



Munich Personal RePEc Archive

## **A modified Kolmogorov-Smirnov test for normality**

Zvi Drezner and Ofir Turel and Dawit Zerom

California State University-Fullerton

22. October 2008

Online at <http://mpra.ub.uni-muenchen.de/14385/>

MPRA Paper No. 14385, posted 1. April 2009 04:38 UTC

# A Modified Kolmogorov-Smirnov Test for Normality

Zvi Drezner, Ofir Turel and Dawit Zerom\*  
Steven G. Mihaylo College of Business and Economics  
California State University-Fullerton  
Fullerton, CA 92834.

## Abstract

In this paper we propose an improvement of the Kolmogorov-Smirnov test for normality. In the current implementation of the Kolmogorov-Smirnov test, a sample is compared with a normal distribution where the sample mean and the sample variance are used as parameters of the distribution. We propose to select the mean and variance of the normal distribution that provide the closest fit to the data. This is like shifting and stretching the reference normal distribution so that it fits the data in the best possible way. If this shifting and stretching does not lead to an acceptable fit, the data is probably not normal. We also introduce a fast easily implementable algorithm for the proposed test. A study of the power of the proposed test indicates that the test is able to discriminate between the normal distribution and distributions such as uniform, bi-modal, beta, exponential and log-normal that are different in shape, but has a relatively lower power against the student  $t$ -distribution that is similar in shape to the normal distribution. In model settings, the former distinction is typically more important to make than the latter distinction. We demonstrate the practical significance of the proposed test with several simulated examples.

*Keywords:* Closest fit; Kolmogorov-Smirnov; Normal distribution.

## 1 Introduction

Many data analysis methods depend on the assumption that data were sampled from a normal distribution or at least from a distribution which is sufficiently close to a normal distribution. For example, one often tests normality of residuals after fitting a linear model to the data in order to ensure the normality assumption of the model is satisfied. Such an assumption is of great importance because, in many cases, it determines the method that ought to be used to estimate the unknown parameters in the model and also dictates the test procedures which the analyst may

---

\*Corresponding author: Mihaylo College of Business and Economics, California State University, Fullerton, CA, 92834-6848, (714) 278-3635, dzerom@fullerton.edu.

apply. There are several tests available to determine if a sample comes from a normally distributed population. Those theory-driven tests include the Kolmogorov-Smirnov test, Anderson-Darling test, Cramer-von Mises test, Shapiro-Wilk test and Shapiro-Francia test. The first three tests are based on the empirical cumulative distribution. Shapiro-Francia test (Shapiro and Francia, 1972 and Royston, 1983) is specifically designed for testing normality and is a modification of the more general Shapiro-Wilk test (Shapiro and Wilk 1965). There are also tests that exploit the shape of the distribution of the data. For example, the widely available Jarque-Bera test (Jarque and Bera, 1980) is based on skewness and kurtosis of the data. To complement the results of formal tests, graphical methods (such as box-plots and Q-Q plots) have also been used and increasingly so in recent years.

In this paper we focus on the Kolmogorov-Smirnov (KS) test. The KS test is arguably the most well-known test for normality. It is also available in most widely used statistical software packages. In its original form, the KS test is used to decide if a sample comes from a population with a completely specified continuous distribution. In practice, however, we often need to estimate one or more of the parameters of the hypothesized distribution (say, the normal distribution) from the sample, in which case the critical values of the KS test may no longer be valid. For the case of normality testing, Massey (1951) suggests using sample mean and sample variance, and this is the norm in the current use of KS test. Lilliefors (1967) and Dallal and Wilkinson (1986) provide a table of approximate critical values for use with the KS statistics when using sample mean and sample variance.

While the use of sample mean and sample variance seems a natural choice, using these fixed values is not necessarily the best available option. When one concludes (after using the KS test) that a sample is not normal, this only means that the data is not normal at the specified sample mean

and sample variance. But it could well be that the data is normal or sufficiently close to normal at other values of the mean and variance of the normal distribution. Although the scope of this paper is limited to the KS test, this drawback is also shared by other tests such as Anderson-Darling and Cramer-von Mises tests. Interestingly, Stephens (1974) writes after comparing several tests (such as KS, Anderson-Darling and so on) “It appears that since one is trying, in effect, to fit a density of a certain shape to the data, the precise location and scale is relatively unimportant, and being tied down to fixed values, even correct ones, is more of a hinderance than a help.” In this paper, we suggest an approach that circumvents the need to use pre-determined values of mean and variance. Instead, we look for mean and variance values such that the resulting normal distribution fits the given sample data. When such values do not exist, we conclude that the sample data is probably not normally distributed. Avoiding the use of fixed parameters, we propose a modified KS test in which we choose data-driven mean and variance values of the normal distribution by minimizing the KS statistics. In the traditional KS test, the data is compared against a normal distribution with fixed parameter values. On the other hand, our approach looks for a normal distribution that fits the data in the best possible way, and hence favors the sample data when passing judgment about its closeness to a normal distribution.

Suppose that the sample consists of  $n$  independent observations. These observations are sorted  $x_1 \leq x_2 \leq \dots \leq x_n$ . The cumulative distribution of the data is a step function (see Figures 1 and 2). At each  $x_k$  the step is between  $\frac{k-1}{n}$  and  $\frac{k}{n}$ . For a given mean  $\mu$  and variance  $\sigma^2$ , the cumulative normal distribution at  $x_k$  is  $\Phi\left(\frac{x_k - \mu}{\sigma}\right)$ . The KS statistics is given by

$$KS(\mu, \sigma) = \max_{1 \leq k \leq n} \left\{ \frac{k}{n} - \Phi\left(\frac{x_k - \mu}{\sigma}\right), \Phi\left(\frac{x_k - \mu}{\sigma}\right) - \frac{k-1}{n} \right\}. \quad (1)$$

The traditional KS statistics is simply  $KS(\bar{x}, s)$  where  $\mu = \bar{x}$  and  $\sigma = s$ . We propose a modified KS

statistics denoted by  $KS(\tilde{\mu}, \tilde{\sigma})$  where the vector  $(\tilde{\mu}, \tilde{\sigma})$  is a solution to the following minimization problem

$$\min_{\mu, \sigma} \{KS(\mu, \sigma)\} \quad (2)$$

where  $KS(\mu, \sigma)$  is as defined in (1). In section 2, we analyze this optimization problem and provide a tractable algorithm for its solution. In section 3, we provide critical values for the modified KS test using 100 million replications. The proposed algorithm is quite efficient and we are able to complete the critical values table (Table 1) in less than 4 days (6000 calculations per second). To facilitate implementation of our test, we also provide approximation formulas (that work for any  $n \geq 20$ ) for finding critical values at typical significant levels.

To best of our knowledge, there has not been any study that extends the KS test by allowing the use of optimized distribution parameters. Closely related to our work is that of Weber *et al* (2006) where they consider the problem of parameter estimation of continuous distributions (not just normal distribution) via minimizing the KS statistics. They use the heuristic optimization algorithm of Sobieszczanski-Sobieski *et al* (1998) to estimate the parameters of a number of widely used distributions and also provide a user-friendly software tool. The practical advantage of this software is that it suggests a best fitted distribution to given data by looking at the minimized KS statistics values among a set of continuous distributions. In this sense, our algorithm of minimizing the KS statistics may also serve the same purpose as that of Weber *et al* (2006) although our paper is wider in scope.

To motivate our modified KS test, we give two Monte Carlo based examples that can highlight the weaknesses of the existing KS and offer interesting practical implications for proper use of the KS

test.

**Example 1:** We generate 999 standard normal random samples of size  $n = 30$ . The choice of 999 samples (instead of say, 1000) is only to facilitate the calculation of the median sample as we will see below. For each sample, we calculate the two KS statistics values,  $KS(\bar{x}, s)$  and  $KS(\tilde{\mu}, \tilde{\sigma})$ , where the algorithm in section 2 is used to compute  $\tilde{\mu}$  and  $\tilde{\sigma}$ . We also compute  $\Delta = KS(\bar{x}, s) - KS(\tilde{\mu}, \tilde{\sigma})$  which is simply the difference between the two KS statistics values. It should be noted that  $KS(\tilde{\mu}, \tilde{\sigma}) \leq KS(\bar{x}, s)$  and hence  $\Delta \geq 0$ . We do the above steps for all 999 samples. Let  $\Delta_j$  denote a value obtained for sample  $j$  where  $j = 1 \dots, 999$ . We select a “typical” sample, say the  $k$ -th sample, to be the one where  $\Delta_k = Median\{\Delta_j\}_{j=1}^{999}$ . Similarly, an “extreme” sample, say the  $\ell$ -th sample, to be the one where  $\Delta_\ell = Max\{\Delta_j\}_{j=1}^{999}$

Based on the typical sample (sample  $k$ ), Figure 1 gives the empirical cumulative distribution (the step-function), the cumulative normal distribution (the dotted line) based on the sample mean ( $\bar{x}_k = 0.1078$ ) and sample variance ( $s_k = 1.022$ ) and the cumulative normal distribution (the solid line) based on  $\tilde{\mu}_k = 0.1712$  and  $\tilde{\sigma}_k = 1.089$ . The subscript  $k$  is attached to estimates to indicate that they correspond to the typical sample  $k$ . For this typical sample,  $KS(\bar{x}_k, s_k) = 0.0954$  and  $KS(\tilde{\mu}_k, \tilde{\sigma}_k) = 0.0704$  which indicate a 26% improvement by the latter. Note from the empirical cumulative distribution plots that the solid line is closer overall to the sample cumulative distribution. Using critical values Table 1 (for  $n = 30$ ), both KS statistics values lead to the non-rejection of the null of normality with p-value  $p > 0.2$ . This conclusion is correct as we know the sample is generated from a normal distribution.

Based on the extreme sample (sample  $\ell$ ), Figure 2 gives the empirical cumulative distribution (the step-function), the cumulative normal distribution (the dotted line) based on the sample mean ( $\bar{x}_\ell = -0.1628$ ) and sample variance ( $s_\ell = 0.9303$ ) and the cumulative normal distribution

(the solid line) based on  $\tilde{\mu}_\ell = -0.3238$  and  $\tilde{\sigma}_\ell = 0.7436$ . For this sample,  $KS(\bar{x}_\ell, s_\ell) = 0.1896$  and  $KS(\tilde{\mu}_\ell, \tilde{\sigma}_\ell) = 0.0951$  which indicate a 50% improvement by the latter. From the empirical cumulative distribution plots, the solid line is much closer to the sample cumulative distribution for data values roughly below -0.5 and these values constitute approximately 80% of the data observations. Using the critical values table for  $n = 30$  (Table 1), the traditional KS test implies that the sample data deviates from normality (at p-value  $p < 0.01$ ). On the other hand, the modified KS test concludes that we can not reject the null of normality at a convincing p-value  $p > 0.2$ . The conclusion from our test proposal is correct as the sample is generated from a normal distribution.

This example illustrates that the sample mean and sample variance do not necessarily provide the closest fit to the empirical distribution of the sample. Our approach shifts and stretches the normal distribution (by looking for data-driven mean and variance values) so that it fits the sample data in the best possible way.

**Example 2:** We consider  $n = 20, 40, \dots, 400$  (in an interval of 20). For each  $n$ , we generate 10,000 standard normal random samples of  $n - 1$  and one outlier. We define an outlier as  $outlier = C$  where the constant  $C$  takes values 4, 5,  $\dots$ , 10. We will only report results for  $C = 4, 6, 8, 10$  as the implications from the other outliers are qualitatively similar.

The purpose of this example is to evaluate the two tests: the traditional KS test (which is based on  $KS(\bar{x}, s)$ ) and the modified KS test (which is based on  $KS(\tilde{\mu}, \tilde{\sigma})$ ), in terms of their size using the level of significance  $\alpha = 0.05$ . When implementing both tests, we use the approximation formula in Table 2 for locating the critical values. Using 10,000 replications, we plot the size of the two tests for each  $n$  in Figure 3. Size is defined as the percentage of times (out of the total 10,000 samples) a test rejects the null hypothesis of normality. If a test is correctly sized, this percentage should be

close to 0.05. The dotted line in the figure corresponds to the size of the modified KS test while the solid lines correspond to the traditional KS test. Interestingly, the modified KS test is always close to 0.05 regardless of the magnitude of the outlier for all  $n$  (the average size from all  $n$  is 0.0508 with standard deviation of 0.0005). However the traditional KS test is very sensitive to outliers leading to clearly wrong conclusions about the distribution of the data. While increasing the sample size  $n$  seems to help minimize the effect of an outlier on the test, we still need unrealistically large sample sizes to get rid off the effect.

This example is only meant to illustrate the danger of using fixed parameter values that do not respond to the structure of sample data. The modified KS test adapts to the data by attempting (via choice of  $\tilde{\mu}$  and  $\tilde{\sigma}$ ) to fit the normal distribution to the majority of the data by weighting down the outlier. In practice, researchers often deal with small data sets with potentially a few outliers. Even if much of the data may be well approximated by a normal distribution, a blind use of traditional KS test will lead to rejection of normality - suggesting use of transformations or complex models. In contrast, the modified KS test is robust to these few outliers and can lead to more nuanced judgments regarding the normality of the data.

## 2 Algorithm

In this section, we analyze the optimization problem given in equation (2) and provide a tractable algorithm for its solution. By (1)

$$KS(\mu, \sigma) \geq \frac{k}{n} - \Phi\left(\frac{x_k - \mu}{\sigma}\right)$$

$$KS(\mu, \sigma) \geq \Phi\left(\frac{x_k - \mu}{\sigma}\right) - \frac{k-1}{n}$$



Let  $L$  be the minimum possible value of  $KS(\mu, \sigma)$ . The solution to the following optimization problem is the minimum possible  $KS(\mu, \sigma)$  and thus is equivalent to (2).

$$\min\{ L \} \tag{3}$$

subject to:

$$\frac{k}{n} - \Phi\left(\frac{x_k - \mu}{\sigma}\right) \leq L \quad \text{for } k > nL \tag{4}$$

$$\Phi\left(\frac{x_k - \mu}{\sigma}\right) - \frac{k-1}{n} \leq L \quad \text{for } k < n(1-L) + 1. \tag{5}$$

Note that if  $\frac{k}{n} - L \leq 0$ , constraint (4) is always true and if  $L + \frac{k-1}{n} \geq 1$ , constraint (5) is always true. We can solve (3-5) by designing an algorithm that finds whether there is a feasible solution to (4-5) for a given  $L$ .

For a given  $L$ , the constraints are equivalent to:

$$\mu \leq x_k - \Phi^{-1}\left(\frac{k}{n} - L\right)\sigma \quad \text{for } k > nL \tag{6}$$

$$\mu \geq x_k - \Phi^{-1}\left(L + \frac{k-1}{n}\right)\sigma \quad \text{for } k < n(1-L) + 1. \tag{7}$$

Constraints (6) and (7) can be combined into one constraint each.

$$\mu \leq \min_{k > nL} \left\{ x_k - \Phi^{-1}\left(\frac{k}{n} - L\right)\sigma \right\} \tag{8}$$

$$\mu \geq \max_{k < n(1-L)+1} \left\{ x_k - \Phi^{-1}\left(L + \frac{k-1}{n}\right)\sigma \right\} \tag{9}$$

For a given  $\sigma$  there is a solution for  $\mu$  satisfying the system of equations (8-9) if and only if

$$\min_{k > nL} \left\{ x_k - \Phi^{-1}\left(\frac{k}{n} - L\right)\sigma \right\} \geq \max_{k < n(1-L)+1} \left\{ x_k - \Phi^{-1}\left(L + \frac{k-1}{n}\right)\sigma \right\} \tag{10}$$

or

$$F(\sigma, L) = \min_{k > nL} \left\{ x_k - \Phi^{-1} \left( \frac{k}{n} - L \right) \sigma \right\} - \max_{k < n(1-L)+1} \left\{ x_k - \Phi^{-1} \left( L + \frac{k-1}{n} \right) \sigma \right\} \geq 0. \quad (11)$$

For a given  $L$ , the function  $F(\sigma, L)$  is a piece-wise linear concave function in  $\sigma$  (see Figure 4). We prove that  $F(\sigma, L)$  is a concave function in  $\sigma$  for a given  $L$ .

**Theorem 1:** *The function  $F(\sigma, L)$  for a given  $L$  is concave in  $\sigma$ .*

**Proof:** All the functions in the braces of (11) are linear in  $\sigma$  and all the other values are constants for a given  $L$ . Furthermore, the minimum of linear functions is concave and the maximum of linear functions is convex. Therefore, the difference  $F(\sigma, L)$  is a concave function in  $\sigma$ .  $\square$

By Theorem 1, for a given  $L$ ,  $F(\sigma, L)$  has only one local maximum which is the global one. The maximum value of  $F(\sigma, L)$  for a given  $L$  can be easily found by a search on  $\sigma$ . For any value of  $\sigma$   $F(\sigma, L)$  can be calculated and if the slope is positive we know that the optimal  $\sigma$  is to the right, and if it is negative we know that it is to the left. The solution is always at the intersection point between two lines, one with a positive slope and one with a negative slope (see figure 4). Megiddo (1983) suggested a very efficient method for solving such a problem.

Note that if  $F(\sigma, L) \geq 0$ , any  $\mu$  in the range

$$\left[ \max_{k < n(1-L)+1} \left\{ x_k - \Phi^{-1} \left( L + \frac{k-1}{n} \right) \sigma \right\}, \min_{k > nL} \left\{ x_k - \Phi^{-1} \left( \frac{k}{n} - L \right) \sigma \right\} \right]$$

(or specifically the midpoint of the range) with the  $\sigma$  used in calculating  $F(\sigma, L)$  yields a KS statistic which does not exceed  $L$ .

Let  $G(L) = \max_{\sigma} \{F(\sigma, L)\}$  found by either the method in Megiddo (1983) or any other search

method. If  $G(L) \geq 0$ , there is a solution  $(\mu, \sigma)$  for this value of  $L$  and if  $G(L) < 0$  no such solution exists. To find the minimum value of  $L$  we propose a binary search. The details of the binary search are now described. The optimal  $L$  must satisfy  $L \leq KS(\bar{x}, s)$ . Also, any KS statistic must be at least  $\frac{1}{2n}$ . Therefore,  $\frac{1}{2n} \leq L \leq KS(\bar{x}, s)$ . A binary search on any segment  $[a, b]$  is performed as follows.  $G(L)$  for  $L = \frac{a+b}{2}$  is evaluated. If  $G(L) \geq 0$ , there is a solution  $(\mu, \sigma)$  for this value of  $L$  and the search segment is reduced to  $[a, \frac{a+b}{2}]$ . If  $G(L) < 0$  no such solution exists and the search segment is reduced to  $[\frac{a+b}{2}, b]$ . In either case the search segment is cut in half. Following a relatively small number of iterations, the search segment is reduced to a small enough range (such as  $10^{-5}$ ) and the upper limit of the range yields a solution  $(\mu, \sigma)$  and its value of  $L$  is within a given tolerance (the size of the final segment) of the optimal value of  $L$ .

### 3 Monte Carlo estimation of test statistics distribution

In this section we provide critical values for the modified KS statistics using Monte Carlo simulation. To derive the distribution of this statistics, we draw a random sample of size  $n$  from a standard normal distribution. We estimate  $\tilde{\mu}$  and  $\tilde{\sigma}$  and compute  $KS(\tilde{\mu}, \tilde{\sigma})$ , and for every sample size  $n$ , we repeat this procedure 100 million times. The critical values are given in Table 1. We also recalculate the critical values for the traditional KS test in the same way and are available in Table 1. Because we use 100 million samples, the critical values we report for the traditional KS test are more accurate than Lilliefors (1967) and Dallal and Washington (1986).

The critical values for both  $KS$  tests can be approximated for  $n \geq 20$  by the formula  $a + \frac{b}{\sqrt{n}} (1 - \frac{c}{n})$  where  $a$ ,  $b$  and  $c$  are functions of  $\alpha$ . These three parameters are given in Table 2. The approximation is very accurate with an error (when compared to Table 1) of not more than 0.0002. So, the approximation formula can replace the tables for  $n \geq 20$ . We obtain the approximation formula

via multiple regression, where for each  $\alpha$ , the critical values in Table 1 are used as the dependent variable, and  $\frac{1}{\sqrt{n}}$  and  $\frac{1}{n\sqrt{n}}$  are the independent variables. We select these two independent variables through experimentation. We begin with a single variable regression involving only  $\frac{1}{\sqrt{n}}$ . We then add variables, one at a time, which are functions of  $n$ . A regression involving  $\frac{1}{\sqrt{n}}$  and  $\frac{1}{n\sqrt{n}}$  provides an excellent fit.

## 4 Power comparisons

In this section we compare the approximate powers of the modified KS test with the traditional KS test for a set of selected distributions. These distributions convey a wide array of shapes where some resemble the normal distribution while others are substantially different. Some of these distributions are also used in Lilliefors (1967) and Stephens (1974), among others. We consider a uniform (0,1) distribution; a bi-modal distribution which is a composite of two normal distributions, one centered at +2 and one at -2 with variance of 1; a beta(1,2) distribution whose density function is a straight line connecting (0, 0) and (1, 1); an exponential distribution with mean and variance of 1; a log-normal distribution with mean  $e^{1/2}$  and variance  $e(e - 1)$  and three  $t$ -distributions with degrees of freedom 1, 2 and 6. We also include the normal distribution where we expect power to be close to  $\alpha$ . To save space, we only report results for  $\alpha = 0.05$  (the behavior is very similar for other values of  $\alpha$ ).

For a given alternative hypothesis (say, a uniform distribution), computation of the power of the modified KS test is done as follows. We draw a random sample of size  $n$  from the distribution specified in the alternative hypothesis. Based on this sample, we estimate the parameters  $\tilde{\mu}$  and  $\tilde{\sigma}$  using the algorithm outlined in section 2 and compute  $KS(\tilde{\mu}, \tilde{\sigma})$ . Then, apply the critical values in Table 2 to test if such sample comes from a normal distribution. Repeating this procedure

10,000 times, and counting the number of correct decisions gives the approximate power. The same approach is followed to compute power for traditional KS test. The complete power results are given in Table 3.

From Table 3 we can see that the power of the modified KS test is consistently better than the traditional KS test for uniform, beta and bi-modal distributions. The improvement is quite large especially for uniform and beta distributions. These power results indicate that the proposed KS test is able to better discriminate between the normal distribution and those distributions that are very different in shape from normal, i.e. those that substantially deviate from normality. For exponential and log-normal distributions, the powers of the two KS tests are quite similar where both achieve reasonably good powers for  $n \geq 40$ . For the  $t$ -distributions, the modified KS test has a much lower power than the traditional KS test. What is common to the  $t$ -distributions is that they resemble the normal distribution except for their heavier tails. In theory, with increasing degrees of freedom, the tails of the  $t$ -distribution get lighter eventually behaving like the normal distribution. The modified KS test has difficulty detecting non-normality when the observed distribution is similar to normal and increasingly so with larger degrees of freedom, i.e. as it gets closer to normal.

On the surface, the low power for the  $t$ -distribution may seem like a weakness of the modified KS test. However, would one expect, with a small  $n$ , that data generated by a  $t_6$  distribution be distinguishable from a normal distribution - thus be identified as non-normal? We argue that the reason the traditional KS test has a higher power is that it rejects data which can be fitted quite well to a normal distribution by a proper selection of  $\mu$  and  $\sigma$ . It is indeed strange that the power of the traditional KS test is higher for a  $t_2$  distribution than it is for the uniform and beta distributions while the latter are substantially different from normality. By construction, the modified KS test tries to look for those mean and variance values that lead to the closest fit to the

data. In a way, we are trying to approximate the reference distribution (the  $t$ -distribution) with a normal distribution. If such a normal approximation exists, the data may be considered sufficiently normal. For example, for  $t_6$ , the powers at  $n \leq 100$  are close to  $\alpha = 0.05$  implying the sample data is hardly distinguishable from the normal distribution (see how close the powers of  $t_6$  are to those of the normal distribution). When the degrees of freedom is made smaller, the power of the modified KS test improves because the deviation from normality gets larger. When normal approximation can not be achieved, the sample data is flagged as non-normal. For  $t_2$ , the modified KS test is able to detect difference from normality at  $n = 200$  while  $t_6$  requires a very large  $n$  to be detected by the modified KS. For  $t_1$ , the power of the proposed KS test gets a lot better reaching decent power at  $n = 100$ . The reason is that  $t_1$  has a much heavier tail than the normal distribution making normal approximation via data driven mean and variance values very difficult.

To see why the modified KS test treats several small data from the  $t$ -distribution as normally distributed, we use the  $t_2$ -distribution as an example. To do so, we repeat the experiments described in Example 1 (see section 1) but draw 999 samples (of  $n = 30$ ) from a  $t_2$  distribution. The odd number of simulation replications has the same purpose as in Example 1. We select a “typical” sample in terms of the difference between the traditional KS statistic and our proposed KS statistic. Similar to Figures 1 and 2, three cumulative distribution are depicted in Figure 5 (the range of  $x$  was truncated for better exposition). For this typical sample,  $\bar{x} = -0.0117$ ,  $s = 2.506$ ,  $\tilde{\mu} = 0.0321$  and  $\tilde{\sigma} = 1.571$ . The traditional KS is  $KS(\bar{x}, s) = 0.1805$  while the modified KS is  $KS(\tilde{\mu}, \tilde{\sigma}) = 0.1003$ . Using the critical value tables in section 3, the traditional KS test rejects the normality with a p-value of  $p = 0.05$ . On the contrary, the modified KS test does not reject normality with p-value  $p > 0.10$ .

## 5 Conclusion

Many data analysis methods (t-test, ANOVA, regression) depend on the assumption that data were sampled from a normal distribution. One of the most frequently used test to evaluate how far data are from normality is the Kolmogorov-Smirnov (KS) test. In implementing the KS test, most statistical software packages use the sample mean and sample variance as the parameters of the normal distribution. However, the sample mean and sample variance do not necessarily provide the closest fit to the empirical distribution of the data. Therefore, we propose a modified KS test in which we optimally choose the mean and variance of the normal distribution by minimizing the KS statistics. To facilitate easy implementation we also provide an algorithm to solve for the optimal parameters.

## References

1. Dallal G. E. and L. Wilkinson (1986) "An analytic approximation to the distribution of Lilliefors's test statistic for normality," *The American Statistician*, 40, 294-296.
2. Jarque, C.M. and A.K. Bera (1980) "Efficient Tests for Normality, Homoscedasticity and Serial Independence of Regression Residuals," *Economics Letters*, 6(3), 255-259.
3. Lilliefors H. W. (1967) "On the Kolmogorov-Smirnov test for normality with mean and variance unknown," *Journal of the American Statistical Association*, 62, 399-402.
4. Massey F. J. (1951) "The Kolmogorov-Smirnov test for goodness of fit," *Journal of the American Statistical Association*, 46, 68-78.
5. Megiddo N. (1983) "Linear-time algorithms for linear programming in  $R^3$  and related problems," *SIAM Journal on Computing*, 12, 759-776.
6. Royston, J. P. (1983) "A Simple Method for Evaluating the Shapiro-Francia W' Test of Non-Normality," *Statistician*, 32(3) (September), 297-300.
7. Shapiro, S. S. and R. S. Francia (1972) "An Approximate Analysis of Variance Test for Normality," *Journal of the American Statistical Association*, 67, 215-216.
8. Shapiro, S. S. and M. B. Wilk (1965) "An Analysis of Variance Test for Normality (Complete Samples)," *Biometrika*, 52(3/4) (December), 591-611.

9. Sobieszczanski-Sobieski, J., Laba, K. and R. Kincaid (1998) “*Bell-curve evolutionary optimization algorithm*,” Proceedings of the 7th AIAA Symposium on Multidisciplinary Analysis and Optimization, St. Louis, MO, 2-4 September, AIAA paper 98-4971.
10. Stephens, M.A. (1974) “EDF statistics for goodness of fit and some comparisons, ” *Journal of the American Statistical Association*, 69, 730-737.
11. Weber, M., Leemis, L. and R. Kincaid (2006) “Minimum Kolmogorov-Smirnov test statistic parameter estimates,” *Journal of Statistical Computation and Simulation*, 76, 3, 195-206.



Table 1: Critical Values for the Traditional and Modified KS Test

$n$	Traditional KS statistics						Modified KS statistics					
	Upper Tail Probabilities						Upper Tail Probabilities					
	0.20	0.15	0.10	0.05	0.01	0.001	0.20	0.15	0.10	0.05	0.01	0.001
4	0.3029	0.3215	0.3453	0.3753	0.4131	0.4327	0.2396	0.2436	0.2474	0.2499	0.2987	0.3518
5	0.2894	0.3027	0.3189	0.3430	0.3967	0.4388	0.2000	0.2108	0.2255	0.2458	0.2763	0.3063
6	0.2687	0.2809	0.2971	0.3234	0.3705	0.4232	0.1962	0.2046	0.2147	0.2286	0.2570	0.2945
7	0.2523	0.2643	0.2802	0.3042	0.3508	0.4011	0.1855	0.1922	0.2006	0.2139	0.2435	0.2708
8	0.2388	0.2503	0.2651	0.2880	0.3328	0.3827	0.1748	0.1810	0.1899	0.2038	0.2281	0.2502
9	0.2272	0.2381	0.2522	0.2741	0.3172	0.3657	0.1661	0.1727	0.1811	0.1932	0.2151	0.2418
10	0.2171	0.2274	0.2410	0.2621	0.3035	0.3509	0.1591	0.1650	0.1725	0.1836	0.2045	0.2324
11	0.2081	0.2181	0.2312	0.2514	0.2914	0.3375	0.1524	0.1578	0.1648	0.1753	0.1972	0.2240
12	0.2003	0.2099	0.2224	0.2420	0.2807	0.3255	0.1462	0.1514	0.1580	0.1681	0.1902	0.2158
13	0.1932	0.2025	0.2146	0.2335	0.2710	0.3146	0.1407	0.1457	0.1521	0.1627	0.1839	0.2081
14	0.1869	0.1958	0.2076	0.2259	0.2623	0.3048	0.1358	0.1406	0.1472	0.1576	0.1780	0.2012
15	0.1811	0.1898	0.2012	0.2189	0.2543	0.2958	0.1314	0.1363	0.1428	0.1528	0.1725	0.1949
16	0.1759	0.1843	0.1954	0.2126	0.2471	0.2875	0.1276	0.1325	0.1388	0.1485	0.1674	0.1893
17	0.1710	0.1793	0.1900	0.2068	0.2404	0.2800	0.1243	0.1290	0.1351	0.1445	0.1628	0.1845
18	0.1666	0.1746	0.1851	0.2015	0.2342	0.2729	0.1211	0.1257	0.1316	0.1407	0.1585	0.1799
19	0.1625	0.1703	0.1806	0.1965	0.2285	0.2663	0.1182	0.1226	0.1284	0.1372	0.1545	0.1756
20	0.1587	0.1663	0.1763	0.1919	0.2232	0.2603	0.1154	0.1198	0.1254	0.1339	0.1510	0.1716
25	0.1430	0.1499	0.1589	0.1730	0.2014	0.2351	0.1040	0.1079	0.1129	0.1207	0.1363	0.1547
30	0.1312	0.1376	0.1458	0.1588	0.1849	0.2161	0.0955	0.0990	0.1036	0.1108	0.1251	0.1422
40	0.1145	0.1200	0.1272	0.1385	0.1614	0.1889	0.0833	0.0864	0.0905	0.0967	0.1092	0.1242
50	0.1029	0.1078	0.1143	0.1245	0.1450	0.1699	0.0749	0.0777	0.0813	0.0869	0.0982	0.1116
60	0.0943	0.0988	0.1047	0.1140	0.1328	0.1556	0.0687	0.0712	0.0745	0.0797	0.0900	0.1023
70	0.0875	0.0917	0.0972	0.1058	0.1233	0.1445	0.0638	0.0661	0.0692	0.0740	0.0835	0.0950
80	0.0820	0.0859	0.0911	0.0992	0.1156	0.1355	0.0598	0.0620	0.0649	0.0694	0.0783	0.0891
90	0.0775	0.0812	0.0860	0.0937	0.1092	0.1279	0.0565	0.0586	0.0613	0.0655	0.0740	0.0841
100	0.0736	0.0771	0.0817	0.0890	0.1037	0.1216	0.0537	0.0557	0.0583	0.0623	0.0703	0.0799
400	0.0373	0.0390	0.0414	0.0450	0.0524	0.0615	0.0273	0.0283	0.0296	0.0316	0.0356	0.0405
900	0.0249	0.0261	0.0277	0.0301	0.0351	0.0411	0.0183	0.0190	0.0198	0.0212	0.0239	0.0271

Table 2: Coefficients for the approximate formulas

$\alpha$	Traditional KS test			Modified Ks test		
	$a$	$b$	$c$	$a$	$b$	$c$
0.20	0.00053	0.73574	0.78520	0.00060	0.53446	0.80443
0.15	0.00049	0.77149	0.78515	0.00068	0.55329	0.76285
0.10	0.00059	0.81689	0.77062	0.00062	0.57999	0.78034
0.05	0.00052	0.89105	0.79780	0.00061	0.62082	0.81183
0.01	0.00054	1.03964	0.84912	0.00055	0.70276	0.85751
0.001	0.00052	1.22182	0.99171	0.00056	0.79997	0.89234

Table 3: Powers (%) of the Traditional and Modified KS tests ( $\alpha = 0.05$ )

$n$	Uniform		Bi-modal		Beta	
	†	‡	†	‡	†	‡
20	9.59	16.28	35.11	44.48	17.73	25.09
40	19.30	29.14	70.57	76.65	36.20	51.99
60	32.15	44.66	90.40	92.48	54.74	75.78
80	46.30	60.53	97.34	97.89	69.53	89.83
100	58.58	74.36	99.43	99.65	81.45	96.57
200	94.53	99.20	100.00	100.00	99.66	99.99
300	99.81	100.00	100.00	100.00	100.00	100.00
400	99.99	100.00	100.00	100.00	100.00	100.00
$n$	Exponential		Log-normal		$t_1$	
	†	‡	†	‡	†	‡
20	58.38	57.06	79.88	67.88	84.86	30.46
40	90.49	91.94	98.27	95.97	98.16	56.05
60	98.66	99.16	99.94	99.68	99.82	76.60
80	99.87	99.96	99.99	99.98	99.99	88.70
100	100.00	100.00	100.00	100.00	100.00	95.16
200	100.00	100.00	100.00	100.00	100.00	99.98
300	100.00	100.00	100.00	100.00	100.00	100.00
400	100.00	100.00	100.00	100.00	100.00	100.00
$n$	$t_2$		$t_6$		Normal	
	†	‡	†	‡	†	‡
20	45.74	9.87	11.40	5.33	5.01	5.03
40	68.93	14.85	15.32	5.40	5.04	4.93
60	84.02	21.15	17.83	5.80	5.13	5.02
80	91.89	29.31	21.78	6.25	5.13	4.92
100	95.86	36.80	25.01	7.16	5.15	4.99
200	99.92	73.15	40.01	9.27	5.45	5.22
300	100.00	92.16	53.43	11.46	5.05	5.36
400	100.00	98.35	64.87	14.30	5.16	4.87

† Traditional KS test

‡ Modified KS test

Figure 1: The Typical Sample

### Cumulative Distributions

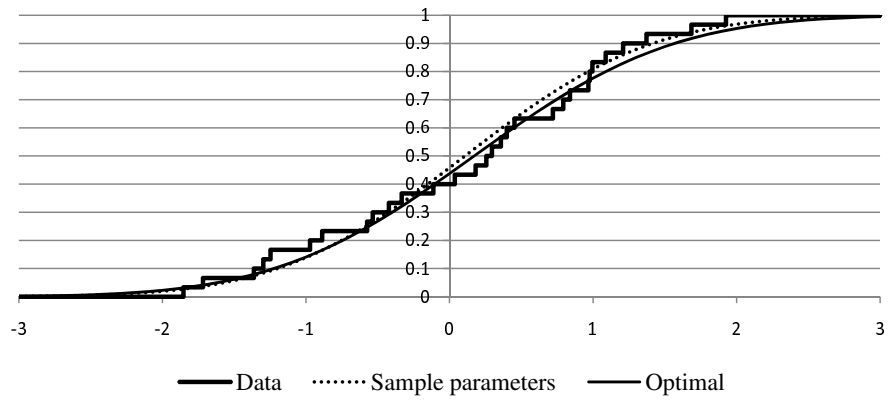


Figure 2: The Extreme Sample

### Cumulative Distributions

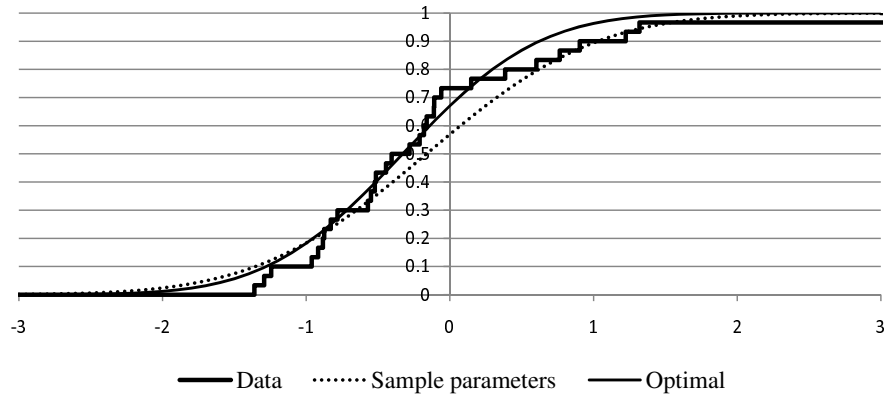


Figure 3: Sizes of the Traditional (solid line) and Modified (dotted line) KS tests ( $\alpha = 0.05$ )

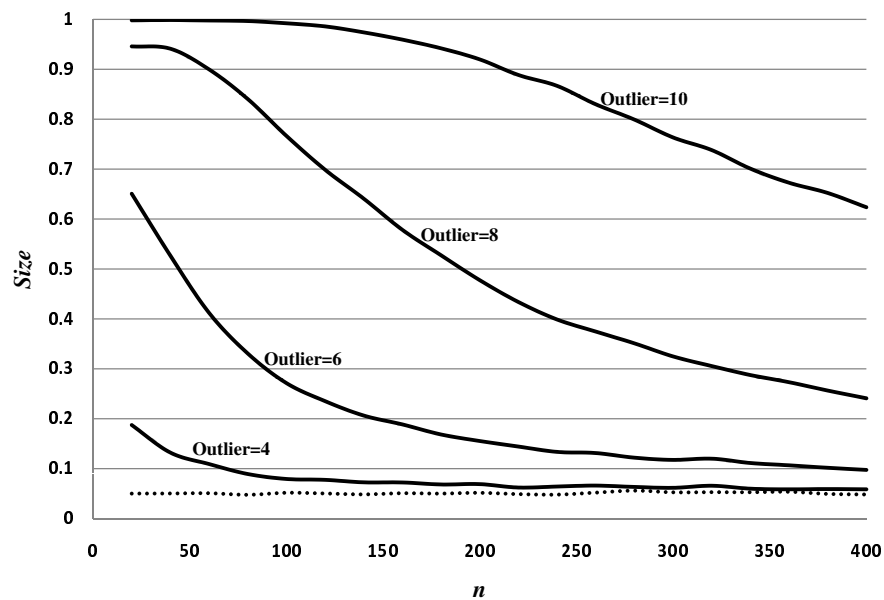


Figure 4: The Function  $F(\sigma, L)$

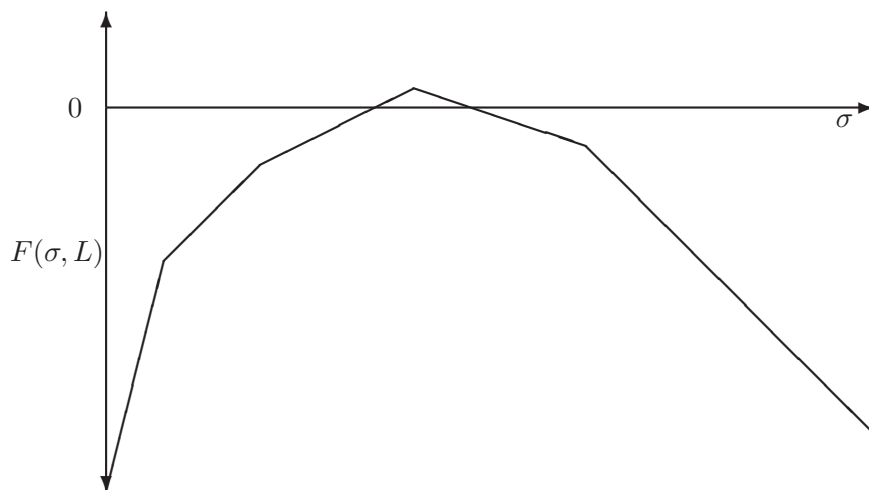


Figure 5: Typical  $t_2$  Samples

### Cumulative Distributions

