

Multidimensional test statistics and statistical comparison of histograms

Sergey I. Bityukov, Nikolai V. Krasnikov, Anastasia V. Maksimishkina, and Vera V. Smirnova

Abstract—Several approaches for the comparison of histograms are considered. A new method for the distinguishing of flows of events via multidimensional comparative analysis of histograms is proposed. The example of the use of the method is presented.

Keywords—Data analysis, flow of events, measurement, Monte Carlo method, theory of errors.

I. INTRODUCTION

The term "histogram" was coined by the famous statistician Karl Pearson to refer to a "common form of graphical representation" [1]. Histograms are very useful in their canonical visual representation, but today histograms are considered as purely mathematical objects.

Histograms are used in different scientific fields. Besides physics data analyses, histograms play a very important role in databases, image processing, computer vision [1]. Correspondingly, goals and methods of the treatment of histograms are varied in dependence to the area of application. Here histograms are considered in frame of tasks related to physical experiments.

II. HISTOGRAM

Let us call the appearance of the realization of the random variable (or random variables) as the event. Suppose, there is given a set of non-overlapping intervals. A histogram represents the frequency distribution of data that populates those intervals. This distribution is obtained during data processing of the sample taken from the flow of events. These intervals usually are called as bins.

The filling procedure of a histogram influences the analysis of histogram. There are two extreme cases.

The first case: one event produces one histogram. For

This work was supported by the Ministry of Education and Science RF (Agreement on October 17, 2014 N 14.610.21.0004, id. PNIER RFMEFI61014X0004).

S. I. Bityukov is with the Department of experimental physics, Institute for high energy physics Natinal research centre "Kurchatov institute", Protvino, Moscow region 142281, Russia (corresponding author to provide e-mail: Sergej.Bitoukov@cern.ch).

N. V. Krasnikov is with the Theory physics department, Institute for nuclear research RAS, Moscow 117312, Russia (e-mail: Nikolai.Krasnikov@cern.ch).

A. V. Maksimushkina is with the Department of general and special physics, National research nuclear university MEPhI, Moscow 115409 (e-mail: AVMaksimushkina@mephi.ru).

V. V. Smirnova is with Department of mathematics and computing, Institute for high energy physics Natinal research centre "Kurchatov institute", Protvino, Moscow region 142281, Russia (e-mail: Vera.Smirnovar@ihep.ru).

example, the distribution of brightness in a photo is a result of data processing of one event. Here one sample consists from one event and one event is one photo.

The second case: one event is one measurement of random variable and resulting value is put to histogram. The filling of the histogram is a chain of independent measurements with gradual filling of the histogram. The second case is used, usually, in physical experimental researches for data processing. Correspondingly, the content of the bin is called the number of events in the bin, the sum of contents of bins in histogram is a volume of the histogram.

Common issues of are in construction of histograms, for example, the choice of optimal binning and the choice of the model for distribution of errors for observed values in the bins.

III. COMPARISON OF HISTOGRAMS

Given two histograms, how do we assess whether they are similar or not? What does it means "similar"? Several standard procedures exist for this task.

Suppose, a reference histogram is known. Usually, the proximity of test histogram and reference histogram is measured via a test statistics, that provides the quantitative expression of the "distance" between histograms [2]. The smaller the distance the more similar they are.

A. "Distance" Between Histograms

There are several definitions of distance in the literature, for example, the Kolmogorov distance [3], the Kullback-Leibner [4] distance, the total variation distance [5], the chi-square distance [6] and so on. Usually, it is the some test statistics, distribution of which can be calculated via formulae or constructed by Monte Carlo method. Other approach is based on the fact that a histogram of a measurement provides the basis for an empirical estimate of the probability density function (pdf) [7]. Computing the distance between two pdfs can be regarded as the same as computing the Bayes (or minimum misclassification) probability. This is equivalent to measuring the overlap between two pdfs as the distance. Sometimes, the Bhattacharyya distance [8] (or Hellinger distance [9]) is used as the distance between two pdfs. Note, that the Kolmogorov distance [3], the Anderson-Darling distance [10], the Kullback-Leibner distance [4] also allow to compare samples of events without their presentation in form of histograms. Recently, the test based on the maximum mean discrepancy (MMD) [11] was appeared. The important approach for comparison of histograms is tests based on ranks and/or permutations (Mann-Whitney [12], ...). In the vector approach, a histogram is treated as a fixed-dimensional vector.

Hence standard vector norms such as city block, Euclidean or intersection can be used as distance measures [13]. Similarity measures can be used in the comparing histograms. For example, the method of modulo similarity [14] is based on Lukasiewicz logic [15].

B. Testing of Consistency of Histograms or Distinguishability of Histograms

Also, a goal of histogram comparison is a testing of their consistency [16] or vice versa of their distinguishability [17]. Consistency here is the statement that both histograms are produced during data processing of independent samples which are taken from the same flow of events (or from the same population of events). In paper [18] is proposed approach that allows to estimate the distinguishability of histograms and, correspondingly, the distinguishability of parent events flows. The method is based on the statistical comparison of histograms. The multidimensional test statistics is used as a distance between histograms. In paper [19] is proposed an approach based on this method [18] for the detection of the changing of parameters in the context of wireless transmission.

If the goal of the comparison of histograms is the check of their consistency, then task is reduced to hypotheses testing: main hypothesis H_0 (histograms are produced during data processing of samples taken from the same flow of events) against alternative hypothesis H_1 (histograms are produced during data processing of samples taken from different flows of events). In principle, the choice between main and alternative hypothesis depends on the task. The determination of critical area allows to estimate Type I error (α) and Type II error (β) in decision about choice between H_0 and H_1 . The Type I error is a probability of mistake if done choice is H_1 , but H_0 is true. The Type II error is a probability of mistake if done choice is H_0 , but H_1 is true. The selection of a significance level (α) allows to estimate the power of the test ($1-\beta$). Usually, values of significance level are 10%, 5%, 1%. If both hypotheses are equivalent, then other combinations of the α and β are used. For example, in task about distinguishability of the flows of events works a relative uncertainty $(\alpha+\beta)/(2-(\alpha+\beta))$ [20]. Under the test of equal tails [21] the mean error $(\alpha+\beta)/2$ can be used.

C. Other Goals of Comparison of Histograms

Many other goals of comparison of histograms exist.

For example, the search for anomalous structures in test histogram in comparison with reference histogram is a very important task in particle physics. Possible solution is the comparison of the contents of two histograms, bin by bin. In this case, the probability that both bins were produced from a distribution with the same mean is calculated.

Also, the method for sorting events of multiparticle production according to the anisotropy of their momentum distribution by the use of histograms is presented in paper [22].

D. Comparison of Normalization and Comparison of Shape

The histograms comparison can usually be decomposed into comparison of normalization and comparison of shape. Sometimes the normalization and the shape are not independent, so the decomposition till works but it becomes more difficult to come up with a meaningful combination of the two tests. In the simplest case, normalization can be estimated by common suppositions. It may be the ratio of the volumes of samples corrected due to any additional knowledge (for example, efficiencies of registration of events). It may be the ratio of times for gathering samples and so on. A vast amount of statistical literature is devoted to the theme of shape comparisons (see, for example, [23]).

E. "Rehistogramming"

The hypotheses testing require the knowledge of the distribution of test statistics. As mentioned above the distribution of test statistics can be constructed by Monte Carlo. Let us consider the simple case of the filling of histograms - event-by-event in frame of the method of statistical comparison of histograms [18, 17]. The number of events in each bin of histogram can be considered as a realization (observed value) of the random variable with parameter "the expected number of events in given bin of histogram for given sample". The knowledge of uncertainty of the observed value in the case of statistically dual distributions [24] allows describing the uncertainty of the corresponding value of parameter. If we work with Poisson flows of events, then uncertainty of the parameter obeys the gamma distribution. If we work with Gaussian approximation, then the distribution of uncertainty obeys the normal distribution. As a result, we can use the Monte Carlo method for construction of two imitation models of the possible histograms sets. These two sets of histograms imitate the two general populations (two models) which provided us two histograms for comparison. This procedure can be named as "rehistogramming", similar to "resampling" in bootstrap technique [25]. The first imitation population (the first set of histograms) is used for construction of the distribution of test statistics for the case of H_0 hypothesis. The second imitation population (the second set of histograms) is used for construction of the distribution of test statistics for the case of H_1 hypothesis. The overlapping of these distributions gives us the estimation of the uncertainty in the hypotheses testing [18, 17]. The similar approaches for histograms comparison is described in papers [26, 27, 28] too.

F. "Significance of the Difference"

The convenient object for comparison of histograms is a distribution of the "significances of the difference". The "significances of the difference" are calculated for corresponding pairs of bins of the comparing histograms. The choice of type of "significance of the difference" depends on the task [29]. If the comparing histograms are taken from the same population of histograms (or the corresponding samples are taken from the same flow of events), the distribution of "significances of the difference" is close to standard normal

distribution.

IV. MULTIDIMENSIONAL COMPARISON

As mentioned above, the method of statistical comparison of histograms (SCH) [18, 17] is a multidimensional method. Let us consider a simple example of the use of the bidimensional test statistics.

A. Bidimensional Test Statistics

Suppose there are two histograms hist1 and hist2 (where M – number of bins) which are produced during processing of two independent samples of events. These histograms can be considered as two sets of numbers --

hist1: $n_{11} \pm \sigma_{11}, n_{21} \pm \sigma_{21}, \dots, n_{M1} \pm \sigma_{M1}$ and

hist2: $n_{21} \pm \sigma_{21}, n_{22} \pm \sigma_{22}, \dots, n_{M2} \pm \sigma_{M2}$.

Let $\hat{S}_i = \frac{\hat{n}_{i1} - K\hat{n}_{i2}}{\sqrt{\hat{\sigma}_{i1} + K^2\hat{\sigma}_{i2}}}$ be the observed significance of difference for bin# i , $i=1, M$. Here K is a coefficient of normalization (for example, it may be the ratio of volumes of histograms). Observed significance of difference \hat{S}_i is a realization of some random variable. If both samples of events are taken from the same flow of events then the test statistics \hat{S}_i obeys the distribution which close to standard normal distribution $N(0,1)$. It means that distribution of observed significances $\hat{S}_i, i=1, M$ also must be close to standard normal distribution. It allows using statistical moments of this distribution as components of multidimensional test statistics. For example, the mean value \bar{S} and the root mean square rms of this distribution [18] is a bidimensional test statistics $SRMS = (\bar{S}, rms)$. If $SRMS=(0, 0)$ then histograms are identical. If $SRMS \approx (0, 1)$ then samples are taken from the same flow of events. If the previous conditions are not valid, the flows have difference.

B. Monte Carlo Experiment

Two pairs (reference pair and test pair) of independent flows of samples with realizations of random variables (each realization is “event”) are produced to estimate the possibility of SCH method for distinguishing of samples from different information flows. The volume of each flow equals 5000 samples. First flow from each pair of flows is a reference flow of samples with 1000 events (1000 realizations of random variable $N(300,50)$). Second flow from first pair of flows also is a reference flow of samples with 2000 events (2000 realizations of the same random variable $N(300,50)$). Second flow from second pair of flows is a test flow of samples with 2000 events (2000 realizations of random variable $N(300,44)$). During data processing, the histogram is constructed for each sample. The examples of histograms for comparison are shown in Fig.1 (histograms from first pair of flows) and Fig. 2 (histograms from second pair of flows).

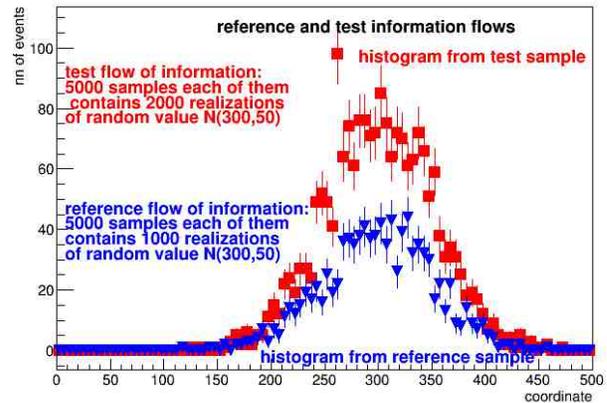


Fig. 1: Histogram from reference sample (1000 events) is a result of data processing of one sample from first flow of first pair. Histogram from test sample (2000 events) is a result of data processing of one sample from second flow of first pair of flows.

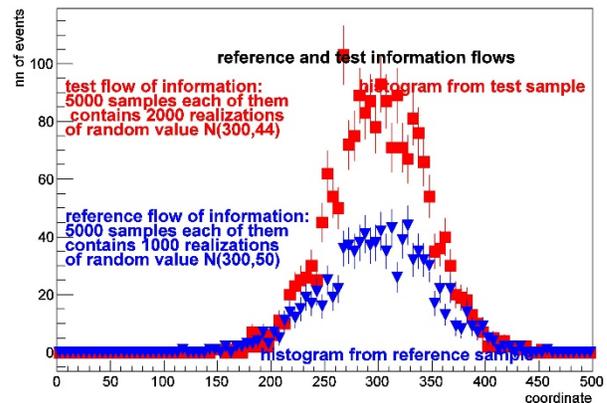


Fig. 2: Histogram from reference sample (1000 events) is a result of data processing of one sample from first flow of second pair of flows. Histogram from test sample (2000 events) is a result of data processing of one sample from second flow of second pair of flows.

After that for each pair of samples from corresponding pair of flows the comparison of histograms is performed with calculation of mean value and root mean square of the distribution of significances of the difference between corresponding bins of histograms. The distribution of the bidimensional test statistics $SRMS$ for comparison of samples from first pair of flows (left spot) and the distribution of the bidimensional test statistics $SRMS$ for comparison of samples from second pair of flows (right spot) is shown in Fig. 3.

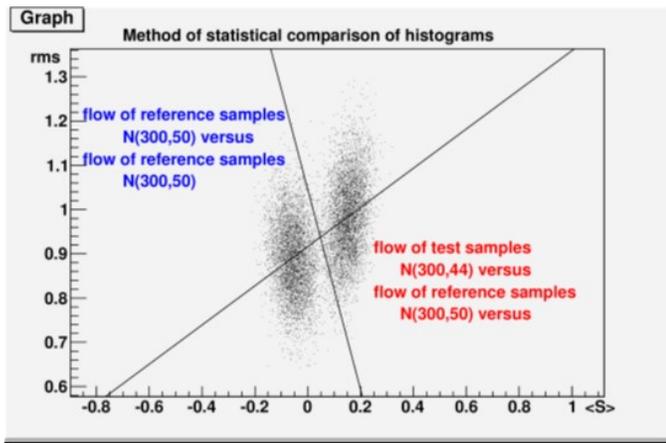


Fig.3: Two distributions SRMS for 5000 comparisons of samples from two similar flows (reference flows) and the case of the comparison of samples from reference flow and samples from test ($N(300,44)$) flow. In the picture are shown two straight lines. One of them connects the mean values of left and right spot. Second one is a critical line for hypotheses testing.

V. CONCLUSION

Several approaches for the comparative analysis of histograms are considered. We propose the method of statistical comparison of histograms for the distinguishing of flows of events under studying. This method uses a multidimensional test statistics based on the distribution of the significances of the difference. In principle, the method allows to include any other one-dimensional test statistics as an additional component of multidimensional test statistics. Also, this method allows comparing multidimensional histograms or the sets of histograms.

ACKNOWLEDGMENT

The authors would like to thank Prof. V. Kachanov, Prof. Yu. Korovin, Dr. S. Gleizer and Dr. N. Korneeva for helpful discussions.

REFERENCES

- [1] Y. Ioannidis, The history of histograms (abridged), Proceedings 2003 VLDB Conference, pp.19-30, 2003.
- [2] S.-H. Cha, S.N. Srihari, On measuring the distance between histograms, *Pattern Recognition*, Vol.35, 2002, pp. 1355-1370.
- [3] A.N. Kolmogorov, Confidence limits for an unknown distribution function, *Ann.Math.Stat.* Vol.12, 1941, pp.461-463.
- [4] S. Kullback, *Information Theory and Statistics*, Wiley, New York, 1959.
- [5] J. Rosenthal, Convergence rates for Markov chains, *SIAM Rev.* Vol.37, 1995, pp. 387-485.
- [6] W. Cochran, The chi-square test of goodness of fit, *Ann.Math.Stat.*, Vol.23, 1952, pp. 315-342.
- [7] R.O. Duda, P.E. Hart, *Pattern Classification and Scene Analysis*, 1st Edition, Wiley, New York, 1973.
- [8] T. Kailath, The divergence and Bhattacharyya distance measures in signal selection, *IEEE Trans.Commun.Technol.* COM-15, Vol.1, 1967, pp. 52-60.
- [9] E. Hellinger, Neue Begründung der Theorie quadratischer Formen von unendlichvielen Veränderlichen, *Journal für die reine und angewandte Mathematik*, (Crelle's Journal) Vol. 1909, Issue 136, pp. 210-271.
- [10] T.W. Anderson, D.A. Darling, Asymptotic theory of certain «goodness of fit» criteria based on stochastic processes, *Ann. Math. Statist.*, Vol.23, 1952, pp. 193-212.

- [11] A. Gretton, K. Borgwardt, M.J. Rasch, B. Scholkopf, A.J. Smola, A Kernel method for two-sample problem, *arXiv:0805.2368*, 2008.
- [12] H.B. Mann, D.R. Whitney, On a test of whether one of two random variables is stochastically larger than the other, *Annals of Mathematical Statistics*, Vol.18, 1947, pp. 50-60.
- [13] H. Bandemer, W. Nather, *Fuzzy data analysis*, Kluwer academic publishers, Dordrecht, 1992.
- [14] P. Luuka, M. Collan, Modulo similarity in comparing histograms, Proc. of IFSA-EUSFLAT2015, Eds. J.M. Alonso, H. Bustince, M. Reformat, Atlantis Press, 2015, pp. 393-397.
- [15] J. Lukasiewicz, *Selected Works*, North-Holland Publishing co., Amsterdam, 1970.
- [16] F. Porter, Testing Consistency of Two Histograms, *arXiv:0804.0380*, 2008.
- [17] S. Bityukov, N. Krasnikov, A. Nikitenko, V. Smirnova, On the distinguishability of histograms, *Eur.Phys.J.Plus*, Vol.128, No.143, 2013, pp.1-6.
- [18] S. Bityukov, N. Krasnikov, A. Nikitenko, V. Smirnova, A method for statistical comparison of histograms, *arXiv:1302.2651*, 2013.
- [19] B. Krupanek, R. Bogacz, Comparison algorithm of multimodal histograms from wireless transmission, *Przeglad Electrotechniczny*, Vol. 11, 2014, pp. 32-34.
- [20] S.I. Bityukov, N.V. Krasnikov, Distinguishability of hypotheses, *Nucl.Instr. & Meth. A*, Vol.534, 2004, pp. 152-155.
- [21] S.I. Bityukov, N.V. Krasnikov, *The use of statistical methods for the search for new physics at the LHC* (in Russian), Second edition, Krasand, Moscow, 2014.
- [22] R. Kopečna, B. Tomasik, Event shape sorting, *arXiv:1506.06776*, 2015.
- [23] O. Thas, *Comparing Distributions*, Springer Series in Statistics, 2010.
- [24] S. Bityukov, N. Krasnikov, S. Nadarajah, V. Smirnova, Statistically Dual Distributions and Estimations, *Applied Mathematics*, Vol.5, 2014, pp. 963-968.
- [25] B. Efron and R. Tibshirani, *An Introduction to the Bootstrap*, Chapman & Hall, New York, 1993.
- [26] Y. Cao, L. Petzold, Accuracy limitations and the measurement of errors in the stochastic simulation of chemically reacting systems, *J. of Computational Physics*, Vol.212, 2006, pp. 6-24.
- [27] K.-M. Xu, Using the bootstrap method for a statistical significance test of differences between summary histograms, *NASA Technical Reports Server*, ID: 20080015431, 2006.
- [28] N. Cardiel, Histogram comparison via Numerical Simulation, Proceedings of Astronomical Data Analysis Software (ADASS XXIV) Conference, eds: A.R. Taylor and E. Rosolowsky, San Francisco: Astronomical Society of the Pacific. 2015, p. 335.
- [29] S. Bityukov, N. Krasnikov, A. Nikitenko, V. Smirnova, Two approaches to Combining Significances, Proceedings of Science (ACAT08), Vol.118, 2008, pp. 1-12.