

Applications of Mathematics

Martin Schindler

Kolmogorov-Smirnov two-sample test based on regression rank scores

Applications of Mathematics, Vol. 53 (2008), No. 4, 297--304

Persistent URL: <http://dml.cz/dmlcz/140323>

Terms of use:

© Institute of Mathematics AS CR, 2008

Institute of Mathematics of the Academy of Sciences of the Czech Republic provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This paper has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://project.dml.cz>

KOLMOGOROV-SMIRNOV TWO-SAMPLE TEST BASED ON
REGRESSION RANK SCORES*

MARTIN SCHINDLER, Praha

(Supplement to the special issue of Appl. Math. 53 (2008), No. 3)

Abstract. We derive the two-sample Kolmogorov-Smirnov type test when a nuisance linear regression is present. The test is based on regression rank scores and provides a natural extension of the classical Kolmogorov-Smirnov test. Its asymptotic distributions under the hypothesis and the local alternatives coincide with those of the classical test.

Keywords: regression rank scores, Kolmogorov-Smirnov test, two sample problem, Cramér-von Mises test

MSC 2010: 62G08, 62G10, 62J05

1. INTRODUCTION

In [3] Hájek extended the Kolmogorov-Smirnov test of the hypothesis of randomness to tests against alternatives of simple linear regression. He expressed the test criterion (see equation (4)) as a functional of a special rank score process (Hájek's rank scores) for which he proved convergence to Brownian bridge. We mention this fact in Subsection 2.1. Similarly he extended the Cramér-von Mises and the Rényi tests. If, instead of Hájek's rank scores, we consider the process of regression rank scores (see e.g. [1]), we can extend the (two-sample) Kolmogorov-Smirnov test also to a nuisance regression.

So here we deal with the tests of Kolmogorov-Smirnov type on one component of the regression parameter β in the linear model $\mathbf{Y} = \mathbf{X}\beta + \mathbf{e}$. These tests, based on regression rank scores, were introduced in Jurečková [5]. We derive the two-sample variant of the test and show that this test represents a straightforward extension of the classical Kolmogorov-Smirnov test, more specifically the variant of the classical Kolmogorov-Smirnov test that is the most sensitive to difference in location.

* This work was supported by the Czech Science Foundation under Grant No. 201/05/H007 and by Research Project LC06024.

Note that already in Gutenbrunner and Jurečková [1] the regression rank score process was studied. Further, in Gutenbrunner et al. [2] a broader class of tests of hypothesis in linear regression model based on regression rank scores was derived. This class represents a generalization of simple linear rank tests.

Consider the linear regression model

$$(1) \quad \mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e} = \mathbf{X}^{(1)}\boldsymbol{\beta}^{(1)} + \mathbf{x}^{(p)}\beta_p + \mathbf{e},$$

where $\mathbf{Y} = (Y_1, \dots, Y_N)'$ is a vector of observations, $\mathbf{X} = \mathbf{X}_{N \times p}$ is a known design matrix, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)' = (\boldsymbol{\beta}^{(1)'}, \beta_p)$ are unknown parameters, $\mathbf{e} = (e_1, \dots, e_N)'$ is the vector of i.i.d. errors, the matrix $\mathbf{X}^{(1)}$ consisting of the first $p - 1$ columns of the matrix \mathbf{X} represents the nuisance regression and $\mathbf{x}^{(p)}$ is the p th column of \mathbf{X} . Here we do not specify the vector $\mathbf{x}^{(p)}$ but later, in Section 2, we will set $\mathbf{x}^{(p)} = (1, \dots, 1, 0, \dots, 0)'$ to derive and describe the two-sample Kolmogorov-Smirnov test.

We want to test the hypothesis

$$H_0: \beta_p = 0, \quad \boldsymbol{\beta}^{(1)} \text{ unspecified.}$$

This problem will be tested by a test of Kolmogorov-Smirnov (K-S) type. In the presence of nuisance regression, regression rank scores (RRS) are employed. RRS (see e.g. [2]) in the submodel of (1) given by H_0 are defined as the vector of solutions $\hat{\mathbf{a}}_N(\alpha) = (\hat{a}_{N1}(\alpha), \dots, \hat{a}_{NN}(\alpha))'$, $0 \leq \alpha \leq 1$ of the linear programming problem ($\mathbf{1}_N$ denotes the $(N \times 1)$ vector of ones):

$$\max \mathbf{Y}'\hat{\mathbf{a}}_N(\alpha)$$

subject to

$$(2) \quad \begin{aligned} \mathbf{X}^{(1)'}\hat{\mathbf{a}}_N(\alpha) &= (1 - \alpha)\mathbf{X}^{(1)'}\mathbf{1}_N, \\ \hat{\mathbf{a}}_N(\alpha) &\in [0, 1]^N. \end{aligned}$$

1.1. Assumptions

We will impose the following conditions on the regression matrix \mathbf{X} and on the underlying distribution function F .

Let \mathbf{x}'_i denote the i th row of the matrix \mathbf{X} , $i = 1, \dots, N$. We assume that the matrix $\mathbf{X} = \mathbf{X}_N$ satisfies the regularity conditions

$$(X.1) \quad x_{i1} = 1, \quad i = 1, \dots, N,$$

$$(X.2) \quad \max_{\substack{1 \leq i \leq N \\ 1 \leq j \leq p}} |x_{ij}| = \mathcal{O}(N^{(2(b-a)-\delta)/(1+4b)})$$

$$\text{for some } a, b, \delta, \quad 0 < a \leq \frac{1}{4} - \varepsilon, \quad 0 < b - a \leq \frac{1}{2}\varepsilon, \quad \varepsilon > 0, \quad \delta > 0,$$

$$(X.3) \quad \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i\|^3 = \mathcal{O}(1) \text{ as } N \rightarrow \infty,$$

$$(X.4) \quad \mathbf{D}_N = N^{-1} \mathbf{X}'_N \mathbf{X}_N \xrightarrow{N \rightarrow \infty} \mathbf{D}, \text{ where } \mathbf{D} \text{ is a positively definite matrix.}$$

Further assume that the errors e_1, \dots, e_N in (1) are i.i.d. with an absolutely continuous distribution function F whose tails are assumed to satisfy the following regularity conditions (F.1)–(F.4) (these conditions are satisfied by many common densities f including t -distributions with 5 or more d.f.):

(F.1) F has an absolutely continuous density f , positive for $A < x < B$ and decreasing monotonously when $x \rightarrow A+$, $x \rightarrow B-$, where $-\infty \leq A = \sup\{x: F(x) = 0\}$ and $+\infty \geq B = \inf\{x: F(x) = 1\}$. The derivative f' of f is bounded a.e.

(F.2) $|F^{-1}(\alpha)| \leq c(\alpha(1-\alpha))^{-a}$ (with a from (X.2)) for $0 < \alpha \leq \alpha_0$, $1 - \alpha_0 \leq \alpha < 1$ for some $0 < \alpha_0 \leq \frac{1}{2}$ and for some $c > 0$.

(F.3) $1/f(F^{-1}(\alpha)) \leq c(\alpha(1-\alpha))^{-1-a}$ for $0 < \alpha \leq \alpha_0$, $1 - \alpha_0 \leq \alpha < 1$, $c > 0$.

(F.4) $\left| \frac{f'(x)}{f(x)} \right| \leq c(|x| + 1)$, $x \in \mathbf{R}^1$, $c > 0$.

1.2. Statistic of K-S type

Consider the model (1) and define the projection matrix

$$\mathbf{H}^{(1)} = \mathbf{H}_N^{(1)} = (h_{ij}^{(1)})_{i=1, \dots, N}^{j=1, \dots, N} = \mathbf{X}_N^{(1)} (\mathbf{X}_N^{(1)'} \mathbf{X}_N^{(1)})^{-1} \mathbf{X}_N^{(1)'}$$

and $\mathbf{x}^* = (x_1^*, \dots, x_N^*)' = \mathbf{H}^{(1)} \mathbf{x}^{(p)}$ the projection of $\mathbf{x}^{(p)}$ into the space spanned by the columns of $\mathbf{X}_N^{(1)}$.

We define the process $\{S_N(t): 0 \leq t \leq 1\}$ on $C[0, 1]$:

$$S_N(t) = \left(\sum_{i=1}^N (x_i^{(p)} - x_i^*)^2 \right)^{-1/2} \sum_{i=1}^N (x_i^{(p)} - x_i^*) \hat{a}_{Ni}(t).$$

It is shown in [5] that under the conditions (X.1)–(X.4) and (F.1)–(F.4) it follows from [2, Theorem 3.2] that

$$(3) \quad \sup_{0 \leq t \leq 1} |S_N(t) - \tilde{S}_N(t)| \xrightarrow{p} O \text{ as } N \rightarrow \infty,$$

where

$$\tilde{S}_N(t) = \left(\sum_{i=1}^N (x_i^{(p)} - x_i^*)^2 \right)^{-1/2} \sum_{i=1}^N (x_i^{(p)} - x_i^*) I[e_i > F^{-1}(t)], \quad 0 \leq t \leq 1$$

and that $S_N(t)$ converges to the Brownian bridge in the uniform topology on $C[0, 1]$. In the next section we show how to construct a test based on this fact in the case of a two sample problem.

2. TWO-SAMPLE PROBLEM

Consider the model (1) and let $\mathbf{x}^{(p)} = (1, \dots, 1, 0, \dots, 0)'$ be the vector with m ones and n zeros, $m + n = N$.

We want to test the hypothesis H_0 of no difference between the samples. This two-sample problem will be tested by a test of Kolmogorov-Smirnov (K-S) type which is a generalization of the classical rank test of K-S type (the variant that is the most sensitive to difference in location) that works in the model (1) without nuisance regression ($\mathbf{X}^{(1)} = \mathbf{1}_N$).

2.1. Classical K-S two-sample test

In the location model (model (1) with $\mathbf{X}^{(1)} = \mathbf{1}_N$) the solution $\hat{\mathbf{a}}_N(\alpha)$ of (2) specializes to Hájek's rank scores $\mathbf{a}_N^*(\alpha) = (a_{N1}^*(\alpha), \dots, a_{NN}^*(\alpha))$ where

$$a_{Ni}^*(\alpha) = a_N^*(R_i, \alpha) = \begin{cases} 1, & 0 \leq \alpha \leq (R_i - 1)/N, \\ R_i - \alpha N, & (R_i - 1)/N < \alpha \leq R_i/N, \\ 0, & R_i/N < \alpha \leq 1, \end{cases}$$

where R_i is the rank of Y_i among Y_1, \dots, Y_N , $i = 1, \dots, N$. Hájek in [3] or Hájek & Šidák in [4] considered the process $T_N = \{T_N(t) : 0 \leq t \leq 1\}$,

$$(4) \quad T_N(t) = \left(\sum_{i=1}^N (c_{Ni} - \bar{c}_N)^2 \right)^{-1/2} \sum_{i=1}^N (c_{Ni} - \bar{c}_N) a_N^*(R_i, t),$$

with a triangular array $\mathbf{c}_N = (c_{N1}, \dots, c_{NN})'$ of constants satisfying

$$\sum_{i=1}^N (c_{Ni} - \bar{c}_N)^2 / \max_{1 \leq i \leq N} (c_{Ni} - \bar{c}_N)^2 \xrightarrow{N \rightarrow \infty} \infty, \quad \bar{c}_N = N^{-1} \sum_{i=1}^N c_{Ni}$$

and showed that T_N converges in the uniform topology on $C[0, 1]$ to the Brownian bridge. We define the empirical distribution functions of the two samples $\hat{F}_m(x) = m^{-1} \sum_{i=1}^m I[Y_i \leq x]$ and $\hat{G}_n(x) = n^{-1} \sum_{i=m+1}^N I[Y_i \leq x]$ and the zero-one quantity V_i , $V_i = 1$ if $Y_{(i)}$ is one of Y_1, \dots, Y_m , $i = 1, \dots, N$.

Setting $\mathbf{c}_N = \mathbf{x}^{(p)}$, $\max_{0 \leq t \leq 1} T_N(t)$ coincides with the classical K-S two-sample test statistic T^+ . We use the fact that (2) implies $\sum_{i=1}^N a_N^*(R_i, j/N) = (1 - j/N)N = N - j$ and that $\max_{0 \leq t \leq 1} T_N(t) = \max_{1 \leq j \leq N} T_N(j/N)$ since the process $T_N(t)$ is linear on every

interval $[(j-1)/N, j/N]$, $j = 1, \dots, N$:

$$\begin{aligned}
T^+ &= \left(\frac{mn}{N}\right)^{1/2} \max_{1 \leq j \leq N} [\hat{G}_n(Y_{(j)}) - \hat{F}_m(Y_{(j)})] \\
&= \left(\frac{mn}{N}\right)^{1/2} \max_{1 \leq j \leq N} \left[\frac{1}{n}((1-V_1) + \dots + (1-V_j)) - \frac{1}{m}(V_1 + \dots + V_j) \right] \\
&= \left(\frac{N}{mn}\right)^{1/2} \max_{1 \leq j \leq N} \left[\frac{jm}{N} - (V_1 + \dots + V_j) \right] \\
&= \left(\frac{N}{mn}\right)^{1/2} \max_{1 \leq j \leq N} \left[\frac{jm}{N} - \sum_{i=1}^m \left(1 - a_N^*\left(R_i, \frac{j}{N}\right)\right) \right] \\
&= \left(\frac{N}{mn}\right)^{1/2} \max_{1 \leq j \leq N} \left[\sum_{i=1}^m a_N^*\left(R_i, \frac{j}{N}\right) - \frac{m}{N}(N-j) \right] \\
&= \left(\frac{N}{mn}\right)^{1/2} \max_{1 \leq j \leq N} \left[\left(1 - \frac{m}{N}\right) \sum_{i=1}^m a_N^*\left(R_i, \frac{j}{N}\right) - \frac{m}{N} \sum_{i=m+1}^N a_N^*\left(R_i, \frac{j}{N}\right) \right] \\
&= \max_{0 \leq t \leq 1} T_N(t).
\end{aligned}$$

2.2. Main result

Let us first recall that $\mathbf{x}^{(p)} = (1, \dots, 1, 0, \dots, 0)'$ and $\mathbf{x}^* = (x_1^*, \dots, x_N^*)' = \mathbf{H}^{(1)}\mathbf{x}^{(p)}$. The projection matrix $\mathbf{H}^{(1)} = (h_{ij}^{(1)})_{i=1, \dots, N}^{j=1, \dots, N}$ corresponding to $\mathbf{X}^{(1)}$ is idempotent, so

$$\begin{aligned}
\sum_{i=1}^N (x_i^{(p)} - x_i^*)^2 &= (\mathbf{x}^{(p)} - \mathbf{x}^*)'(\mathbf{x}^{(p)} - \mathbf{x}^*) = \mathbf{x}^{(p)'}(\mathbf{I}_N - \mathbf{H}^{(1)})(\mathbf{I}_N - \mathbf{H}^{(1)})\mathbf{x}^{(p)} \\
&= \mathbf{x}^{(p)'}(\mathbf{I}_N - \mathbf{H}^{(1)})\mathbf{x}^{(p)} = m - \sum_{i=1}^m \sum_{j=1}^m h_{ij}^{(1)} > 0.
\end{aligned}$$

Theorem 1. Assume that \mathbf{X}_N satisfies (X.1)–(X.4) and F satisfies (F.1)–(F.4). Let $\hat{\mathbf{a}}_N(\alpha) = (\hat{a}_{N1}(\alpha), \dots, \hat{a}_{NN}(\alpha))'$, $0 \leq \alpha \leq 1$ be the regression rank scores corresponding to the submodel of the model (1), i.e. under H_0 . Then the process $\{S_N(t): 0 \leq t \leq 1\}$,

$$\begin{aligned}
S_N(t) &= \left(\sum_{i=1}^N (x_i^{(p)} - x_i^*)^2 \right)^{-1/2} \sum_{i=1}^N (x_i^{(p)} - x_i^*) \hat{a}_{Ni}(t) \\
&= \left(m - \sum_{i=1}^m \sum_{j=1}^m h_{ij}^{(1)} \right)^{-1/2} \left[\sum_{i=1}^m \left(1 - \sum_{j=1}^m h_{ij}^{(1)} \right) \hat{a}_{Ni}(t) \right. \\
&\quad \left. + \sum_{i=m+1}^N \left(- \sum_{j=1}^m h_{ij}^{(1)} \right) \hat{a}_{Ni}(t) \right],
\end{aligned}$$

converges to the Brownian bridge in the uniform topology on $C[0, 1]$. Thus, for $K_N^+ = \max_{0 \leq t \leq 1} S_N(t)$ and $K_N = \max_{0 \leq t \leq 1} |S_N(t)|$ we can write, under H_0 ,

$$\lim_{N \rightarrow \infty} P(K_N^+ < x) = \begin{cases} 1 - e^{-2x^2}, & x \geq 0, \\ 0, & x < 0, \end{cases}$$

$$\lim_{N \rightarrow \infty} P(K_N < x) = \begin{cases} 1 - 2 \sum_{z=1}^{\infty} (-1)^{z+1} e^{-2z^2 x^2}, & x \geq 0, \\ 0, & x < 0. \end{cases}$$

Proof. It follows from (3) and from the properties of the Brownian bridge. \square

The statistics K_N^+ and K_N , similarly to the classical K-S statistics, can be used for testing the two sample problem with nuisance regression against one-sided and two-sided alternatives.

By Theorem 1 the test based on K_N^+ rejects H_0 on the asymptotic significance level α provided $K_N^+ \geq (-\frac{1}{2} \log \alpha)^{1/2}$.

The asymptotic power of the test based on K_N^+ , against the local alternative

$$H_N: \beta_p = N^{-1/2} \Delta, \quad \text{with } \beta_1, \dots, \beta_{p-1} \text{ unspecified,}$$

can be obtained from the following theorem.

Theorem 2. *Under the conditions of Theorem 1 and under H_N , the process*

$$S_N(t) - \left[\left(m - \sum_{i=1}^m \sum_{j=1}^m h_{ij}^{(1)} \right)^{1/2} \Delta N^{-1/2} f(F^{-1}(t)) \right]$$

converges to the Brownian bridge $\{Z(t): 0 \leq t \leq 1\}$ in the uniform topology on $C[0, 1]$ from which it follows that

$$\lim_{N \rightarrow \infty} P(K_N^+ \geq x | H_N) = P \left(\max_{0 \leq t \leq 1} \left\{ Z(t) + \left(m - \sum_{i=1}^m \sum_{j=1}^m h_{ij}^{(1)} \right)^{1/2} \Delta N^{-1/2} f(F^{-1}(t)) \right\} \geq x \right)$$

for any $x > 0$. Additionally,

$$\begin{aligned} \lim_{N \rightarrow \infty} P \left(K_N^+ \geq \left(-\frac{\log \alpha}{2} \right)^{1/2} | H_N \right) - \alpha &= \left[2 \left(m - \sum_{i=1}^m \sum_{j=1}^m h_{ij}^{(1)} \right)^{1/2} \Delta N^{-1/2} \alpha \left(-\frac{\log \alpha}{2} \right)^{1/2} \right. \\ &\quad \left. \times \int_0^1 -\frac{f'(F^{-1}(u))}{f(F^{-1}(u))} \psi(u, \alpha) du \right] (1 + o(1)) \end{aligned}$$

holds for $\left(m - \sum_{i=1}^m \sum_{j=1}^m h_{ij}^{(1)}\right)^{1/2} \Delta N^{-1/2} \rightarrow 0$, where

$$\psi(u, \alpha) = 2\Phi\left[\left(-\frac{\log \alpha}{2}\right)^{1/2} (2u - 1)(u(1 - u))^{-1/2}\right] - 1,$$

$0 < \alpha < 1$ and Φ is the standard normal distribution function.

Proof. It follows from (3) and from [4, Theorem VI.3.2]. The last assertion follows from [4, Theorem VI.4.5]. \square

Remark 1 (2-way ANOVA model). For example, in two-way layout, we can use this two-sample K-S test, similarly to e.g. the Friedman test, for comparing two treatments applied on I blocks. The effects of the blocks would represent the nuisance regression here.

2.3. Cramér-von Mises type test

Similarly to the Kolmogorov-Smirnov test, we can generalize also the Cramér-von Mises type two-sample test for nuisance regression.

We first look at the location model. With the same notation as in Subsection 2.1, we put again $\mathbf{c}_N = \mathbf{x}^{(p)} = (1, \dots, 1, 0, \dots, 0)'$, and for $T_N(t)$ from (4) we have that the classical Cramér-von Mises two-sample test statistic M equals (see [4, III.1.3.11 and V.3.8])

$$M = \frac{1}{mn} \sum_{j=1}^{N-1} \left[\frac{jm}{N} - (V_1 + \dots + V_j) \right]^2 = \int_0^1 T_N^2(t) dt + \frac{1}{6N}.$$

In model (1) (for nuisance linear regression) the test criterion of the Cramér-von Mises type two-sample test is then $\int_0^1 S_N^2(t) dt$, where $S_N(t)$ is the process from Theorem 1, and it can be seen from the form of the test statistic that this test with a critical region $\{\int_0^1 S_N^2(t) dt \geq C\}$ is suitable only for two-sided alternatives (similarly to the statistic K_N). The critical values can be obtained from the following theorem.

Theorem 3. *Under the conditions of Theorem 1 we have*

$$\lim_{N \rightarrow \infty} P\left(\int_0^1 S_N^2(t) dt < x\right) = P\left(\sum_{j=1}^{\infty} \frac{X_j^2}{j^2 \pi^2} < x\right),$$

where X_1, X_2, \dots are independent standardized normal random variables.

Proof. It follows from Theorem 1 and from the property of the Brownian bridge stated in [4, Theorem V.3.3.c]. \square

All the tests proposed in this paper are based on the regression rank scores and their construction is inspired by the structure of the classical Kolmogorov-Smirnov (Cramér-von Mises) test. Therefore, they do not require a preliminary estimation of the nuisance parameter and their asymptotic distributions coincide with the classical tests.

Acknowledgement. The author would like to thank his supervisor Prof. RNDr. Jana Jurečková, DrSc. for helpful discussion.

References

- [1] *C. Gutenbrunner, J. Jurečková*: Regression rank scores and regression quantiles. *Ann. Stat.* 20 (1992), 305–330.
- [2] *C. Gutenbrunner, J. Jurečková, R. Koenker, S. Portnoy*: Tests of linear hypotheses based on regression rank scores. *J. Nonparametric Stat.* 2 (1993), 307–331.
- [3] *J. Hájek*: Extension of the Kolmogorov-Smirnov test to regression alternatives. *Proc. Bernoulli-Bayes-Laplace Seminar (L. LeCam, ed.)*. Univ. of California Press, 1965, pp. 45–60.
- [4] *J. Hájek, Z. Šidák*: *Theory of Rank Tests*. Academia, Praha, 1967.
- [5] *J. Jurečková*: Tests of Kolmogorov-Smirnov type based on regression rank scores. *Information Theory, Statistical Decision Functions, Random Processes*. Trans. 11th Prague Conf., Prague 1990, Vol. B (J. Á. Vášek, ed.). Academia & Kluwer, Praha & Dordrecht, 1992, pp. 41–49.
- [6] *R. Koenker, G. Bassett (1978)*: Regression quantiles. *Econometrica* 46 (1978), 33–50.

Author's address: *M. Schindler*, Charles University in Prague, Faculty of Mathematics and Physics, Department of Probability and Statistics, Sokolovská 83, 186 75 Prague 8, Czech Republic, e-mail: schindle@karlin.mff.cuni.cz, and Technical University in Liberec, Studentská 2, 461 17 Liberec 1, Czech Republic.