher 2. Two criteria, AUROC and TVAFV, are used to evaluate the qualities of matchers. AUROC has a standard error associated with [3], but TVAFV does not. To be consistent between these two criteria as well as for simplicity, the standard error of AUROC is not used in this article. The values of AUROC and TVAFV of the baseline for Matcher 1 and 2 are shown in Table 1.

Matcher	AUROC	TVAFV
1	0.997170	0.991167
2	0.983862	0.892333

**Table 1 The baseline results of Matcher 1 and 2.**

Generally speaking, the smaller the sample size of 35994000 non-match similarity scores, the wider the Chebyshev's greater-than-95% interval. The error bars, i.e.,  $4.48 \hat{\sigma}$ , are relatively small for the sample size greater than 100000, and relatively large for the sample size less than 10000. Therefore, results are presented with the sample sizes decreasing from

100000 non-match similarity scores down to 10000 by every 10000 for both Matcher 1 and 2. The Monte Carlo calculation was run for 500 iterations in each case. Chebyshev's greater-than-95% interval is expressed in terms of sample mean, the error bar, the upper bound (sample mean plus error bar), and the lower bound (sample mean minus error bar). The results of AUROC and TVAFV for two matchers are shown, respectively, from Table 2 to Table 5. The trend of variations of Chebyshev's greater-than-95% interval in each case, and the relationship between the intervals and the baseline results are accordingly depicted, respectively, from Figure 3 to Figure 6.

100000	70000	60000	50000	40000	30000	20000	10000
0.99717	0.99717 0	0.99717 1	0.99717 1	0.99717 0	0.99717 0	0.99717 1	0.99717 1
0.00002	0.00003 0	0.00003 1	0.00003 3	0.00003 8	0.00004 2	0.00005 2	0.00007 0
0.99719	0.99720 0	0.99720 1	0.99720 4	0.99720 7	0.99721 2	0.99722 3	0.99724 1
0.99714	0.99714 0	0.99714 0	0.99713 7	0.99713 2	0.99712 9	0.99711 9	0.99710 1

**Table 2 Matcher 1's Chebyshev's greater-than-95% intervals in terms of sample mean, error bar, upper bound and lower bound of 500 Monte Carlo iterations for different sample sizes of non-match similarity scores using AUROC.**

100000	70000	60000	50000	40000	30000	20000	10000
0.99117	0.99117 8	0.99117 2	0.99119 3	0.99119 6	0.99120 4	0.99122 4	0.99125 7
0.00041	0.00043 4	0.00045 3	0.00060 9	0.00072 0	0.00089 8	0.00127 3	0.00177 8
0.99159	0.99161 2	0.99162 5	0.99180 2	0.99191 6	0.99210 2	0.99249 8	0.99303 5
0.99076	0.99074 4	0.99071 9	0.99058 4	0.99047 6	0.99030 5	0.98995 1	0.98947 8

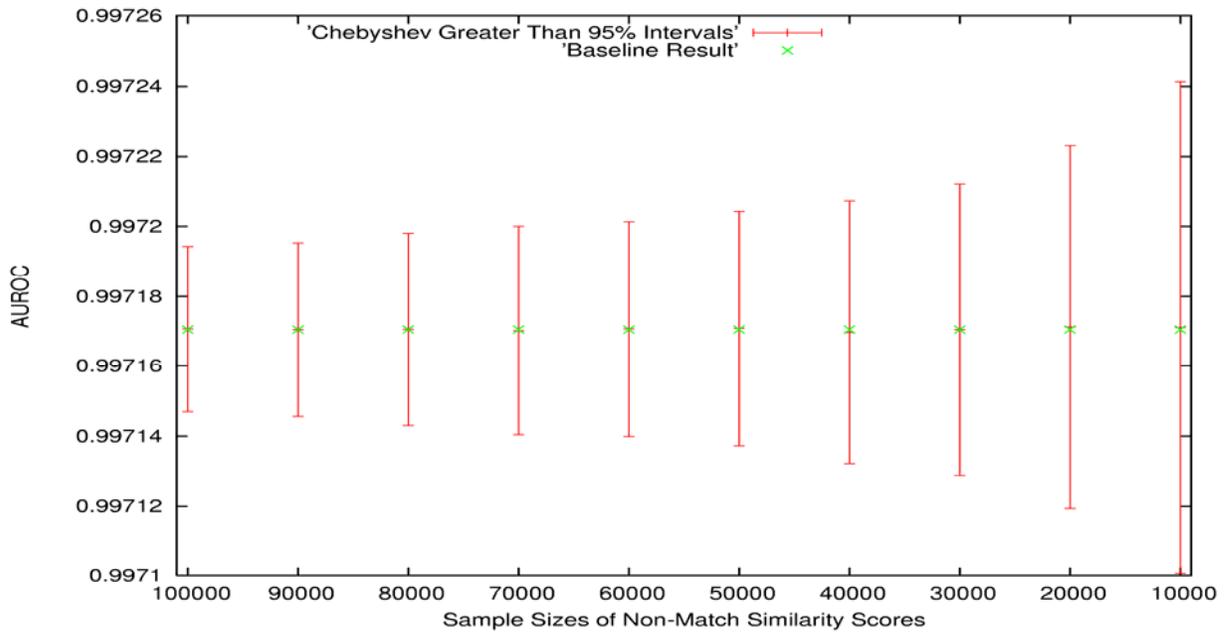
**Table 3 Matcher 1's Chebyshev's greater-than-95% intervals in terms of sample mean, error bar, upper bound and lower bound of 500 Monte Carlo iterations for different sample sizes of non-match similarity scores using TVAFV.**

0000	70000	60000	50000	40000	30000	20000	10000
0.98386	0.98386 4	0.98386 3	0.98386 2	0.98387 1	0.98386 3	0.98385 7	0.98387 5
0.00024	0.00025 5	0.00028 1	0.00031 7	0.00034 4	0.00040 1	0.00050 9	0.00071 6
0.98410	0.98411 9	0.98414 4	0.98417 9	0.98421 5	0.98426 3	0.98436 6	0.98459 2
0.98362	0.98360 8	0.98358 2	0.98354 5	0.98352 7	0.98346 2	0.98334 8	0.98315 9

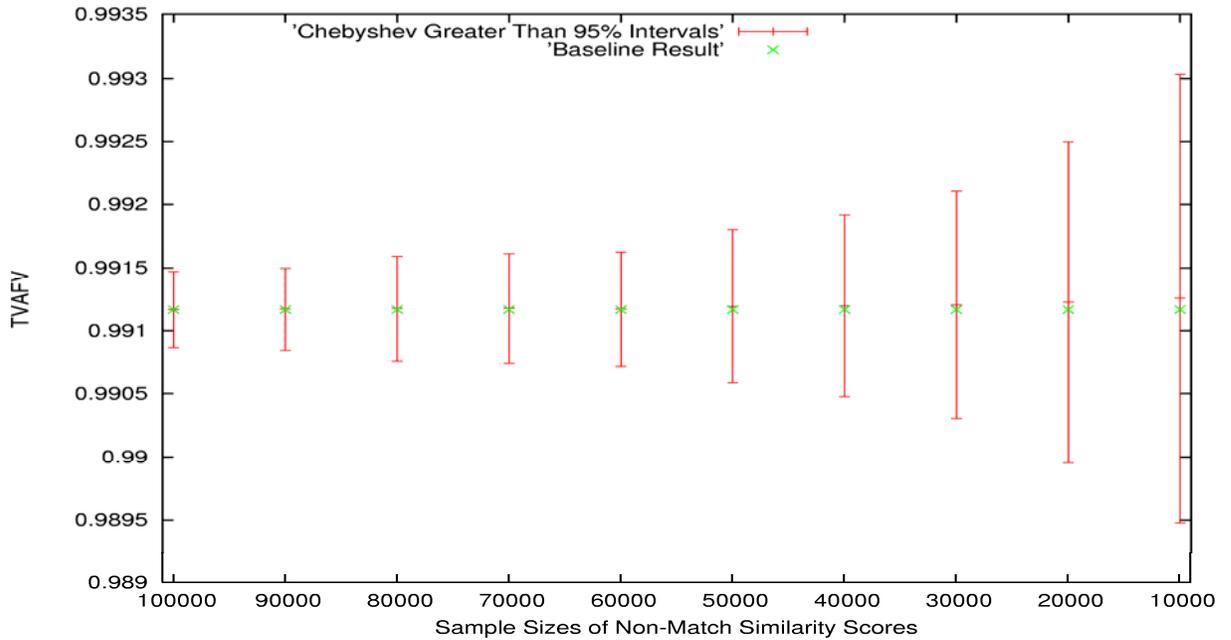
**Table 4 Matcher 2's Chebyshev's greater-than-95% intervals in terms of sample mean, error bar, upper bound and lower bound of 500 Monte Carlo iterations for different sample sizes of non-match similarity scores using AUROC.**

0000	70000	60000	50000	40000	30000	20000	10000
0.89000	0.88993 5	0.88987 6	0.88981 9	0.88999 9	0.88964 9	0.88949 3	0.88907 1
0.01311	0.01330 9	0.01348 6	0.01452 7	0.01499 0	0.01648 6	0.02072 2	0.03058 9
0.90311	0.90324 3	0.90336 2	0.90434 6	0.90498 8	0.90613 5	0.91021 5	0.91966 0
0.87689	0.87662 6	0.87639 0	0.87529 2	0.87500 9	0.87316 3	0.86877 1	0.85848 2

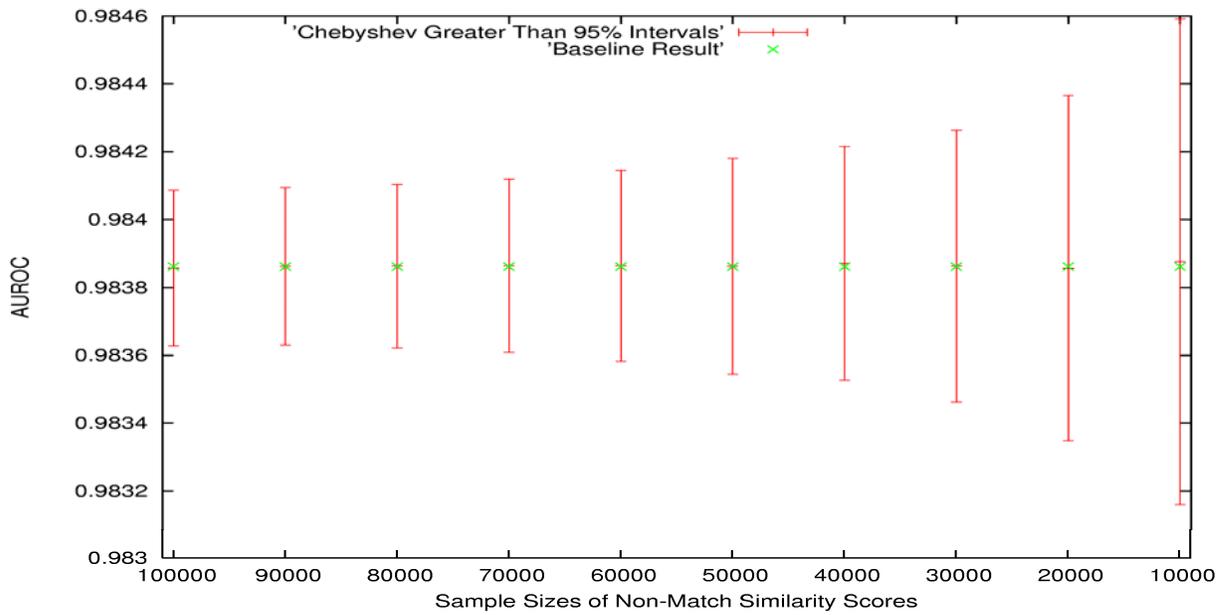
**Table 5 Matcher 2's Chebyshev's greater-than-95% intervals in terms of sample mean, error bar, upper bound and lower bound of 500 Monte Carlo iterations for different sample sizes of non-match similarity scores using TVAFV.**



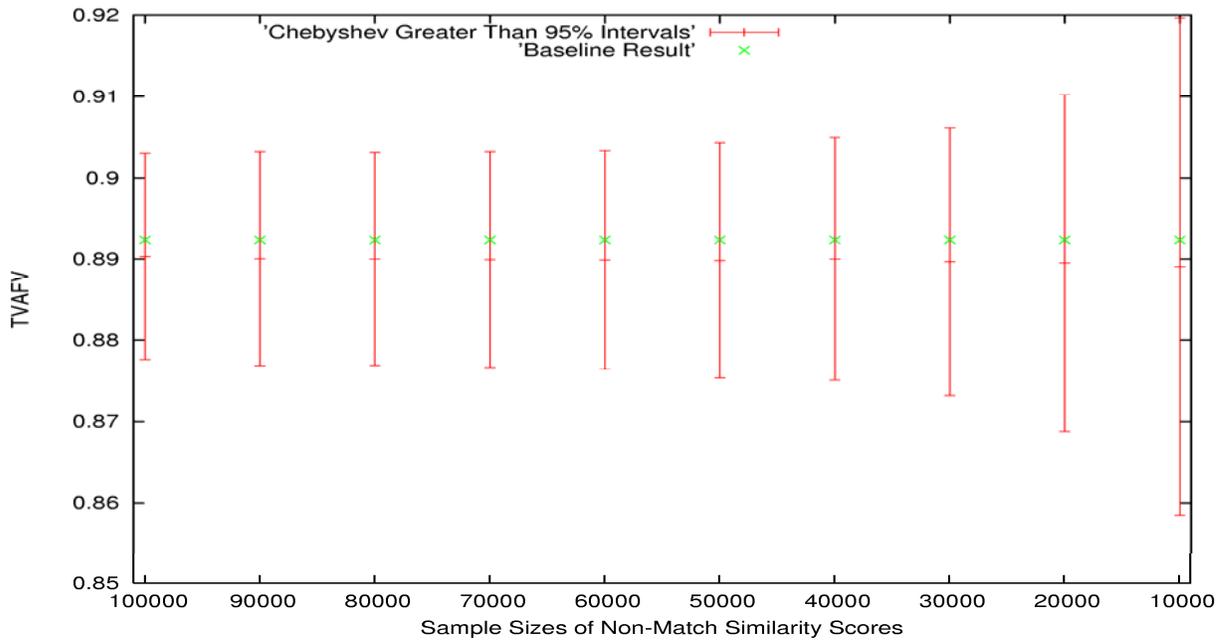
**Figure 3** Matcher 1's Chebyshev's greater-than-95% intervals of 500 Monte Carlo iterations for different sample sizes of non-match similarity scores using AUROC, along with the baseline result of AUROC.



**Figure 4** Matcher 1's Chebyshev's greater-than-95% intervals of 500 Monte Carlo iterations for different sample sizes of non-match similarity scores using TVAFV, along with the baseline result of TVAFV.



**Figure 5** Matcher 2's Chebyshev's greater-than-95% intervals of 500 Monte Carlo iterations for different sample sizes of non-match similarity scores using AUROC, along with the baseline result of AUROC.



**Figure 6** Matcher 2's Chebyshev's greater-than-95% intervals of 500 Monte Carlo iterations for different sample sizes of non-match similarity scores using TVAFV, along with the baseline result of TVAFV.

As illustrated in the tables, if using the criterion of AUROC, the error bars monotonically increase from 0.000024 to 0.000070 while the sizes of simple random samples selected from 35994000 non-match similarity scores decrease from 100000 down to 10000 for high-quality Matcher 1, but from 0.000230 to 0.000716 for low-quality Matcher 2. If using the

criterion of TVAFV, the error bars vary between 0.000301 and 0.001778 within the same range of sample sizes for Matcher 1, but between 0.012728 and 0.030589 for Matcher 2. As shown in the figures, all Chebyshev's greater-than-95% intervals contain the corresponding baseline results, and no upper bound of Chebyshev's greater-than-95% interval is exceeding one.

The error bars as well as the deviations between the sample means and the corresponding baseline results all depend on the qualities of the fingerprint-image matchers and the criteria invoked for evaluation of ROC curve. For high-quality matcher such as Matcher 1 as opposed to low-quality matcher such as Matcher 2, the error bars are smaller and the sample means are closer to the corresponding baseline results. The same relationship exists between AUROC and TVAFV.

The higher the matcher's quality is, the more convergent the outcome is, therefore the less the variance will be. To reach the same error bar, the higher-quality matcher needs much smaller number of non-match similarity scores than the lower-quality matcher. As for comparing AUROC with TVAFV, the former is taking the whole ROC curve into account [3], but the latter is only picking a TAR value on an ROC curve at an operational FAR value. Hence, TVAFV is more sensitive to SRS than AUROC.

The tolerances used to determine the sample size for high-quality matchers must be smaller than the ones for low-quality matchers, since the values of AUROC and TVAFV of high-quality matchers are very close to 1. Therefore, if invoking the criterion of AUROC, 10000 non-match similarity scores are enough for both Matcher 1 and 2, once the tolerances for Matcher 1 and 2 are set to be 0.0001 and 0.001, respectively. If using the criterion of TVAFV, 30000 non-match similarity scores are enough for both Matcher 1 and 2, while the tolerances for Matcher 1 and 2 are set to be 0.001 and 0.02, respectively.

The choice of sample size is dependent on the qualities of fingerprint-image matchers as well as on which criterion is invoked. To be more conservative, as well as to get balance among different qualities of matchers and between two different criteria of AUROC and TVAFV, in general, it seems that for 6000 match similarity scores, 50000 to 70000 non-match similarity scores randomly selected from 35994000 non-match similarity scores are enough to ensure that the error bars of Chebyshev's intervals are within the accepted tolerance range with greater-than-95% probability.

### **3.3 The Stability of Monte Carlo Calculation**

The above results were derived from 500 iterations of Monte Carlo calculations. How stable are the results with respect to the number of Monte Carlo iterations? The smaller the sizes of simple random samples selected from 35994000 non-match similarity scores are, the larger deviation the distributions of selected non-match similarity scores have from the distribution in the baseline, therefore the less stable the outcome is. As a result, the case of 10000 non-match similarity scores is chosen to show the stability.

The stability metrics from 100 to 500 Monte Carlo iterations for sample size of 10000 are presented in Table 6, for Matcher 1 and 2 as well as for two different criteria of AUROC and TVAFV, respectively. As expected, the stability metric of Matcher 1 is smaller than the one of Matcher 2 for a fixed criterion, and the stability metric of AUROC is smaller than the one

of TVAFV for a specified matcher. This indicates again that the higher-quality matcher has less variance, and AUROC criterion is more convergent than TVAFV criterion.

Criterion	Matcher	Monte Carlo Iterations				
		100	200	300	400	
AUROC	1	0.000083	0.000072	0.000068	0.000073	
	2	0.000689	0.000721	0.000669	0.000673	
TVAFV	1	0.001686	0.001738	0.001764	0.001910	
	2	0.035012	0.032547	0.035804	0.032913	

**Table 6 The stability metrics of 10000 non-match similarity scores.**

In Table 6, it shows that the outcome of Monte Carlo calculation for 10000 non-match similarity scores is very stable from 100 iterations up to 500 iterations with respect to specified matcher and criterion. The worst deviations of the one-trial test result from the baseline result with greater-than-95% probability vary by no-larger-than 0.000015, 0.000061, 0.000224, and 0.003257 (the maximum minus the minimum in each row of Table 6) from 100 to 500 iterations for Matcher 1 and 2 and for two different criteria, respectively, even when the sample size is down to only 10000 non-match similarity scores. As a consequence, the results presented above out of 500 Monte Carlo iterations are reliable.

#### 4. Discussion

The methodology of invoking Chebyshev’s greater-than-95% interval in combination with simple random sampling serves our objective well. The result of taking one trial with reduced number of non-match similarity scores must be close to the baseline result. In terms of evaluation of ROC curves, that is,  $\Delta(ROC\ curve)$  must be within an accepted tolerance with greater-than-95% probability. The half of Chebyshev’s greater-than-95% interval is  $4.48 \hat{\sigma}$ . And the margin of error of 95% confidence interval estimate of the population mean is  $1.96 \hat{\sigma} / \sqrt{n}$ , where n is 500 in our case. The former is about 51 times larger than the latter. Therefore, the sampling error of the sample mean, namely, the absolute value of the difference between the unbiased point estimator of the population mean (i.e., the sample mean) and the population mean, is relatively negligible in each case, while using Chebyshev’s greater-than-95% interval.

In this article, only the case is explored, in which the number of non-match similarity scores needs to be reduced. As a matter of fact, the same technique can be applied to other scenarios of the biometric evaluation of fingerprint data, as long as the standard deviation of the population is small and the objective is only taking *one* trial instead of taking average of many trials. For instance, what if we want to see the results while the numbers of both non-match similarity scores and match similarity scores are reduced? As a matter of fact, the requirement that the standard deviation of the population be small is the disadvantage of employing Chebyshev’s inequality.

As has been demonstrated, the outcome is very much dependent on the qualities of fingerprint-image matchers. The higher-quality matchers are more convergent, thus have less variance than the lower-quality matchers. Hence, the higher-

quality matchers need fewer number of non-match similarity scores than the lower-quality matchers in our application. Accordingly, the accepted tolerance is also dependent on the quality of matchers. Presented in this article are only two matchers, namely, Matcher 1 and 2. And Matcher 1's quality is higher than Matcher 2's. In our tests for this article, four fingerprint-image matchers were used, two of which were high-quality matchers, and the other two were low-quality matchers. They exhibit the similar behavior to that shown in this article.

In biometric evaluation of fingerprint data, the sample sizes are also determined by other factors. For instance, if using the TVAFV criterion to evaluate ROC curve and setting the operational FAR value to be 0.001, for very high-quality fingerprint-image matchers, the TAR value could reach as high as 0.999. If there are only 6000 match similarity scores, then the number of failures related to Type I error is only about 6, which is very much less significant. For such quality of matchers, in order to increase the significance of test, the number of match similarity scores must increase accordingly.

## 5. Conclusion

In conclusion, in the current framework of SDK tests, with respect to 6000 match similarity scores, it seems that 35994000 non-match similarity scores are much more than what is needed. The number of non-match similarity scores can be dramatically reduced down to 50000 to 70000, as long as that amount of non-match similarity scores is a simple random sample of 35994000 non-match similarity scores. It holds good for different qualities of fingerprint-image matchers as well as for criteria of both AUROC and TVAFV. And it is valid with greater-than-95% probability.

## References

1. C.L. Wilson, *et al.*, Fingerprint vendor technology evaluation 2003: summary of results and analysis report, NISTIR 7123, National Institute of Standards and Technology, June 2004.
2. C. Watson, C. Wilson, K. Marshall, M. Indovina, R. Snelick, Studies of one-to-one fingerprint matching with vendor SDK matchers, NISTIR 7119, National Institute of Standards and Technology, May 2004.
3. J.C. Wu, C.L. Wilson, Nonparametric Analysis of Fingerprint Data, NISTIR 7226, National Institute of Standards and Technology, May 2005.
4. S.D. Walter, The partial area under the summary ROC curve, *Statist. Med.* 24 (2005) 2025-2040.
5. G.S. Gazelle, P.M. McMahon, U. Siebert, M.T. Beinfeld, Cost-effectiveness analysis in the assessment of diagnostic imaging technologies, *Radiology* 235 (2005) 361-370.
6. J.A. Hanley, B.J. McNeil, A method of comparing the area under two ROC curves derived from the same cases, *Radiology* 148 (1983) 839-843.
7. J.A. Hanley, B.J. McNeil, The meaning and use of the area under a receiver operating characteristic (ROC) curve, *Radiology* 143 (1982) 29-36.
8. B. Ostle, L.C. Malone, *Statistics in research: basic concepts and techniques for research workers*, fourth ed., Iowa State University Press, Ames, 1988.

9. P.J. Bickel, K.A. Doksum, *Mathematical statistics: basic ideas and selected topics*, Holden-Day, Inc., San Francisco, 1977.
10. W.H. Press, S.A. Teukolsky, W.T. Vetterling, B.P. Flannery, *Numerical recipes in C++: the art of scientific computing*, second ed., Cambridge University Press, New York, 2002, pp. 647-648.