

Jana Jurečková

Trimmed polynomial regression

Commentationes Mathematicae Universitatis Carolinae, Vol. 24 (1983), No. 4, 597--607

Persistent URL: <http://dml.cz/dmlcz/106259>

Terms of use:

© Charles University in Prague, Faculty of Mathematics and Physics, 1983

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>

TRIMMED POLYNOMIAL REGRESSION
Jana JUREČKOVÁ

Abstract : Robust test about the degree of the polynomial regression is suggested. The test is based on the trimmed least-squares estimator due to Koenker and Bassett [4] and has an asymptotically distribution-free critical region for a general class of distributions. The Pitman efficiency of the test coincides with the relative asymptotic efficiency of the trimmed least-squares estimator to the ordinary least-squares estimator.

Key words : Polynomial regression, regression quantile, trimmed least-squares estimator.

Classification: 62G10, 62J05, 62G20

1. **Introduction.** Let us consider the polynomial regression model

$$(1.1) \quad Y_{ni} = \beta_0 + \beta_1 x_{ni} + \dots + \beta_p x_{ni}^p + E_{ni}, \quad i=1, \dots, n$$

where $\underline{Y}_n = (Y_{n1}, \dots, Y_{nn})'$ is the vector of independent observations, $\underline{x}_n = (x_{n1}, \dots, x_{nn})'$ is a given vector, $\underline{\beta} = (\beta_0, \beta_1, \dots, \beta_p)'$ is a $(p+1) \times 1$ vector of unknown parameters and $\underline{E}_n = (E_{n1}, \dots, E_{nn})'$ is the vector of errors which are independent and identically distributed (i.i.d.) random variables with a continuous distribution function (d.f.) F . In the subsequent text, we shall omit the subscript n in Y_{ni} , x_{ni} , E_{ni} , etc., unless it causes a confusion. Our main

interest is in robust testing the hypothesis about the degree of the polynomial regression, i.e.,

$$(1.2) \quad H_0: \beta_j = 0, \quad j = p+1-m, \dots, p, \quad 1 \leq m \leq p.$$

Koenker and Bassett [4] introduced the concept of regression quantile, which seems to provide a basis for L-estimation and L-testing in the general linear model. The same authors proposed the trimmed least-squares estimator (trimmed LSE) as an extension of the trimmed mean to the linear model. This estimator was later on studied by Ruppert and Carroll [8] ; they considered a special design with an intercept and such that the slope-columns of the design matrix sum-up to zero.

The idea to use the regression quantiles for testing the linear hypothesis was mentioned in Ruppert and Carroll [8] who proposed a test based on the trimmed LSE under the special design mentioned above. Koenker and Bassett [5] studied several robust tests of linear exclusion hypothesis based on l_1 -estimation, i.e. on minimizing the sum of absolute residuals. However, neither of the mentioned procedures covers the general polynomial regression. Jurečková [3] derived the Bahadur-type representation and the asymptotic distribution of the regression quantiles and of the trimmed LSE under a more general design. The model considered in [3] covers the polynomial regression and thus enables to construct a robust test of H_0 . The test statistic is asymptotically distributed according to χ^2 distribution with m degrees of freedom under H_0 and according to noncentral χ^2 distribution under the contiguous alternatives. The Pitman efficiency of the test with respect to the classical one based on the or-

dinary LSE coincides with the relative asymptotic efficiency of the trimmed LSE to the ordinary LSE.

2. Notation and preliminary results. Let us fix α_1, α_2 , $0 < \alpha_1 < \alpha_2 < 1$. We shall start from the polynomial regression model (1.1); we shall assume that E_1, \dots, E_n are i.i.d. with the d.f. $F(x)$ which satisfies the following set of conditions (A):

- (A.1) F is absolutely continuous with the density f .
- (A.2) $0 < f(x) < \infty$ for $\xi_1 - \varepsilon < x < \xi_2 + \varepsilon$, $\varepsilon > 0$
where $\xi_i = F^{-1}(\alpha_i)$, $i=1,2$.
- (A.3) The derivative f' of f exists and is bounded in neighborhoods of ξ_1 and ξ_2 .

Denote

$$(2.1) \quad \tilde{X}_n = \tilde{X} = \begin{pmatrix} 1 & x_1 & \dots & x_1^p \\ \dots & \dots & \dots & \dots \\ 1 & x_n & \dots & x_n^p \end{pmatrix}$$

the $n \times (p+1)$ matrix. The i -th row of \tilde{X}_n will be denoted by \tilde{x}_i' . We shall assume that the sequence $\{\tilde{X}_n\}_{n=1}^{\infty}$ satisfies the following set of conditions (B):

- (B.1) $\lim_{n \rightarrow \infty} \frac{1}{n} \tilde{X}_n' \tilde{X}_n = Q$, $Q = (q_{jk})_{j,k=0, \dots, p}$
where Q is a positively definite $(p+1) \times (p+1)$ matrix.
- (B.2) $\max_{\substack{1 \leq i \leq n \\ 1 \leq j \leq p}} |x_i|^j = O(n^{1/4})$, as $n \rightarrow \infty$.
- (B.3) $\max_{1 \leq i \leq n} |x_i|^j \left(\sum_{k=1}^n |x_k|^j \right)^{-1} \rightarrow 0$, as $n \rightarrow \infty$; $j=1, \dots, p$.

The condition (B.1) means that

$$(2.2) \quad \frac{1}{n} \sum_{i=1}^n x_i^j \rightarrow q_j, \quad j=0, 1, \dots, 2p$$

where q_0, \dots, q_{2p} satisfy

$$(2.3) \quad q_{k,j} = q_{k+j}, \quad k, j=0, \dots, p.$$

As an example of $\sum_{i=1}^n$ satisfying (B) may serve the following sequence

$$(2.4) \quad x_{ni} = J\left(\frac{i}{n+1}\right), \quad i=1, \dots, n$$

where $J(t) : [0, 1] \rightarrow R^1$ is a bounded function; then

$$\frac{1}{n} \sum_{i=1}^n x_{ni}^k \rightarrow \int_0^1 (J(t))^k dt, \quad \text{as } n \rightarrow \infty.$$

For $\alpha = \alpha_1, \alpha_2$, denote

$$(2.5) \quad \varphi_\alpha(x) = \alpha - I[x < 0], \quad x \in R^1$$

and

$$(2.6) \quad \varrho_\alpha(x) = x \cdot \varphi_\alpha(x), \quad x \in R^1.$$

The α -regression quantile $\hat{\beta}_n(\alpha)$ is then defined as the $(p+1) \times 1$ vector $\tilde{t} = (t_0, \dots, t_p)'$ which solves

$$(2.7) \quad \sum_{i=1}^n \varrho_\alpha(Y_i - \tilde{x}_i' \tilde{t}) = \min.$$

The solution of (2.7) is generally not uniquely determined; suppose that a rule is given which selects a unique element of the set of solutions for $\alpha = \alpha_1, \alpha_2$.

It then follows from Theorem 2.1 of Jurečková [3] that

$$(2.8) \quad \begin{aligned} & n^{1/2} (\hat{\beta}_n(\alpha) - \beta_0 - g_1 F^{-1}(\alpha)) \\ &= n^{-1/2} [F(F^{-1}(\alpha))]^{-1} Q^{-1} \sum_{i=1}^n x_i \varphi_\alpha(x_i - F^{-1}(\alpha)) + o_p(n^{-1/4}) \end{aligned}$$

for $\alpha = \alpha_1, \alpha_2$, where $g_1 = (1, 0, \dots, 0)'$ is a $(p+1) \times 1$ vector,

Let A_n be the diagonal $n \times n$ matrix with the diagonal

$$(2.9) \quad a_{ii} = a_i = \begin{cases} 0 & \text{if } Y_i < \tilde{x}_i' \hat{\beta}_n(\alpha_1) \text{ or } Y_i > \tilde{x}_i' \hat{\beta}_n(\alpha_2) \\ 1 & \text{otherwise, } i=1, \dots, n. \end{cases}$$

The trimmed LSE is then defined as the ordinary LSE calcula-

ted after trimming-off Y_i with $a_i = 0, i=1, \dots, n$, i.e.,

$$(2.10) \quad \underline{L} = \underline{L}_{\underline{n}}(\alpha_1, \alpha_2) = (\underline{X}' \underline{A} \underline{X})^{-1} (\underline{X}' \underline{A} \underline{Y}).$$

It then follows from Theorem 3.1 of Jurečková [3] that

$$(2.11) \quad \begin{aligned} & n^{1/2} [\underline{L}_{\underline{n}}(\alpha_1, \alpha_2) - \underline{\beta} - \underline{e}_1 \delta] \\ &= n^{-1/2} (\alpha_2 - \alpha_1)^{-1} Q^{-1} \sum_{i=1}^n x_i (\psi(E_i) - E\psi(E_i)) + o_p(n^{-1/4}) \end{aligned}$$

where

$$(2.12) \quad \psi(x) = \begin{cases} \xi_1 & \text{if } x < \xi_1 \\ x & \text{if } \xi_1 \leq x \leq \xi_2 \\ \xi_2 & \text{if } \xi_2 < x \end{cases}$$

$$\text{and } \delta = (\alpha_2 - \alpha_1)^{-1} \int_{\alpha_1}^{\alpha_2} F^{-1}(u) du.$$

Consequently (cf. Theorem 3.2 of [3]),

$$(2.13) \quad \begin{aligned} & \mathcal{L} \{ n^{1/2} (\underline{L}_{\underline{n}}(\alpha_1, \alpha_2) - \underline{\beta} - \underline{e}_1 \delta) \} \\ & \rightarrow N_{p+1}(0, \sigma^2(\alpha_1, \alpha_2, F) Q^{-1}), \text{ as } n \rightarrow \infty \end{aligned}$$

where

$$(2.14) \quad \begin{aligned} \sigma^2(\alpha_1, \alpha_2, F) &= (\alpha_2 - \alpha_1)^{-2} \left\{ \int_{\alpha_1}^{\alpha_2} (F^{-1}(u) - \delta)^2 du \right. \\ & \left. + \alpha_1 (\xi_1 - \delta)^2 + (1 - \alpha_2) (\xi_2 - \delta)^2 - [\alpha_1 (\xi_1 - \delta) + (1 - \alpha_2) (\xi_2 - \delta)]^2 \right\}. \end{aligned}$$

We may see from (2.11) that only the first component L_0 of \underline{L} is generally asymptotically biased. The asymptotic variance (2.14) coincides with that of the trimmed mean in the location model.

3. Test of H_0 . Let us turn back to the polynomial regression model (1.1). The distribution function F of the

errors E_1, \dots, E_n is generally unspecified; we shall only assume that F satisfies the condition (A) for fixed α_1, α_2 , $0 < \alpha_1 < \alpha_2 < 1$. We wish to construct a test of the hypothesis

$$(3.1) \quad H_0 : \beta_j = 0, \quad j=p+1-m, \dots, p \quad (1 \leq m \leq p),$$

which is insensitive to the special shape of F .

Assume that the matrix \underline{X}_n satisfies the condition (B) of Section 2. Denote

$$(3.2) \quad \underline{X}_n^* = \begin{pmatrix} 1 & x_1 & \dots & x_1^{p-m} \\ \dots & \dots & \dots & \dots \\ 1 & x_n & \dots & x_n^{p-m} \end{pmatrix}$$

the $n \times (p+1-m)$ submatrix of \underline{X}_n . Then, according to (B.1),

$$(3.3) \quad \frac{1}{n} \underline{X}_n^{*'} \underline{X}_n^* \rightarrow \underline{Q}^*, \quad \text{as } n \rightarrow \infty,$$

where \underline{Q}^* is a positively definite $(p+1-m) \times (p+1-m)$ submatrix of \underline{Q} . Denote

$$(3.4) \quad \underline{Q}^* = \left\{ \begin{array}{cc} \underline{Q} & \underline{0} \\ \underline{0} & \underline{0} \end{array} \right\} \begin{array}{l} p+1-m \\ m \end{array}$$

the $(p+1) \times (p+1)$ matrix. Let $\underline{L} = \underline{L}_n(\alpha_1, \alpha_2)$ denote the trimmed LSE defined in Section 2; let $\underline{L}^* = \underline{L}_n^*(\alpha_1, \alpha_2)$ denote the trimmed LSE calculated under the assumption that H_0 is true. Then \underline{L}^* is a $(p+1-m) \times 1$ vector; denote

$$(3.5) \quad \underline{L}^* = (L_0^*, \dots, L_{p-m}^*, 0, \dots, 0)'$$

its extension to $(p+1) \times 1$ vector. Consider the statistic

$$(3.6) \quad T_n = (\underline{L} - \underline{L}^*)' \underline{X}' \underline{X} (\underline{L} - \underline{L}^*) / S_n^2$$

where

$$(3.7) \quad S_n^2 = (\alpha_2 - \alpha_1)^{-2} \left\{ (n-m)^{-1} Z_n^2 + \alpha_1 (\hat{\beta}_0(\alpha_1) - L_0)^2 + (1 - \alpha_2) (\hat{\beta}_0(\alpha_2) - L_0)^2 - [\alpha_1 (\hat{\beta}_0(\alpha_1) - L_0) + (1 - \alpha_2) (\hat{\beta}_0(\alpha_2) - L_0)]^2 \right\}$$

and

$$(3.8) \quad Z_n^2 = \underset{\sim}{Y}' \underset{\sim}{A} [I_{p+1} - \underset{\sim}{X}(\underset{\sim}{X}' \underset{\sim}{X})^{-1} \underset{\sim}{X}] \underset{\sim}{A} \underset{\sim}{Y} ;$$

$\hat{\beta}_0(\alpha_i), L_0$ are the first components of $\hat{\beta}(\alpha_i), L$, respectively, $i=1,2$.

We propose T_n as a test criterion for testing H_0 . The corresponding asymptotic critical region is given in the following theorem.

Theorem 3.1. Let Y_1, \dots, Y_n be independent observations satisfying the model (1.1) with the i.i.d. errors distributed according to the d.f. F satisfying the condition (A). Let X_n satisfy the condition (B). Then, under H_0 , the statistics T_n are asymptotically distributed, as $n \rightarrow \infty$, according to χ^2 distribution with m degrees of freedom.

The following theorem gives the asymptotic distribution of T_n under the local alternatives.

Theorem 3.2. Let Y_1, \dots, Y_n and X_n satisfy the assumptions of Theorem 3.1. Then, under the sequences of alternatives:

$$(3.9) \quad K_n : \beta_{nk} = \sqrt{n} b_k ; b_k \in R^1, \quad k=p+1-m, \dots, p,$$

the statistics T_n are asymptotically distributed according to noncentral χ^2 distribution with m degrees of freedom and with the noncentrality parameter

$$(3.10) \quad \eta^2 = \underset{\sim}{\bar{b}}' \underset{\sim}{Q} \underset{\sim}{\bar{b}} / \sigma^2(\alpha_1, \alpha_2, F)$$

where

$$(3.11) \quad \underset{\sim}{\bar{b}} = (0, \dots, 0, \underbrace{b_{p+1-m}, \dots, b_p}_m)'$$

It follows from Theorem 3.2 that the Pitman efficiency of

the test based on T_n with respect to the classical F-test coincides with the relative asymptotic efficiency of the trimmed LSE to the ordinary LSE.

4. Proofs of Theorems 3.1 and 3.2. The theorems will be proved with the aid of three lemmas.

Lemma 4.1. Under the assumptions of Theorem 3.1, the statistics

$$(4.1) \quad (\hat{\sigma}^2(\alpha_1, \alpha_2, F))^{-1} V_n$$

where

$$(4.2) \quad V_n = (\underline{L} - \underline{L}^*)' \underline{X}' \underline{X} (\underline{L} - \underline{L}^*)$$

are asymptotically distributed as χ^2 with m degrees of freedom under H_0 and as noncentral χ^2 with m degrees of freedom and noncentrality parameter η^2 of (3.10) under K_n , respectively.

Proof. Let us first consider the asymptotic distribution under H_0 . Consider the partitions

$$(4.3) \quad \underline{Q} = \left(\begin{array}{cc} Q_{11} & Q_{12} \\ \underbrace{Q_{21}}_{p+1-m} & \underbrace{Q_{22}}_m \end{array} \right) \left. \begin{array}{l} \} \\ \} \end{array} \right\} \begin{array}{l} p+1-m \\ m \end{array}$$

and

$$(4.4) \quad \underline{Q}^{-1} = \left(\begin{array}{cc} Q^{11} & Q^{12} \\ \underbrace{Q^{21}} & \underbrace{Q^{22}} \end{array} \right)$$

Moreover, denote

$$(4.5) \quad \underline{B} = \left(\begin{array}{cc} \underline{Q}^{*-1} & \underline{0} \\ \underline{0} & \underline{0} \end{array} \right)$$

the $(p+1) \times (p+1)$ matrix and

$$(4.6) \quad \underline{C} = \underline{Q}^{-1} - \underline{B}.$$

It follows from the symmetry of \underline{C} that

$$(4.7) \quad \underline{C} = \underline{Q}^{-1} \underline{Q} \underline{C} = \begin{pmatrix} \underline{Q}^{12} (\underline{Q}^{22})^{-1} \underline{Q}^{21} & \underline{Q}^{12} \\ \underline{Q}^{21} & \underline{Q}^{22} \end{pmatrix}$$

and

$$(4.8) \quad \underline{C}' \underline{Q} \underline{C} = \underline{C}; \quad \underline{C} \text{ is of rank } m.$$

It follows from (B.1), (4.8) and (2.11) that, under H_0 ,

$$(4.9) \quad (\underline{L} - \underline{L}^*)' \underline{X}' \underline{X} (\underline{L} - \underline{L}^*) \\ = (\alpha_2 - \alpha_1)^{-2} \left[n^{-1/2} \sum_{i=1}^n \underline{x}_i (\psi(E_i) - E\psi(E_i)) \right]' \underline{C} \\ \cdot \left[n^{-1/2} \sum_{i=1}^n \underline{x}_i (\psi(E_i) - E\psi(E_i)) \right] + o_p(1)$$

as $n \rightarrow \infty$, and it follows from the classical central limit theorem that

$$(4.10) \quad \mathcal{L} \left\{ (\alpha_2 - \alpha_1)^{-1} n^{-1/2} \sum_{i=1}^n \underline{x}_i (\psi(E_i) - E\psi(E_i)) \right\} \\ \rightarrow N_{p+1}(0, \sigma^2(\alpha_1, \alpha_2, F) \underline{Q}), \text{ as } n \rightarrow \infty.$$

(4.8), (4.9) and (4.10) together with Proposition VIII of Section 8a.2 of Rao [7] imply that (4.1) is, under H_0 , asymptotically χ^2 distributed with m degrees of freedom.

Proceeding quite analogously, we get that, under K_n ,

$$(4.11) \quad (\sigma^2(\alpha_1, \alpha_2, F))^{-1} (\underline{L} - \underline{L}^* - n^{-1/2} \underline{b})' \underline{X}' \underline{X} (\underline{L} - \underline{L}^* - n^{-1/2} \underline{b}) \\ \text{is asymptotically } \chi^2\text{-distributed with } m \text{ degrees of freedom;} \\ \text{this completes the proof of the lemma.}$$

Lemma 4.2. Under the conditions (A) and (B),

$$(4.12) \quad S_n^2 \xrightarrow{P} \sigma^2(\alpha_1, \alpha_2, F), \quad \text{as } n \rightarrow \infty.$$

Proof. The estimator S_n^2 of σ^2 was suggested by Ruppert

and Carroll [8] who proved its consistency under a special design. The proof of the lemma rests on the following lemma which could be proved quite analogously as Lemma 3.1 of Jurečková [3].

Lemma 4.3. Let U_1, \dots, U_n be i.i.d. random variables with the d.f. F satisfying the condition (A). Denote

$$(4.13) \quad T_n(\underline{v}) = n^{-1/2} \sum_{i=1}^n U_i^2 I[U_i \leq F^{-1}(\alpha) + n^{-1/2} \underline{x}_i' \underline{v}] \quad , \quad \underline{v} \in R^{p+1}$$

with $\alpha = \alpha_1, \alpha_2$. Then, provided the matrix X_n with the rows \underline{x}_i' , $i=1, \dots, n$ satisfy the condition (B),

$$(4.14) \quad \sup_{|\underline{v}| \leq K} |T_n(\underline{v}) - T_n(0) - n^{-1} (F^{-1}(\alpha))^2 f(F^{-1}(\alpha)) \sum_{i=1}^n \underline{x}_i' \underline{v}| \xrightarrow{P} 0$$

as $n \rightarrow \infty$, for any $K > 0$.

It follows from (4.14) and from (2.8) that

$$(4.15) \quad (n-m)^{-1} \sum_{i=1}^n a_i E_i^2 \xrightarrow{P} \int_1^2 x^2 dF(x) \quad , \quad \text{as } n \rightarrow \infty .$$

Moreover, by Lemma 3.2 of [3],

$$(4.16) \quad n^{-1} \underline{X}' \underline{A} \underline{X} = (\alpha_2 - \alpha_1) \underline{Q} + o_p(1)$$

and thus, by (2.10),

$$(4.17) \quad \begin{aligned} & (n-m)^{-1} \underline{E}' \underline{A} \underline{X} (\underline{X}' \underline{A} \underline{X})^{-1} \underline{X}' \underline{A} \underline{E} \\ &= (n-m)^{-1} (\underline{L} - \underline{\beta})' (\underline{X}' \underline{A} \underline{X}) (\underline{L} - \underline{\beta}) + o_p(1) \\ &= (n-m)^{-1} (\alpha_2 - \alpha_1) (\underline{L} - \underline{\beta})' \underline{Q} (\underline{L} - \underline{\beta}) + o_p(1) \end{aligned}$$

and this, by (2.11), can be rewritten as

$$(4.18) \quad \begin{aligned} & (\alpha_2 - \alpha_1)^{-1} \left[n^{-1} \sum_{i=1}^n \underline{x}_i (\psi(E_i) - E\psi(E_i)) \underline{Q}^{-1} + e_1 (\alpha_2 - \alpha_1) \delta \right]' \underline{Q} \\ & \cdot \left[n^{-1} \sum_{i=1}^n \underline{x}_i (\psi(E_i) - E\psi(E_i)) \underline{Q}^{-1} + e_1 (\alpha_2 - \alpha_1) \delta \right] + o_p(1) \\ &= (\alpha_2 - \alpha_1) \delta^2 + o_p(1). \end{aligned}$$

Combining (4.15), (4.17) and (4.18), we get

$$(4.19) \quad (n-m)^{-1} Z_n^2 \xrightarrow{P} \int_{\sum_1}^2 (x-\delta)^2 dF(x).$$

Moreover, it follows from (2.8) and (2.11) that

$$(4.20) \quad \hat{\beta}_0(\alpha_i) - L_0 = F^{-1}(\alpha_i) - \delta + o_p(1), \quad i=1,2.$$

(4.19) and (4.20) then complete the proof of Lemma 4.2.

Theorems 3.1 and 3.2 then follow from Lemmas 4.1 and 4.2.

5. References.

- [1] J.JUREČKOVÁ: Robust estimators of location and regression parameters and their second order asymptotic relations, Trans. 9th Prague Conf.on Inform.Theory,Statist.Decis. Functions and Random Proc.,pp.19-32. Academia,Praha (1983).
- [2] J.JUREČKOVÁ: Winsorized least-squares estimator and its M-estimator counterpart, Contributions to Statistics: Essays in Honour of Norman L.Johnson (P.K.Sen,ed.), pp.237-245. North Holland (1983).
- [3] J.JUREČKOVÁ: Regression quantiles and trimmed least squares estimator under a general design. Submitted.
- [4] R.KOENKER and G.BASSETT: Regression quantiles, Econometrica 46(1978), 33-50.
- [5] R.KOENKER and G.BASSETT: Tests of linear hypotheses and l_1 estimation, Econometrica 50(1982), 1577-1583.
- [6] E.L.LEHMANN: Testing Statistical Hypotheses. J.Wiley (1959).
- [7] C.R.RAO: Linear Statistical Inference and Its Application. J.Wiley (1973).
- [8] D.RUPPERT and R.J.CARROLL: Trimmed least-squares estimation in the linear model. J.Amer.Statist.Assoc.75(1980), 828-838.

Matematicko-fyzikální fakulta
Universita Karlova
Šokolovská 83, 186 00 Praha 8
Československo

(Oblatum 18.7. 1983)