

1 Estimating error rates for bullet comparisons in forensic 2 science

3 Yawei Ge^{a,b}, Heike Hofmann^{a,b}

4 ^a*Department of Statistics, Iowa State University,*

5 ^b*Center for Statistics and Applications of Forensic Evidence, Iowa State University*

6 Abstract

7 This is the abstract.

It consists of two paragraphs.

8 *Text based on elsarticle sample manuscript, see [http://www.elsevier.com/](http://www.elsevier.com/author-schemas/latex-instructions#elsarticle)*
9 *author-schemas/latex-instructions#elsarticle*

10 1. Introduction

11 Firearm examination is an important topic in forensic science to help reveal
12 the underlying pattern of firearm evidence. XXX go into a bit more detail
13 on what you mean by the underlying pattern Currently, those examinations
14 are executed in forensic labs by Firearms and Toolmark Examiners (FTEs)
15 under the regulation of AFTE (AFTE Criteria For Identification Committee,
16 1992). However, this process is in its in nature based on human decisions and
17 therefore subjective. XXX Change from passive voice to active, i.e. The PCAST
18 report and the NRC called out this subjectivity and ... The subjectivity in the
19 conventional firearm evidence identification processes is called to be reduced or
20 to be complemented by more objective procedures by the National Research
21 Council (National Research Council, 2009). The President's Council of Advisers
22 on Science and Technology (PCAST) also emphasized the importance to establish
23 the scientific validity and reliability of subjective forensic feature-comparison
24 methods by blind empirical test (President's Council of Advisors on Science and
25 Technology (PCAST), 2016). The PCAST also recognized the importance of
26 developing of objective computer-based algorithms in its following addendum
27 (President's Council of Advisors on Science and Technology (PCAST), 2017). In
28 cartridge case comparisons, the congruent matching cells (CMC) was developed
29 to conduct the comparisons automatically (Song, 2015). And the error rates
30 are estimated in following papers by establishing distributional models (Song
31 et al., 2018; Zhang, 2019). In bullet comparisons, an automatic matching
32 algorithm was proposed by Hare (Hare et al., 2016). The discussion about
33 the degraded land comparison was also made in the following paper (Hare et
34 al., 2017). The empirical error rates associated with this method based on
35 three test sets are discussed by Vanderplas (Vanderplas et al., 2020). The most
36 prominent single feature in the random forest proposed by Hare (Hare et al.,

2016) is the cross-correlation function (Vorburger et al., 2011) according to the empirical studies (Hare et al., 2016; Vanderplas et al., 2020). Other automated comparison methods or improvement such as Robust LOESS in groove engraved area identification (Rice et al., 2020), Chumbley score method (Krishnan and Hofmann, 2019) are also proposed.

However, due to the limitation of sample sizes and designs of study, a sound distributional assumption about the resulting similarity scores was not well established in the bullet comparison field for automated matching algorithms. The error rates are estimated empirically from dozens of bullets which were originally designed for examiner proficiency tests. And many of those tests are also criticized for the closed-set design and potential information provided to the examiners as pointed out by PCAST (President’s Council of Advisors on Science and Technology (PCAST), 2017, 2016). It’s important to understand the underlying distributions of the similarity scores to estimate the theoretical error rates of any automatic algorithm. And that also helps later to choose proper thresholds for the methods instead of simply choosing 0.5 (in a 0-1 range of similarity scores) as a belief.

Generally, in forensic science, we consider the problem of distinguishing two pieces of evidence coming from a same source or different sources. Specifically, in bullet comparisons, we want to know if two bullets are coming from a same gun or not. When shooting a bullet, the barrel will leave markings on the surface of bullets due to the imperfectness of manufacturing. We distinguish land engraved areas (LEA) and groove engraved error areas. And the former is compared by examiners by inspecting the striations under confocal microscopes in labs (AFTE Criteria For Identification Committee, 1992). And the LEA is also used in the developed computer-based methods. A fired bullet has several LEAs according to the barrels’ design. So, we are not doing one step to identify bullets, instead we consider each LEA and then combine the land level conclusions to a bullet level one. Accordingly, we start our discussion about underlying distributions and error rates from the land level.

There are two types of distributions we would consider in the bullet comparison cases. One is called known-match (KM) distributions for those actually matched lands. The other is called known-non-match distributions (KNM) for those actually not matched lands. Those two distributions are only available in experimental settings where we know the ground truth. They are the starting point for any discussion about theoretical error rates. And particularly, we consider the cross-correlation function (CCF) for bullet signatures as a typical similarity score with range 0-1 used to distinguish between same source and different source comparisons. The cross-correlation function (CCF) is calculated after extracting the signatures of bullet lands by maximizing the cross correlation of two sequence of signatures by horizontally adjusting the signatures (Vorburger et al., 2011).

In section 2... In section 3...

2. The distributional forms of similarity scores

The quantitative methods used to objectively measure the similarity between LEAs report various quantities, such as counts, correlations, distances, probabilities and more general similarity scores (Hare et al., 2016; Ma et al., 2004; Nichols, 2003, 1997). To understand how those quantities reflect the strength of evidence and to study the underlying error rates in making decisions based on those quantities, distributional forms are usually set up (Song et al., 2018; Zhang, 2019). Particularly, we are focusing on the similarity scores which range from 0 to 1 and the probabilities reported as the likelihood of an actual match. The similarity scores reported in the forensic researches are classified into two categories as known matches (KM) and known non-matches (KNM), and the corresponding distributions are named as KM distributions and KNM distributions. When we make any decisions based on any quantitative measurement, we are actually making a distinguish between those two potential distributions. The strength of any identification process is also measured by the disparities of those distributions. However, in practice, we can hardly discriminate those two distributions entirely, thus, we are never 100% sure which distribution the observed score comes from. This is where the identification error raises.

The cross-correlation function (CCF) is the most prominent single feature in the automatic random forest algorithm (Hare et al., 2016) which has theoretical range from -1 to 1. But in real applications, this similarity score always has values in 0 to 1, e.g. in our case, none of 121992 comparisons has negative values. So, it is selected as a representative of similarity scores with range from 0 to 1 in this paper for further analysis. The similarity scores in 0 to 1 can be explained as probabilities that quantify the likelihood that a pair of LEAs are actually a match. Or we can think of them as general similarity measurement. As the name indicated, the higher the similarity score is, the stronger evidence is to support the same source assumption. For different combination of ammunition and firearms, the scores are distributed differently. It is expected that systematic differences exist there for different cases (Vanderplas et al., 2020). So, it is necessary to study the scores under controlled conditions.

We can see from the Figure ... (a figure of nonparametric fit of the CCF?), which is a typical one we usually have for the similarity scores, that the distributions of KM and KNM are apart for the majorities. In the bullet LEA comparison problems, we usually have a well separated bullet scores but for the land scores, there are usually some overlaps (Vanderplas et al., 2020). We propose beta distributions for those scores. Because the beta distribution is a well-used distribution in statistics to describe a quantity from 0 to 1 which is usually a probability or proportion quantifying our knowledge for another distribution in Bayesian analysis. And it is very flexible to capture unimodal asymmetric shapes in 0 to 1. However, it may not be adequate to explain a heavy tail or even a second mode. Thus, we further consider the beta mixture distribution which is a more complex distribution than the beta distribution as a special case. In the beta mixture distribution, we introduce a hierarchical structure with a prior probability to combine a few beta distributions as one.

125 The two-component beta mixture distributions are defined below.

126 The reported similarity scores are denoted as Y_{ij} for j^{th} LEA comparison
 127 within class i , where $i = 1$ is KM and $i = 2$ is KNM. Y_{ij} 's are considered
 128 independent and identically distributed within each class, i.e.

$$Y_{1j} \stackrel{iid}{\sim} \text{Betamix}(p_1, \mu_{11}, \phi_{11}, \mu_{12}, \phi_{12})$$

$$Y_{2j} \stackrel{iid}{\sim} \text{Betamix}(p_2, \mu_{21}, \phi_{21}, \mu_{22}, \phi_{22})$$

129 where μ_{ik} and ϕ_{ik} are distribution parameters for i^{th} class and k^{th} component,
 130 $k = 1$ or 2 . And p_i is the prior probability for the first component in i^{th} class,
 131 thus, $1 - p_i$ is the prior probability for the second component in i^{th} class. Note
 132 that we are using the mean and precision parameterization of beta distributions
 133 which simplifies the math in calculation and is more intuitive (μ is the mean,
 134 and the variance is roughly proportional to the reciprocal of ϕ). It's equivalent
 135 to the usual α and β parameterization through the following transformation:

$$\mu = \frac{\alpha}{\alpha + \beta}$$

$$\phi = \alpha + \beta$$

136 3. LAPD data set and the estimated distributions

137 For the following sections of the paper, we will base our analysis on the
 138 LAPD data sets. This is a large data set of ... It is the first time such a large
 139 data set available to the researchers, which makes it possible for a statistical
 140 analysis for the distribution of similarity scores. (More on the data. Is there any
 141 formal documentation for LAPD? And number of land comparisons for different
 142 cases should be discussed here)

143 We consider cross correlation functions (CCF) for land comparisons, which is
 144 produced by calculating the maximized CCF between two signatures extracted
 145 from a pair of bullet LEAs.

146 The estimation was done using Nelder-Mead algorithm in R [reference].
 147 This is a general purpose numeric method which works reasonably well for
 148 multidimensional optimization problems. In our case, the objective function
 149 is the log likelihood function of the beta mixture distribution. Therefore, we
 150 finally found the maximum likelihood estimates (MLE). Instead of a simpler
 151 beta model, we start with the more complex two-component beta mixture model,
 152 see how well it fits the data and test if some components are necessary. We
 153 estimated beta mixture distributions for both KM and KNM as in Table 1 and
 154 Table 2 respectively. For the beta distributions, the KM distribution has a mean
 155 at 0.635 and the KNM distribution has a mean at 0.404. We can also see that
 156 for the two component beta mixture distributions for both KM and KNM, there
 157 are components with mean around 0.5 and the other components more separated
 158 from each other. This indicates the model successfully accounted for different

159 situations of comparisons involving tank rash and random identification etc.
160 This is an ideal property we would like to see and make use in explaining the
161 similarity scores. Another point about the estimates that worth mentioning is
162 that the ϕ of the KM beta distribution is much smaller than any of the two
163 components of the two-component KM beta mixture distribution. This indicates
164 that the two component distribution could be better because of smaller variances
165 for components.

Full Data KM Distribution Estimation						
Model	Component Prior Probability	μ	ϕ	logLik	p-value	BIC
3-comp	0.247	0.406	24.097	6278	0	-12479
	0.521	0.664	13.342			
	0.232	0.819	33.823			
2-comp	0.419	0.466	15.812	6243	0	-12438
	0.581	0.759	18.476			
Beta		0.635	6.529	5712	0	-11404

Table 1: Parameter estimations for the Beta distribution (1-component beta mixture distribution), 2-component and 3-component distributions for the KM CCF. "2-comp" refers to the 2-component beta mixture distribution, and the same for "3-comp". Column "logLik" is the maximized log likelihood for each distribution. Column "p-value" is the p-value for asymptotic likelihood ratio tests between the current model and one-step more complex model, where 0 indicates that we would reject the hypothesis that the current one is sufficient to describe the data

Full Data KNM Distribution Estimation						
Model	Component Prior Probability	μ	ϕ	logLik	p-value	BIC
3-comp	0.650	0.345	35.017	78961	0	-157830
	0.314	0.490	24.979			
	0.032	0.667	18.454			
2-comp	0.674	0.358	39.908	78808	0	-157558
	0.336	0.494	13.324			
Beta		0.404	15.801	75349	0	-150675

Table 2: Parameter estimations for the Beta distribution (1-component beta mixture distribution), 2-component and 3-component distributions for the KNM CCF. "2-comp" refers to the 2-component beta mixture distribution, and the same for "3-comp". Column "logLik" is the maximized log likelihood for each distribution. Column "p-value" is the p-value for asymptotic likelihood ratio tests between the current model and one-step more complex model, where 0 indicates that we would reject the hypothesis that the current one is sufficient to describe the data

166 The estimated beta distributions are shown in the Figure 1 and the estimated
167 two-component beta mixture distributions are also shown in the Figure 2. As

168 expected, the estimated distributions show the properties we desire for the
 169 similarity scores. The majorities of the estimated distributions are apart from
 170 each other, while the minority part between the two distributions has some
 171 overlap. It's worth to note that both curves have heavy tails to the farther
 172 boundaries and the KM curve has the heavier tail compared to KNM. We can
 173 see the two-component distributions fit the data very well, while the single
 174 beta distributions are not as good as the two-component distributions for both
 175 KM and KNM. The single beta distribution for KM clearly failed to capture
 176 the potential second mode of the histogram indicating it is not sufficient. The
 177 single beta distribution for KNM is not too bad but still failed to capture the
 178 distributional information for some parts. This indicates that even though beta
 179 distributions are flexible for variables in 0 to 1, they are still restricted too
 180 much to be able to describe the cases here. The two-component beta mixture
 181 distributions look more promising.

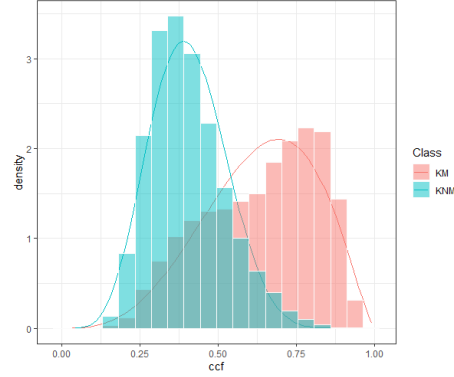


Figure 1: Estimated beta distributions for full data

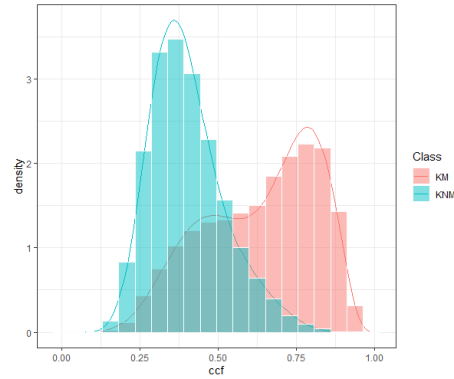


Figure 2: Estimated two-component beta mixture distributions for full data

182 Three candidate distributions for each of KM and KNM are considered. For

183 the increasing of complexity, we have single beta, two-component beta mixture
 184 and three-component beta mixture distributions. And they are also nested in
 185 that order. Naturally, we first look at the maximized log-likelihood of each model,
 186 and we can do asymptotic log-likelihood ratio chi-square tests for single beta
 187 against two-component beta mixture distributions, and two-component against
 188 three-component beta mixture distributions. The p-values of those test are
 189 shown in the Table 1 and Table 2 as the column “p-value”. Surprisingly (or not),
 190 we found all those p-values are 0 which strongly suggests a more complex model
 191 when there is one. However, considering the size of the data we are using to fit
 192 these models, we can expect the statistical significance will be easily achieved
 193 since any small difference of the sufficiency will be detected. So we have to take
 194 the sample size effect into consideration. As Bayesian information criterion (BIC)
 195 is a well used criterion which takes the model complexity, sufficiency and the
 196 sample size (by a log function) into account. As shown in the column “BIC” in
 197 the tables, the BICs for the single beta distributions are larger than that of the
 198 two-component beta mixtures by relatively large proportions. And the BICs for
 199 the two-component beta mixture distributions are a little larger than that of the
 200 three-component beta mixture distributions for both KM and KNM. Obviously,
 201 we would prefer the two-component beta mixture distributions, but we would
 202 cast a doubt when it comes to the three-component beta mixture distributions.
 203 We still prefer the two-component one instead of the three-component one. The
 204 reasons are: 1) these differences of BICs between two and three component
 205 distributions are really small in proportion (by 0.3% for KM and 0.017% for
 206 KNM), 2) the estimation cost is not accounted, which will be higher for a more
 207 complex model, 3) the BICs still don’t take the sample size effect fully into
 208 account since the log function for sample size goes to flat when the sample sizes
 209 are large. By simple calculation, we can see that the log-likelihood increased
 210 by a multiplicative factor more than 10 when the sample sizes increased from
 211 KM case to KNM case, however, at the same time, the log of sample size
 212 only roughly increased from 9 to 11. So we choose the two-component beta
 213 mixture distributions for both KM and KNM. And also as we have seen, these
 214 distributions have potential good forensic interpretations.

215 It’s also helpful to see how the individual components look like in the beta
 216 mixture distributions as in Figure 3. The KNM and KM distributions seem to
 217 share a common component while keeping the other components far apart from
 218 each other. The separated components represent the ideal cases where the bullet
 219 land engraved areas preserved the information of source well and result in clearly
 220 distinct separation. The shared components represent the cases where the KM
 221 results in lower scores because of some degree of tank rash, pitting, breakoff or
 222 other damages on the bullets and the KNM results in higher score because of the
 223 random identification effect. According to the estimated prior probabilities, both
 224 distributions put less weight on the common component while putting larger
 225 weight on the the components characterizing the differences of KM and KNM
 226 respectively, which agrees on our expectation that majorities of the distributions
 227 are separated while the minorities overlapped. These properties together well
 228 explained the observed empirical distribution of similarity scores in our cases.

229 Particularly, the heavier tail of the KM comparisons is explicitly included in the
 230 form of the model by one of the components.

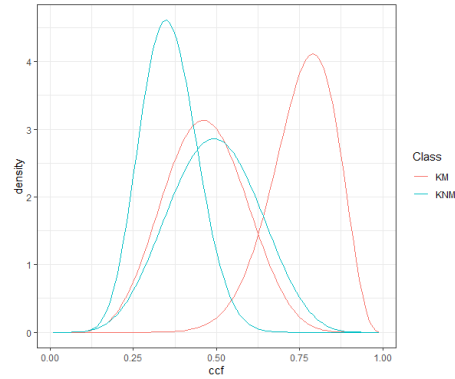


Figure 3: Estimated components

231 4. Evaluate the error rates

232 (report theoretical error rates)

- 233 • given the cutoff point where the two distributions have equal density (also
 234 a likelihood ratio with value 1, which does an optimal classification under
 235 0-1 loss? or at the point where the probabilities equal to each other instead
 236 of the densities)
- 237 • Consider four types of error rates (will need a table to illustrate this): false
 238 positive error rate, false negative error rate, false identification error rate,
 239 false exclusion error rate.
- 240 • Given the iid assumption, could we develop a further model for the bullet
 241 level distribution and estimate the corresponding error rates?

242 5. Estimations with changing sample size

243 (quantify the variation, reproduce some results in the previous sections with
 244 changing sample size)

245 6. Conclusion

246 References

247 AFTE Criteria For Identification Committee, 1992. Theory of identification,
 248 range striae comparison reports and modified glossary definition afte criteria for
 249 identification committee report. AFTE Journal 24, 336-340.

250 Hare, E., Hofmann, H., Carriquiry, A., 2016. Automatic matching of bullet
251 land impressions. *Annals of Applied Statistics* 11, 2332–2356.

252 Hare, E., Hofmann, H., Carriquiry, A., 2017. Algorithmic approaches to
253 match degraded land impressions. *Law, Probability and Risk* 16, 203–221.
254 doi:10.1093/lpr/mgx018

255 Krishnan, G., Hofmann, H., 2019. Adapting the chumbley score to match
256 striae on land engraved areas of bullets. *Journal of Forensic Sciences* 64, 728–740.
257 doi:10.1111/1556-4029.13950

258 Ma, L., Song, J.-F., Whitenton, E.P., Vorburger, T.V., Zhou, J., Zheng, A.,
259 2004. NIST Bullet Signature Measurement System for RM (Reference Material)
260 8240 Standard Bullets 49(4).

261 National Research Council, 2009. Strengthening forensic science in the united
262 states: A path forward, Strengthening Forensic Science in the United States: A
263 Path Forward. National Academies Press. doi:10.17226/12589

264 Nichols, R., 2003. Firearm and toolmark identification criteria: A review of
265 the literature, part II. *Journal of forensic sciences* 48, 318–27. doi:10.1520/JFS14149J

266 Nichols, R.G., 1997. Firearm and Toolmark Identification Criteria: A Review
267 of the Literature. *Journal of Forensic Science* 42, 466–474. doi:10.1520/JFS14149J

268 President’s Council of Advisors on Science and Technology (PCAST), 2016.
269 Report to the president, forensic science in criminal courts: Ensuring scientific
270 validity of feature-comparison methods. Executive Office of the President of the
271 United States, President’s Council

272 President’s Council of Advisors on Science and Technology (PCAST), 2017.
273 An Addendum To The PCAST Report On Forensic Science In Criminal Courts.

274 Rice, K., Genschel, U., Hofmann, H., 2020. A robust approach to automati-
275 cally locating grooves in 3D bullet land scans. *Journal of Forensic Sciences* 65,
276 775–783. doi:10.1111/1556-4029.14263

277 Song, J., 2015. Proposed “Congruent Matching Cells (CMC)” Method for
278 Ballistic Identification and Error Rate Estimation 47, 9.

279 Song, J., Vorburger, T.V., Chu, W., Yen, J., Soons, J.A., Ott, D.B., Zhang,
280 N.F., 2018. Estimating error rates for firearm evidence identifications in forensic
281 science. *Forensic Science International* 284, 15–32. doi:10.1016/j.forsciint.2017.12.013

282 Vanderplas, S., Nally, M., Klep, T., Cadevall, C., Hofmann, H., 2020. Com-
283 parison of three similarity scores for bullet lea matching. *Forensic Science*
284 *International* 308. doi:10.1016/j.forsciint.2020.110167

285 Vorburger, T.V., Song, J.F., Chu, W., Ma, L., Bui, S.H., Zheng, A., Renegar,
286 T.B., 2011. Applications of cross-correlation functions. *Wear, The 12th Inter-*
287 *national Conference on Metrology and Properties of Engineering Surfaces* 271,
288 529–533. doi:10.1016/j.wear.2010.03.030

289 Zhang, N.F., 2019. The use of correlated binomial distribution in estimating
290 error rates for firearm evidence identification. *Journal of Research of the National*
291 *Institute of Standards and Technology* 124. doi:10.6028/jres.124.026