

An Unbinned Goodness-of-Fit Test Based on the Random Walk

K. Kinoshita
University of Cincinnati, Cincinnati, OH 45221 USA

We describe a test statistic for unbinned goodness-of-fit of data in one dimension. The statistic is based on the two-dimensional Random Walk. The rejection power of this test is explored both for simple and compound hypotheses and, for the examples explored, it is found to be comparable to that for the χ^2 test. We discuss briefly how it may be possible to extend this test to multi-dimensional data.

1. INTRODUCTION

This search for an unbinned goodness-of-fit test has been motivated by the widespread use of unbinned maximum likelihood fitting for determining CP -violating parameters at Belle. While there are many cross-checks to insure that there are no spurious signals and biases, the fits tend to be complicated and not very transparent. They often involve probability density functions (PDF's) that differ with every event, based on measured quantities that add dimensions to the data that are not explicit in the fits. As there is no widely accepted unbinned goodness-of-fit test that applies to such fits, testing for statistical consistency of results has been uneven. The tests that have been done, resorting to binned χ^2 or toy Monte Carlo, have their place but have not been entirely satisfactory in addressing the question.

A common technique of unbinned tests involves first transforming the measured quantities to a variable in which the null hypothesis has a uniform distribution, where the PDF is flat, and then to test this "flattened" distribution for consistency with uniformity. There exists a variety of tests for uniformity, but most are not readily extended to multidimensional data, and they do not address compound hypotheses. A review of methods is given in [1].

In this report, we explore a test statistic that is based on the two-dimensional Random Walk. To begin, its distribution in the case of a flat PDF is discussed. The ensemble distribution is then found for several alternate hypotheses, and the rejection power is calculated for comparison with other goodness-of-fit tests. As the aim of a goodness-of-fit test as it would be applied at Belle is to test the validity of the parametrization used in fitting, it is also important to examine how the test is modified under compound hypotheses. The discussion is thus expanded to include data which are fitted to determine one or more parameters. Finally, we discuss the possibility of extending to multidimensional data.

2. RANDOM WALK AS A TEST OF FLATNESS

A data set consisting of N measurements of the one-dimensional quantity x lying in the interval $[0, 1]$ may be mapped trivially to points on a unit circle with polar angle ϕ on the interval $[0, 2\pi]$, so that each point is considered to be a unit vector with direction defined by ϕ . If the PDF in x is flat, the vector sum of the corresponding unit vectors in two dimensions corresponds to the net displacement, D , after a two-dimensional Random Walk of N steps with unit step size. For sufficiently large N , this distribution converges to a well-known form (Rayleigh, 1888) and the distribution in D^2 is an exponential decay with mean equal to N . We take D^2/N as the test statistic. A deviation of the root distribution from the hypothesis will result in a bias of the ensemble distribution of this test statistic away from the origin. This statistic is mathematically equivalent to the first order term in the Fourier series that describes the distribution of the data:

$$\begin{aligned} \mathcal{F}(k=1) &= \int_0^{2\pi} d\phi \sum_{j=1}^N e^{ik\phi} \delta(\phi - \phi_j) \quad (1) \\ &= \sum_{j=1}^N e^{i\phi_j} \end{aligned}$$

where one can see that $D^2 \propto |\mathcal{F}(1)|^2$. One would expect this distribution to be most sensitive to an overall imbalance of the PDF in generally opposite ϕ directions. To obtain sensitivity to higher order differences, one could thus take successively higher order terms in the series, for $k = 2, \dots$. In practice it may not be useful to examine terms above $k = 3$. In this study we look at $k = 1$ ($d = 1$) and define $K_k \equiv \frac{|\mathcal{F}(k)|^2}{N}$. What we have defined as K_1 appears in the review of D'Agostino and Stephens[2] as R in the context of the Von Mises test, a test for uniformity on a circle.

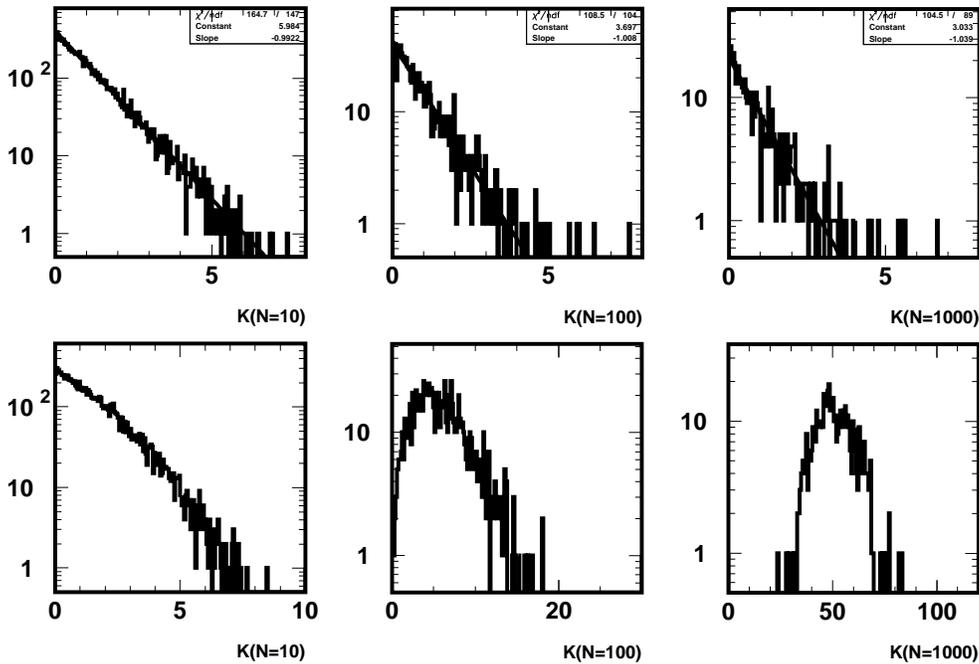


Figure 1: (top row) Distributions in K_1 for flat PDF: experiments with $N = 10$, $N = 100$, and $N = 1000$, shown with fits to an exponential form. (bottom row) Distributions in K_1 for PDF with the form $0.3 + 1.4X$ with $N = 10$, $N = 100$, and $N = 1000$.

Table I Rejection power for functions \mathcal{A}_1 , \mathcal{A}_2 , and \mathcal{A}_3 with a flat null hypothesis.

Function	Rejection Power		
	$N = 10$	$N = 100$	$N = 1000$
\mathcal{A}_1 (Linear)	0.117	0.824	1.00
\mathcal{A}_2 (Wide Gaussian)	0.152	0.910	1.00
\mathcal{A}_3 (Narrow Gaussian)	0.102	0.672	1.00

3. FLAT PDF

As mentioned above, the K_1 distribution for a flat PDF converges rapidly to an exponential with a decay constant of unity. Figure 1 (top row) shows the distributions in K_1 for ensembles of randomly generated experiments containing $N = 10$, 100, and 1000 events. Each of the three distributions is fitted via binned maximum likelihood to an exponential form. The fitted inverse decay constants (“slopes”) are 0.992 ± 0.010 , 1.008 ± 0.033 , and 1.039 ± 0.049 , respectively, in excellent agreement with the expectation.

To evaluate rejection power, these distributions may be compared with those obtained for PDF’s that are not flat. The alternative hypotheses used in a study

by Aslan and Zech [3] provide a convenient range of function types and allow for a direct comparison with the range of tests reviewed in their work. In that paper the rejection power of the alternative hypothesis was defined as one minus the probability for an error of the second kind, given a criterion that yields a 5% significance for the null hypothesis. Since in this case the null hypothesis gives an exponential distribution with unit decay constant, the 5% criterion is $K_1 > 3.0$. Ensembles of experiments were generated for each of three functions used in Ref. [3]:

$$\mathcal{A}_1(X) = 0.3 + 1.4X \quad (2)$$

$$\mathcal{A}_2(X) = 0.7 + 0.3[n_2 e^{-64(X-0.5)^2}] \quad (3)$$

$$\mathcal{A}_3(X) = 0.8 + 0.2[n_3 e^{-256(X-0.5)^2}] \quad (4)$$

where the n_i are normalization constants for the associated Gaussians. All functions are defined in the interval $[0, 1]$. The resulting K_1 distributions for \mathcal{A}_1 are shown in Figure 1 (bottom row). The values for rejection power are summarized in Table I. For comparison, the values for the χ^2 method ($N = 100$) given by Ref. [3] are approximately 0.81, 0.85, and 0.81, respectively, so our method is comparable in power, at least in the case of these three functions.

In order to apply this method as a goodness-of-fit test for non-uniform null hypotheses the PDF, $f(X)$, must first be transformed to a “flat” variable, Y ,

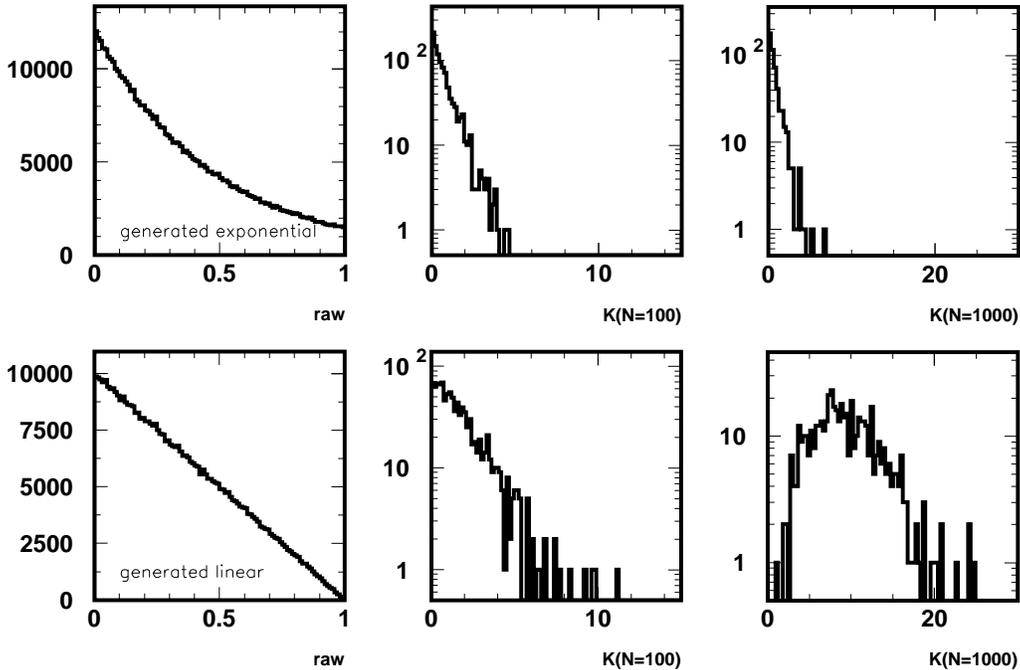


Figure 2: Determination of rejection power for a compound hypothesis: ensembles fitted for decay constant of exponentially decaying form. (top row) PDF matches fit parametrization: (left) Raw distribution, (center, right) distributions in K_1 of fitted, flattened experiments, $N = 100$ and $N = 1000$. (bottom row) PDF inconsistent with parametrization: (left) Raw distribution, (center, right) distributions in K_1 of fitted, flattened experiments, $N = 100$ and $N = 1000$.

where the probability distribution is flat. To form a uniform null hypothesis on a circle one could, for example, construct Y as :

$$Y_i = 2\pi \int_{X_-}^{X_i} f(X) dX \quad (5)$$

where the integer subscript i denotes the i^{th} data point and X_- is the lowest possible value of X .

4. COMPOUND HYPOTHESES

The examples considered thus far have been ones where no parameter fitting has occurred. While this has been an instructive exercise, it has limited application, as most measurements in particle physics involve the fitting of measured distributions to determine shapes and to derive some physics quantity or conclusion. We now look at compound hypotheses.

In evaluating rejection of alternative hypotheses via toy MC in the compound case, it is important that the fitting process be integrated into the evaluation procedure. Consider a data set $\{\phi_i\}$ where the PDF is assumed to be parametrizable as $f(\phi; \alpha)$ and the unbinned likelihood is maximum for $\alpha = \alpha_{max}$. The data are then flattened assuming the PDF is

$f(\phi; \alpha_{max})$, and the associated K_1 is evaluated. The confidence level of this K_1 value may then be found by referencing the ensemble distribution of K_1 when the true PDF is $f(\phi; \alpha_{max})$, and each experiment of the ensemble is treated as data, fitted and flattened according to the fit.

This procedure was used to evaluate rejection power for pairs of similarly shaped PDF's. Here we show one such result, for the hypothesis $n_4(\alpha)e^{-10X\alpha}$, where n_4 is a normalization constant, the measured quantity is X , and experiments are fitted for α . The alternative PDF was the linear form $f(X) = 2(1 - X)$. Experiments were generated according to the alternative PDF (A), and each was fitted to the hypothesis. The mean maximum likelihood value of α was approximately 4.7. Ensembles (B) were generated according to the hypothesis, with $\alpha = 4.7$, and fitted in the same way. The 5% confidence criterion on K_1 for (B) and acceptance of this criterion for (A) were estimated by counting (Figure 2). The rejection powers were found to be 28% and 99% for $N = 100$ and $N = 1000$, respectively. For comparison we also calculated by the same procedure the rejection of the χ^2 test, using 20 bins in the interval $[0,1]$ and found powers of 13% and 100%, respectively.

We also examined the two-dimensional distribution of fitted α_{max} values *vs.* K_1 . Any dependence of the

Table II Inverse decay constants of K_1 distribution for several generated forms, flattened after fitting for parameter(s) $\{\alpha_i\}$. The n_i are normalization constants, which may depend on the parameters α_j . No entry is made for samples where low statistics resulted in best fits which were at the limits of the parametrization.

Form	Generated Fitted		K_1 (Decay Constant) $^{-1}$ (χ^2/ndf)		
			$N = 10$	$N = 100$	$N = 1000$
$(1 - \alpha) + \alpha(2X)$	$\alpha = 0.7$	α	–	–	1.28 ± 0.07 (70/67)
$(1 - \alpha) + \alpha[n_2 e^{-64(X-0.5)^2}]$	$\alpha = 0.3$	α	–	1.90 ± 0.06 (230/80)	1.94 ± 0.09 (223/65)
$(1 - \alpha) + \alpha[n_3 e^{-256(X-0.5)^2}]$	$\alpha = 0.2$	α	–	1.56 ± 0.05 (203/82)	1.56 ± 0.07 (82/68)
$n_4 e^{-10X/\alpha}$	$\alpha = 1.0$	α	1.23 ± 0.01 (147/133)	1.28 ± 0.04 (68/85)	1.28 ± 0.06 (75/76)
$n_5 e^{-[X-(0.5+\alpha_2)]^2/2(\alpha_1/8)^2}$	$\alpha_1 = 1.0,$	α_1	1.36 ± 0.01 (176/131)	1.38 ± 0.05 (93/85)	1.50 ± 0.07 (56/65)
	$\alpha_2 = 0$	α_2	1.22 ± 0.01 (154/135)	1.25 ± 0.04 (122/96)	1.28 ± 0.06 (73/72)
		α_1, α_2	1.84 ± 0.019 (148/90)	2.00 ± 0.065 (53/59)	2.13 ± 0.095 (47/47)

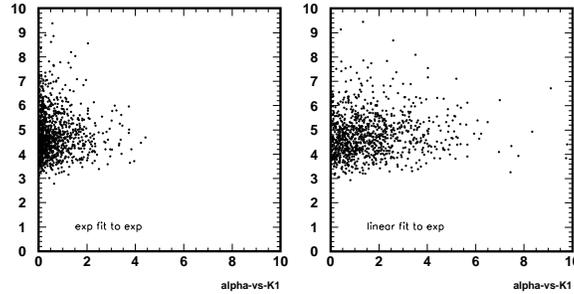


Figure 3: Scatter plots of fitted parameter α_{max} vs. K_1 for ensembles shown in Figure 2 ($N = 100$).

test on the fitted rather than underlying parameter value reduces its utility as a goodness-of-fit test; for example, the maximum likelihood value, \mathcal{L}_{max} , is not usable as a goodness-of-fit statistic because it depends strongly on the fitted parameter value(s) α_{max} – for a certain class of fitting functions, the correlation is 100%[4]. Figure 3 shows scatter plots of α_{max} and K_1 , where the data were generated with $N = 100$ and the generated and fitted forms are those from the example of Figure 2. There appears to be no strong dependence.

In any determination of rejection power with a compound hypothesis, it is necessary to determine the distribution of K_1 for the correct hypothesis. It does not appear that there is a simple ansatz as in the case of binned least squares fitting, where the chisquare converges to a chisquare distribution with the number of degrees of freedom reduced by one unit for each linear fitted parameter. We study this question empirically by generating MC ensembles for a variety of shapes. Each ensemble was generated according to the fitted functional form with parameter value(s) fixed. Each experiment was fitted with parameter(s) floating, and

the K_1 value was obtained from the data flattened according to the best fit. The distribution of resultant K_1 values for each ensemble was fitted for the decay constant, assuming an exponentially decaying form. Ensembles with $N = 10$, $N = 100$, and $N = 1000$ were generated. The results are summarized in Table II. There are several notable features. First, while all of the K_1 distributions had a decaying form, as one might expect, and a fit that converged, not all yielded good fits; the exponential form is not preserved under compound hypotheses. Secondly, all inverse decay constants are greater than unity, indicating that the K_1 distribution moves toward zero with fitting. This is not surprising; fitting identifies for each experiment the shape that is “closest” to the data, giving in general a better goodness-of-fit than the generator shape. Finally, there is no obvious pattern in the value of the decay constant with number of floated parameters. However, it is seen that for a given PDF and set of fitted parameters, the shape of the K_1 distribution shows remarkably little change as N is changed by two orders of magnitude.

5. EXTENSION TO MULTIDIMENSIONAL DATA: SPECULATION

Our goal in this investigation has been to arrive at a multidimensional unbinned goodness-of-fit test, one that has rejection power in all dimensions, not just in one-dimensional projections, for multidimensional data. Many unbinned tests depend on the integrated sum of or spacings between neighboring data points, quantities which are not well-defined when extended to more than one dimension. Although the K_1 statistic does not have this property, it is yet to be determined whether there exists an extension that is fully multidimensional; for example, in two dimensions, two components each mapped to a circle corresponds to a data space that is the surface of a toroid, for which there is no obvious nontrivial vector sum that maps to the Random Walk. A fully general extension to multidimensional data will additionally require a flattening algorithm and provisions for data spaces of arbitrary shape. We will continue to explore the possibilities for extending K_1 for use with multidimensional data.

6. SUMMARY

We have explored an unbinned goodness-of-fit test for data in one dimension that is based on the mapping of flattened distributions to a two-dimensional random walk. This method is truly binning-free and scale-independent, and the ensemble distribution for the null hypothesis is well-defined. For a compound hypothesis we specify a procedure to determine the ensemble distribution of the test statistic via Monte Carlo so that rejection power may be readily determined. The distribution is found for several different

parametrized forms and shown to be largely independent of statistics. We examine several samples for dependence between the test statistic and fitted parameter values, and find no evidence of any. The rejection power for alternate hypotheses is demonstrated for a few examples and is found to be comparable to that of the chisquare method.

Acknowledgments

The author would like to thank R. Cousins, G. Zech, and B. Yabsley for useful discussions and suggestions, and the organizers of PHYSTAT2003 for a stimulating and interesting conference. This work is supported by Department of Energy grant #DE-FG02-84ER40153.

References

- [1] B. Aslan and G. Zech, in *Proc. Conf. on Advanced Statistical Techniques in Particle Physics*, M.R. Whalley and L.Lyons, eds. (2002).
<http://www.ipp.dur.ac.uk/Workshops/02/statistics/proceedings.shtml>
- [2] *Goodness-Of-Fit Techniques (Statistics: Textbooks and Monographs Series, Vol. 68)*, R.B. D'Agostino and M.A. Stephens, eds., Marcel Dekker, Inc (1986).
- [3] B. Aslan and G. Zech, hep-ex/0203010 (2002).
- [4] J. Heinrich, "Can the Likelihood Function Value Be Used to Measure Goodness-of-Fit?" /CDF/MEMO/BOTTOM/CDFR/5639 (unpublished); K. Kinoshita, in *Proc. Conf. on Advanced Statistical Techniques in Particle Physics*, M.R. Whalley and L.Lyons, eds. (2002).