

Süleyman Demirel Üniversitesi
İktisadi ve İdari Bilimler Fakültesi
Y.2005, C.10, S.1 s.263-279.

**AŞIRI DEĞER İÇEREN VERİ KÜMELERİNDE HATA
TERİMLERİNİN BİNOM DAĞILIMA UYDUĞU
DURUMDA EKK VE ROBUST LTS REGRESYON
TAHMİNCİLERİNİN SİMULASYON ÇALIŞMASI İLE
KARŞILAŞTIRILMASI**

**THE COMPARISON OF OLS METHOD WITH ROBUST
LTS REGRESSION ESTIMATORS SIMULERTION IN
CASE OF BINOMIAL DISTRIBUTION OF ERROR
TERMS IN DATA SETS INCLUDING OUTLIERS**

**Dr.Latif ÖZTÜRK
Dr. Hakan TÜRKAY***

ÖZET

Doğrusal regresyon modellerinin çözümü için kullanılacak bir çok alternatif yöntem vardır. Bu yöntemler arasında en çok tanınan, kullanılan ve uygulamada neredeyse standart bir mekanizma halinde gelen yöntem EKK(En Küçük Kareler)'dir. Yöntemin dayandığı, aralarında hata terimlerinin normal dağılıma uyduğu varsayımı da olan önemli bir çok varsayım söz konusudur. Bu varsayımların sağlandığı durumda şüphesiz EKK tahmincisi en iyidir. Ancak, uygulamada sıklıkla karşılaşılan durum varsayımların bazılarının sağlanamaması halidir ve bu durumda EKK istenen özelliklere uygun parametre tahminleri yapamaz. Bu sorun karşısında (varsayımlardan sapma halinde) varsayımlardan sapmaya duyarlı olan robust regresyon yöntemleri geliştirilmiştir. Bu çalışmada, hata terimlerinin Binom dağıldığı durumda robust regresyon yöntemlerinden biri olan LTS (Least Trimmed Squares) tahmincisinin EKK ile karşılaştırılması simülasyon çalışması yardımıyla yapılmıştır. Karşılaştırma, veri kümesinin aşırı değerler (outliers) içerdiği ve içermediği durumlarda gerçekleştirilmiştir. Yöntemlerden elde edilen parametre tahminleri başlangıçta atanan parametre değerlerinden gösterdikleri sapma miktarlarına göre iki yöntemin performansı değerlendirilmiştir.

ABSTRACT

There are many alternative methods which can be used for the solution of linear regression models. Among these methods, LS is the best known, the most commonly used method, and it is assumed as the standard mechanism in applications. There are many assumptions, one of which is that

* İnönü Üniversitesi İ.İ.B.F. Ekonometri Bölümü

“error terms are distributed normal”, on which the method is based. When the assumptions are proved, LS is the best method to be used. However, in practice, all of the assumptions mostly can not be proved, and in such a situation, LS can not estimate the parameters with required properties. Towards this problem (deviations from assumptions) robust regression methods which are insensitive to deviations from assumptions, are developed. In this study, LTS (Least Trimmed Regression), one of the robust regression methods, is compared with LS by using simulation approach, within the case that error terms are distributed binomial. This comparison is held within the situations where data include outliers and vice-versa. The performance of the methods are evaluated according to the bias between the parameter estimate which are obtained by using the regression methods and appointed as initial parameter values.

Robust Regresyon, LTS, Aşırı Değer, Simulasyon
Robust Regression, LTS, Outliers, Simulation

1. GİRİŞ

Regresyon analizinde öncelikli amaç, örneklem regresyon fonksiyonunu temel alarak anakütle regresyon fonksiyonunu olabildiğince doğruya yakın tahmin etmektir. Örneklem regresyon fonksiyonunu bulmak için kullanılan çeşitli yöntemler arasında en yaygın tanınan ve kullanılan yöntem, Carl Friedrich Gauss tarafından bulunan Sıradan En Küçük Kareler (EKK) yöntemidir. Yöntemin ana fikri kalıntı karelerinin toplamını minimum

yapmak ($\text{Min}_{\hat{\theta}} \sum_{i=1}^n r_i^2$) ve bu yolla en iyi tahmin değerlerine ulaşmaktır. Bu

minimizasyon işlemi sonucunda bulunan tahminciler en küçük kareler tahmincileri diye adlandırılmakta ve klasik istatistiğin çok önemli bir temel taşını oluşturmaktadır. Bu tahminciler, klasik doğrusal regresyon modelinin varsayımları veri iken sapmasız doğrusal tahminciler içerisinde en küçük varyanslı olanlardır; yani en iyileridir. Bu ifade, Gauss-Markov teoremi diye adlandırılır.

EKK yönteminin bu ölçüde yaygın kullanılmasının nedeni, kavramın kolaylıkla anlaşılabilmesi ve uygulanabilmesidir. Yöntemin bulunduğu zamanlarda bilgisayarların olmaması ve EKK tahmincilerinin veriler üzerinde matris cebri uygulamalarıyla kolayca hesaplanabilir olması, yalnızca EKK tahmincisini uygulanabilir bir yaklaşım yapmıştır. Gauss, hata teriminin normal (Gaussyen) dağıldığı durumda EKK'nın en iyi olduğunu gösterdikten sonra, Gaussyen varsayımlar ve EKK yöntemi uygulamada standart bir mekanizma haline getirilmiştir. Ancak, verilerin çoğu zaman klasik varsayımları tamamen sağlamadığı uzun yıllardır bilinen bir gerçektir. (örneğin, hata teriminin dağılımı her zaman normal değildir.) Buna bağlı olarak, bu sorunu aşabilmek ve varsayımların sağlanmadığı durumlarda da regresyon doğrusunu gerçeğe en yakın biçimde tahmin etmek amacıyla robust (dayanıklı) regresyon tahminciler önerilmiştir.

Regresyon analizi uygulamalarında, hata terimi sadece normal dağılıma uymamakla kalmayıp, bazı durumlarda kesikli dağılımlara uygun bir dağılımda göstermektedir. Örneğin, gölge bağımlı değişkenli doğrusal olasılık modellerinde hata terimi binom dağılıma uymaktadır. Bundan dolayı, bu çalışmada hata teriminin Binom dağılıma uyduğu durumda EKK ve yüksek bozulma sınırına sahip dayanıklı tahmincilerden biri olan LTS (Least Trimmed Square : En küçük Kırılmış Kareler) tahmincilerinin EKK yöntemiyle karşılaştırılması yapılarak, yöntemlerin üstünlük ve zaafı ortaya çıkarılmaya çalışılmıştır.

2. EKK YÖNTEMİNİN AŞIRI DEĞERLERE KARŞI DUYARLILIĞI

Bir veri kümesi dört tip gözlem değeri içerebilir : Düzgün (regular) gözlemler, dikey aşırı değerler, iyi kaldıraç noktaları ve kötü kaldıraç noktaları (yüksek etkili gözlemler).¹

Aşırı değerler, tipik olmayan ve verinin geri kalan kısmının karakteristik dağılımına uygun görüntü vermeyen gözlem değerleridir. Bu değerler, temel olgunun (değişkenin) özelliklerini yansıtmıyor olabileceği gibi, regresyon analizi uygulamalarında sıklıkla karşılaşılan bir durum olan bazı gözlem değerlerinin ölçme, kaydetme, kopyalama veya aktarma hataları nedeniyle yanlış bir değer olarak alınmasından da kaynaklanabilir. Sonuç olarak; “Regresyon aşırı değerleri, verilerin çoğunluğu tarafından oluşturulan doğrusal kalıba uymayan gözlemlerdir. Bu tip bir gözlem olağanüstü bir Y_i veya olağanüstü bir (X_{i1}, \dots, X_{ip}) veya da her ikisi birlikte olabilir.”²

Bağımlı değişken Y_i 'nin gözlem değerlerinden (X_k, Y_k) gibi bir nokta Y_i değerlerinin çoğunluğundan oldukça uzak bir yerde yer almakta ise, bu nokta Y yönünde aşırı değer olarak adlandırılır. Söz konusu nokta, EKK doğrusu üzerinde oldukça büyük bir etkiye sahip olacak ve doğruyu tamamen değiştirebilecektir. Bu tip dikey aşırı değerler büyük pozitif veya büyük negatif kalıntı (r_i) değerlerine sahiptirler. Ancak, bu durumda elde edilen kalıntı değeri (r_i) şüpheli bir büyüklüktür. Bu tip aşırı değerler kalıntıların listesinden veya kalıntıların grafiğinin çizimi yardımıyla tespit edilebilir.

Ancak, büyük hataların sadece bağımlı değişkende ortaya çıkması için herhangi bir neden yoktur, bu durum mutlak anlamda, açıklayıcı değişkenlerin birinde aşırı değerlerin bulunmasına çok benzerdir. Genelde örnekleme ki X_i değerlerinin büyük çoğunluğundan çok uzakta yer alan

¹ Peter J. Rousseeuw, Katrien Van Driessen, “Computing LTS Regression for Large Data Sets” (Düzenlenmiş), *Technical Report*, University of Antwerp, 1999, s.2.

² Peter J. Rousseeuw, Stefan Van Aelst, “Positive-Breakdown Robust Methods in Computer Vision” *Computing Science and Statistics*, Vol:31,1999, s. 451.

X_k gibi bir değer olduğunda (X_k, Y_k) değerine X yönünde aşırı değer denir. Veri kümesinin merkezine uzak olan bu nokta, regresyon doğrusunu büyük bir güçle çektiğinden mekanikteki kaldıraç kavramına benzer olarak kaldıraç noktası (leverage point) olarak adlandırılmaktadır. (X_k, Y_k) noktası verilerin çoğunluğu tarafından oluşturulan örüntüye uygunsuz yani çoğunluk veri tarafından belirlenen regresyon doğrusuna yakın bir yerde ise iyi kaldıraç noktası olarak adlandırılır.³ Çoğu zaman regresyon aşırı değerlerinin EKK kalıntılarının incelenmesiyle tespit edilebileceği düşünülmektedir. Ne yazık ki, aşırı değerlerin kaldıraç noktası olduğu durumlarda bu doğru değildir. Kaldıraç noktaları, çok büyük kalıntı değerlerine sahip olacağından $\sum_{i=1}^n r_i^2$ değerine büyük bir ilave katkıda

bulunur. Bu durumda, gerçek regresyon doğrusu EKK yaklaşımıyla seçilemez, çünkü EKK doğrusu söz konusu noktaya doğru kayacağından büyük kalıntı değeri gerçekte olduğundan küçük, diğer noktaların kalıntıları ise olması gerekenden büyük görünecektir. Eğer “en büyük EKK kalıntılarına sahip olan noktaların atılması” şeklinde bir ilke uygulanırsa, iyi veri noktalarının atılması söz konusu olabilir. Regresyon aşırı değerleri (ister X, isterse Y yönünde olsun) EKK yöntemi için ciddi bir tehdit oluşturmaktadırlar.⁴

Sonuç olarak, EKK gibi klasik parametre tahmin yöntemlerinin anlaşılması ve hesaplanması kolay olmasına rağmen, verinin içerdiği az sayıda aşırı değer bile söz konusu yöntemler tarafından tahmin edilen doğruyu keyfi bir büyüklükte etkilemekte ve parametre tahmin değerleri üzerinde önemli bir rol oynamaktadırlar. Bundan dolayı, dayanıklı tahminci olarak adlandırılan ve bahsedilen sorundan fazlaca etkilenmeyen tahmincilere ilgi hızla artmaktadır.

3. ROBUST (DAYANIKLI) REGRESYON TAHMİNCİLERİNİN KULLANILMA NEDENLERİ

Doğrusal regresyon analizi uygulamalarının bir çoğunda, serilerde aşırı değerlerin olması veya hata teriminin normal dağılımdan uzaklaşması gibi sorunlar ortaya çıkmaktadır. Söz konusu sorunların varlığı durumunda, daha önce belirtilen varsayımlardan küçük sapmalara bile aşırı duyarlılık göstermesi özelliğinden dolayı, EKK tahmincileri istenen özelliklere uygun tahminler verememektedir. Bu nedenle, EKK tahmincilerinde ortaya çıkan sorunların üstesinden gelebilecek ve aşırı değerlere karşı duyarsız olan dayanıklı regresyon yöntemlerine ihtiyaç duyulmuştur. Bu yöntemler, hata terimi normal dağılmadığında ve aşırı değerlerin varlığı durumunda klasik

³ Peter J. Rousseeuw, Annick M. Leroy, **Robust Regression and Outlier Detection**, New York, John Wiley, 1987, s.6.

⁴ Peter J. Rousseeuw, Annick M. Leroy, **a.g.e.**, s.7-8.

yöntemlere göre daha uygun sonuçlar vermektedir. Daha önce bahsedildiği gibi, verinin aşırı değerler içermesi durumunda, öncelikle verinin çeşitli dışlama kuralları kullanılarak aşırı değerlerden temizlenmesi, daha sonra kalan veri üzerinde klasik tahmin ve test yöntemlerinin uygulanmasının yeterli olacağı düşünülebilir. Ancak, EKK yöntemiyle yapılan tahminde regresyon doğrusu aşırı değerlere doğru çekileceğinden bu değerlere ait kalıntılar oldukça küçük olabilecektir. Bu durum, söz konusu aşırı değerleri maskeleyebilecektir. EKK yönteminin veri kümesi içinde aşırı değerlerin varlığı halinde kullanılamamasındaki diğer nedenler aşağıdaki gibi sıralanabilir:⁵

1. Yukarıdaki iki aşamayı (aşırı değerlerin teşhisi ve parametrelerin tahmini aşamalarını) birbirinden ayırabilmek nadiren mümkündür. Örneğin, çok değişkenli regresyon problemlerinde, parametreler için dayanıklı tahminlere başvurmadan aşırı değerlerin (doğru bir şekilde) tespit edilmesi zordur.
2. Orijinal gözlem değerleri kümesi, arasına büyük hatalar karıştırılmış normal gözlemlerden oluşmuş olsa bile, temizlenmiş veri normal olmayacaktır ve istatistiksel hatalar meydana gelecektir (yanlış dışlamalar veya alkoymalar nedeniyle). Sonuç, gerçekte normal olmayan bir dağılımdan türetildiği durumdakinden daha da kötü olacaktır. Dolayısıyla, klasik normal teorisinin temizlenmiş örneklemelere uygulanabilirliği söz konusu değildir. Bu durumdaki, iki aşamalı yöntemin işleyişi, doğru bir dayanıklı yöntemin işleyişinden çok daha zor olabilir.
3. En iyi dışlama yöntemlerinin performanslarının bile, iyi bir dayanıklı yöntemin performansına yeterli düzeyde ulaşamadığı da deneysel olarak gösterilebilen bir gerçektir. Dayanıklı yöntemler gözle görülür biçimde çok daha iyidir, çünkü bunlar değerlerin tamamen kabulü ile tamamen dışlanması arasında sorunsuz bir geçişi sağlayabilirler.

Ayrıca, aşırı değerlerin teşhisi, “ne kadar uç noktada” olduğunu ölçmek için bir ölçüye ihtiyaç duymaktadır. Bu ölçü, veri için bazı modellerden ve bu modelden ayrılığın bazı ölçülerinden ortaya çıkmaktadır. Çoklu aşırı değerler, söz konusu ölçünün kendisinin bir aşırı değer tarafından bozulmuş olabileceği tehdidini taşımaktadır. Aşırı değerler ölçüsünün kirlenmesi, aşırı değeri tespit edicisini bozmaktadır ve tabii ki aynı zamanda bu aşırı değeri tespit etme usulü üzerine kurulan herhangi bir tahminciyi de bozmaktadır.⁶

⁵ Peter J. Huber, **Robust Statistics**, New York, John Wiley and Sons, 1981, s.4-5.

⁶ N. Billor, A.S. Hadi, P.F. Vellman, “BACON : Blocked Adaptive Computational Efficient Outlier Nominators”, **Computational Statistics and Data Analysis**, Vol.34, s.280.

4. ROBUST REGRESYON TAHMİNCİLERİNİN AMAÇLARI

Robust (dayanıklı) regresyon tahmincileri, bir çok amaca yönelik olarak son yıllarda giderek artan bir hızda istatistik ve ekonometri de regresyon modellerinin parametrelerinin tahmin edilmesi için kullanılmaktadır. Bu yöntemlere başvurulmasındaki amaçlar aşağıdaki gibi özetlenebilir:⁷

1. Veri kümesine en iyi uyumu gösteren yapıyı tespit etmek,
2. Eğer gerekli görülürse, uzaktaki veri noktalarını (aşırı değerleri) veya daha ileri analizler için yapının temelindeki sapmaları teşhis etmek,
3. Yüksek düzeyde etkili veri noktalarını teşhis etmek ve bunlar üzerine dikkat çekmek,
4. KuşkuLANılmayan ardışık bağımlılık veya daha da genel olarak varsayılan korelasyon yapısındaki sapmalar ile ilgilenmektir.

Gözlem değerlerinin küçük bir kısmında büyük sapmaların bulunması küçük sapmalar olarak kabul edilebilirler; dayanıklı yöntemlerin öncelikli hedefi, bazı klasik yöntemlerin aşırı duyarlılığı da göz önüne alındığında, büyük hatalara karşı koruma sağlamasıdır. Literatürde dayanıklı yöntemler için açıkça veya ima yolu ile çok sayıda başka amaçlara da (örneğin, yüksek asimptotik görece etkinlik veya yüksek mutlak etkinlik gibi) yer verilmektedir. Ne var ki, bu amaçlar ikinci derecede önemlidir ve bunlara yukarıda bahsedilen amaçlardan daha fazla önem verilmemelidir.⁸

5. BOZULMA SINIRI (BREAKDOWN POINT)

Aşırı değerlere karşı nicel dayanıklılığın bir ölçüsü bozulma sınırıdır. “İlk olarak Hodges tarafından 1967’de kullanılan bozulma sınırı kavramı tek boyutlu konum tahmincileri ile sınırlıydı. Daha sonra Hampel tarafından 1968 ve 1971’de yapılan çalışmalarla genelleştirilmiştir.”⁹ Hampel 1971’de çok daha genel bir tanımlama yapmıştır, bu tanımlamada asimptotik olduğundan ve dolayısıyla matematiksel bir ifade olduğundan çok fazla yaygınlaşamadı. Ancak, 1983’de Donoho ve Huber tarafından yaygın bir şekilde tanınır hale getirilmişlerdir. Bozulma sınırının Donoho ve Huber tarafından tanımlanan basit bir sonlu örneklem versiyonu şöyledir:¹⁰

⁷ Frank R. Hampel, Elvezio M. Ronchetti, Peter J. Rousseeuw, Werner A. Stahel, **Robust Statistics**, New York, John Wiley & Sons, 1986, s.11.

⁸ Peter J. Huber, **a.g.e.**, s.5-6.

⁹ Yijun Zuo, “Some Quantitative Relationships Between Two Types of Finite Sample Breakdown Point”, **Statistics and Probability Letters**, vol.51, s.369-375, 2001.

¹⁰ Peter J. Rousseeuw, Annick M. Leroy, **a.g.e.**, s.9.

Z , n tane gözlemden oluşan bir örneklem ve T bir regresyon tahmincisi varsayılırsa;

$$Z = \{(X_{11}, \dots, X_{1p}, Y_1), \dots, (X_{21}, \dots, X_{np}, Y_n)\} \quad (1)$$

Z örneğine T yöntemi uygulanarak $T(Z) = \hat{\theta}$ şeklinde katsayılar vektörü bulunabilir. Burada, orijinal veri değerlerinden herhangi m tanesinin yerine keyfi bir değer yazılarak elde edilen kirlenmiş bir Z' düşünelim. Böyle bir kirlenme ile neden olunabilen sapma $\text{bias}(m; T, Z)$ ile gösterilirse;

$$\text{Sapma}(m, T, Z) = \sup_{Z'} \|T(Z') - T(Z)\| \quad (2)$$

Eğer sapma $(m; T, Z)$ sonsuz ise, m tane aşırı değer T üzerinde büyük bir etkiye sebep olduğu, hatta tahmincinin bozulduğu söylenebilir. Bu nedenle, Z sonlu örneğinde T tahmincisinin bozulma sınırı,

$$\varepsilon_n^*(T, Z) = \min \left\{ \frac{m}{n}, \text{sapma}(m; T, Z) \text{ sonsuz} \right\} \quad (3)$$

şeklinde tanımlanır. Bozulma değeri veya sınırı bir anlamda en kötü durum senaryosudur. “Kabaca ifade edilecek olursa, bozulma sınırı tahmincilerin başa çıkabilecekleri kötü aşırı değerlerin oranının sınırını verir.”¹¹ Diğer bir deyişle; bozulma sınırı T tahmincisinin $T(Z)$ den oldukça uzak keyfi değerler almaya başlamasına neden olabilecek en düşük kirlenme oranıdır.¹² Yani, bir tahmincinin bozulma sınırı, tahmincinin keyfi olarak büyük sapma değerleri almasına yol açabilen aşırı değerlerin oranı olarak tanımlanır.¹³ EKK'nın bozulma sınırı, örneklemdaki tek bir aşırı değerden bile etkilendiği için $1/n$ 'dir, dolayısıyla asimptotik durumda n sonsuza giderken sifıra yaklaşır.

Bir parametre tahmin yönteminin çöktüğü nokta, yani bozulma sınırı, sonuçta elde edilecek tahmin değerini anlamsızlaştıracak kadar kirlitebilen veya yönlendirebilen örneklemdaki noktaların en küçük kısmıdır.¹⁴ Sonuç olarak, bozulma sınırı bir tahmincinin anlamsız sonuçlar vermeden kullanılabileceği, gerçekte “kötü” kirlenmenin ne kadar olabileceği konusunda kaba bir fikir verir. Bozulma sınırı, tahmincinin kararlılığının global bir ölçüsüdür. Bundan dolayı, T tahmincisi tamamen işe yaramaz hale gelmeden önce, kirlenmiş modelin varsayılan modelden ne kadar

¹¹ Peter J. Huber, **a.g.e.**, s.13.

¹² Peter J. Rousseeuw, Annick M. Leroy, **a.g.e.**, s.10.

¹³ J. Matoušek, D.M. Mount, N.S. Natenyahu, “Efficient Randomized Algorithms For The Repeated Median Line Estimator”, **Algorithmica**, 20, 1998, s.136.

¹⁴ Feifang Hu, Jianhua Hu, “A Note On Breakdown Theory For Bootstrap Methods”, **Statistics and Probability Letters**, 50, 2000, s.49-53.

uzaklaşabileceğini gösterir. Aynı zamanda, bozulma sınırı, modele bağlı olmaması anlamında da global bir ölçüdür.¹⁵

6. ROBUST REGRESYON TAHMİNCİLERİNİN ÖZELLİKLERİ

Ele alınan herhangi bir istatistiksel yöntem aşağıdaki özelliklere sahip olmalıdır:¹⁶

1. Yöntem, varsayılan modelde makul bir ölçüde iyi (optimal veya optime yakın) etkinliğe sahip olmalıdır.
2. Model varsayımlarından küçük sapmaların yöntemin performansını sadece düşük bir düzeyde zayıflatması anlamında dayanıklı olmalıdır. Yani, sonraki (bir tahminin asimptotik varyansı veya testin gücü ve seviyesi ile tanımlanan) modelde hesaplanan nominal değere yakın olmalıdır.
3. Modelden biraz daha büyük sapmalar, büyük bir yıkıma neden olmamalıdır.

Yukarıda sayılan özellikler ile daha önce belirtilen dayanıklı yöntemlerin kullanım amaçları da dikkate alındığında, dayanıklı bir yöntem, veri kümesinin çoğunluğuna uygun yapıyı tanımlama ve aşırı değerleri teşhis etme özelliklerine de sahip olmalıdır. Sonuç olarak, dayanıklı bir tahmincinin sahip olması gereken özellikler aşağıdaki gibi özetlenebilir:¹⁷

1. Bu tahminciler, küçük kirlenmelere küçük tepkiler göstermelidir.
2. Büyük kirlenmelerin varlığında bile güvenilir olmalıdır.
3. Herhangi bir miktardaki kirlenmenin etkisi üzerinde bir sınıra sahip olmalıdır.
4. Yuvarlamaya ve gruplamaya düzgün olarak tepki göstermelidir.
5. Veriden açıkça görülen aşırı değerleri ayıklayabilmelidir.
6. Olabildiğince küçük bir varyansa sahip olmalıdır.

¹⁵ Robert G. Staudte, Simon J. Sheather, **Robust Estimation and Testing**, New York, John Wiley, 1990, s.56.

¹⁶ Peter J. Huber, **a.g.e.**, s.5.

¹⁷ Frank R. Hampel, "The Influence Curve and Its Role in Robust Estimation", **Journal of the American Statistical Association**, 69, 1974, s.383-393.

7. EN KÜÇÜK KIRPILMIŞ KARELER (LTS :LEAST TRIMMED SQUARES)

Yüksek bozulma sınırına sahip tahminciler hem X hem de Y yönündeki aşırı değerlerin varlığı durumunda (çok sayıda olsa bile) güvenilir parametre tahminleri elde etme sorununa çözüm getirmektedirler. Diğer dayanıklı regresyon tahmincileri (M, L ve R tipi tahminciler) Y yönündeki aşırı değerlere karşı direnç gösterebilirler bile, X yönündeki aşırı değerler (kaldıraç noktaları), bu tahmincileri baskı altına alabilmekte ve yönlendirebilmektedir. Yüksek bozulma sınırı ile ifade edilen husus, ele alınan tahmincinin bozulma sınırının asimptotik olarak %50 sınırına ulaşmasıdır. Bu oran bozulma sınırı için mümkün olan en iyi sonuçtur. Bundan dolayı, %50 bozulma sınırına sahip tahminciler, yüksek bozulma sınırlı tahminciler (HBE: High Breakdown Estimators) olarak adlandırılırlar. Bu tahminci grubunun öncülüğünü Rousseeuw tarafından 1984 de öne sürülen LMS (Least Medyan Squares : En küçük Medyan Kareler) ve LTS (Least Trimmed Squares : En küçük Kırpılmış Kareler) tahmincileri yapmıştır. Daha sonra, S tahminciler, MM tahminciler ve LTA (Least Trimmed Absolute Deviations : En küçük Kırpılmış Mutlak Sapmalar) tahmincileri gibi yüksek bozulma sınırına sahip diğer tahminciler önerilmiştir.

LTS tahmincisi, EKK'ye benzer bir minimizasyon işlemi ile elde edilir.

$$\min_{\hat{\theta}} \sum_{i=1}^h (r^2)_{i:n} \quad (4)$$

Burada, $(r^2)_{1:n} \leq (r^2)_{2:n} \leq \dots \leq (r^2)_{n:n}$ küçükten büyüğe doğru sıralanmış kalıntı kareleridir ve h toplama dahil edilen kalıntı karelerinin sayısıdır. Bu formülün EKK ile tek farkı en büyük kareli kalıntıların toplama işlemi dışında bırakılmasıdır. Bu nedenle, LTS, (n-h) sayıda en büyük kalıntıya sahip gözlemin aşırı değer olarak belirlenip, hesaplamanın dışında bırakılmasıyla bir dayanıklı tahmin elde etme yöntemidir. Bu, en küçük EKK amaç fonksiyonlu h birimden oluşan alt küme bulmayla eş değerdir. LTS tahmini, bu h noktaya uyan EKK tahminidir. "Görülmektedir ki, LTS tahmin sürecinde iyi noktaların büyük çoğunluğu kullanılması gerektiğinden, doğru bir tahmin elde etmek için h sayısı veri kümesindeki iyi noktaların sayısına olabildiğince yakın olmalıdır. Bu gerçekleştirildiğinde, LTS mümkün en iyi tahmini verir."¹⁸ Uygulamada, h'nin değerinin belirlenmesi oldukça zordur ve subjektiftir. Eğer varsayılandan daha fazla veri kirliliği ile karşı karşıya kalırsa LTS dayanıklılık zaafi gösterir. Bunun aksine, eğer kirlilik varsayılan düzeyden daha az ise, parametre tahminleri doğruluk bakımından sorun çıkarır ve tüm iyi noktalar hesaba katılmadığından dolayı etkinlik zaafi

¹⁸ Hakan Türkay, "Doğrusal Regresyon Modellerinin Robust (Dayanıklı) Yöntemlerle Tahmini ve Karşılaştırmalı Uygulamaları", Yayınlanmamış Doktora Tezi, İstanbul Üniversitesi, 2004, s.89.

gösterirler. “h değerine ve aşırı değerli veri yapısına bağlı olarak LTS oldukça etkin olabilir. Aslında, aşırı değerli noktaların tamamı kırılırsa, hesaplama açısından EKK’ya eşittir.”¹⁹

Daha önce yapılan simülasyon çalışmalarında, LTS oldukça iyi sonuçlar vermektedir. Örneğin, LTS, LMS ve EKK’nın karşılaştırıldığı bir simülasyon çalışmasında, normal dağılan hata terimlerinde aşırı değerlerin varlığı durumunda LTS’nin diğerlerine göre daha sapmasız parametre tahminleri ürettiği gösterilmiştir.²⁰ LTS yönteminin amaç fonksiyonu düzgün bir fonksiyondur ve buna bağlı olarak yerel etkilere karşı daha az duyarlıdır. LTS, asimptotik olarak normal dağılıma uyar ve yakınsama oranı $n^{-1/3}$ ’tür. Bu avantajlarına rağmen, LTS’nin hesaplanması oldukça zordur. Bunun dışında, LTS amaç fonksiyonunun matematiksel optimizasyonlara elverişli olmaması ve buna bağlı olarak tahminlerin hesaplanmasının zorluğu önemli sorunlar olarak görülmektedir. Ayrıca, LTS amaç fonksiyonu, sıkı bir dışlama usulü üzerine kurulmuştur. Yani, eldeki bir veri noktası tahmin sürecinde ya tamamen hesaba dahil edilir, yada tamamen sürecin dışında bırakılır. Bu durum, eğer şüpheli bölgelerde veri noktaları var ise, iyi bir yol değildir.

7. 1. LTS Tahmincisinin Hesaplanması

LTS tahmincisinin kesin hesaplamalarında, öncelikle örneklemeden h boyutlu alt örneklerin tümü çekilir. Bu, C_n^h kadar alt kümenin elde edilmesi anlamına gelir. Sonra, bu h boyutlu C_n^h tane alt kümenin her biri için regresyon doğrusu uydurulur. Buradan, amaç fonksiyonunda belirtildiği gibi, her bir örnek için hesaplanan kalıntıların en küçük h tanesinin toplamı bulunur. Bu toplamlar arasından minimumu seçilerek LTS çözümü olarak alınır. Burada karşılaşılan en büyük sorun, örneklemeden çekilecek alt kümelerin sayısının devasa boyutlara ulaşabilmesidir. Sadece küçük örneklemlerde bu yolla LTS çözümünün elde edilebileceği açıktır. Örnekleme boyutu büyüdükçe işlem yoğunluğu ve harcanan zaman çok artacak ve belli bir aşamadan sonra LTS çözümüne ulaşmak imkansız hale gelecektir.

LTS tahmincisi, söz konusu imkansızlıklar nedeniyle kesin algoritma yerine genellikle yaklaşık çözümler veren algoritma ile elde edilir. Progress programında kullanılan yaklaşık algoritma şöyle işlemektedir : Parametre sayısı p boyutlu alt örneklemler çekilerek bunların her biri için bir regresyon doğrusu uydurulur ve bu doğrular yardımıyla bulunan kalıntıların kareleri küçükten büyüğe doğru sıralanarak en küçük h tanesinin toplamı alınır. Bu şekilde elde edilen toplam değerlerinden minimumu veren regresyon doğrusu LTS tahmini olarak belirlenir.

¹⁹ Randall E.Schumacker, Michael P. Monahan, Robert E. Mount, “A Comparison of OLS to LTS and MM Robust Regression in S-PLUS”, **Southwest Educational Research Association 25th Annual Meeting**, Austin, Texas,2002,s.3.

²⁰ Latif Öztürk, “Doğrusal Regresyonda Sağlam Kestirim Yöntemleri ve Karşılaştırılmaları” Yayınlanmamış Doktora Tezi, Mimar Sinan Üniversitesi, İstanbul, 2003, s.118.

7.2. LTS Tahmincisinin Özellikleri

LTS tahmincisinin özellikleri aşağıdaki gibi özetlenebilir: ²¹

1. LTS tahmincisi için her zaman bir çözüm vardır.
2. LTS tahmincisi, regresyon eş değişim (equivariant), ölçek eş değişim ve affine eş değişim özelliklerine sahiptir.

$$\sum_{i=1}^h [(Y_i + \mathbf{x}_i \mathbf{V}) - \mathbf{x}_i (\boldsymbol{\theta} + \mathbf{V})]_{i:n}^2 = \sum_{i=1}^h [(Y_i - \mathbf{x}_i \boldsymbol{\theta})^2]_{i:n} \quad (6)$$

$$\sum_{i=1}^h [cY_i + \mathbf{x}_i (c\boldsymbol{\theta})]_{i:n}^2 = c^2 \sum_{i=1}^h [(Y_i - \mathbf{x}_i \boldsymbol{\theta})^2]_{i:n} \quad (7)$$

$$\sum_{i=1}^h [Y_i - (\mathbf{x}_i \mathbf{A})(\boldsymbol{\theta} \mathbf{A}^{-1})]_{i:n}^2 = \sum_{i=1}^h [(Y_i - \mathbf{x}_i \boldsymbol{\theta})^2]_{i:n} \quad (8)$$

Eğer $p > 1$, $h = \lfloor n/2 \rfloor + \lfloor (p+1)/2 \rfloor$ ve gözlemler genel durumda ise, (yani, gözlemlerin p tanesi $\boldsymbol{\theta}$ için tek bir çözüm veriyorsa) LTS yönteminin bozulma sınırı;

$$\varepsilon^* = \frac{\lfloor (n-p)/2 \rfloor + 1}{n} \quad (9)$$

şeklinde tanımlanır. “Daha büyük h değerleri için, LTS’nin bozulma sınırı, $\varepsilon^* \cong (n-h)/n$ ’dir. Örneğin, $h \cong 0,75n$ seçildiğinde yöntemin bozulma sınırı $\varepsilon^* \cong 0,25$ olacaktır.” ²² LTS yöntemi, LMS’de olduğu gibi verilen şartlar altında %50 bozulma sınırına ulaşmaktadır. Genel olarak, h kırpm oranı, α ’ya bağlanır. Örneğin, $h = \lfloor n(1-\alpha) \rfloor + 1$ olduğunda bozulma sınırına $\varepsilon^* = \alpha$ ’dır.

3. Eğer $p > 1$ ve gözlemler genel durumda ve $y_i = X_i \boldsymbol{\theta}$ ’yi sağlayan gözlemlerin $(n+p-1)/2$ den kesinlikle daha çok olan bazı $\boldsymbol{\theta}$ ’lar var ise, o zaman LTS çözümü $\boldsymbol{\theta}$ ’ya eşit olur.
4. LTS tahmincisi, normal dağılım durumunda da aynı asimptotik etkinliğe sahiptir. M tahmincileri gibi,

²¹ Janne Heikkilä, **Robust Regression**, çevrimiçi: <http://www.ee.oulu.fi/~jth/robust.pdf>, 05.01.2003.

²² Mia Hubert, Peter J. Rousseeuw, Stefan Van Aelst, “Robustness”, *Encyclopedia of Actuarial Sciences*, Edited by B. Sundt, J. Teugels, John Wiley, New York, 2003, s.6.

$$\psi(t) = \begin{cases} t, & |t| < \Phi^{-1}(1 - \alpha/2) \text{ ise} \\ 0, & \text{diğer durumda} \end{cases} \quad (10)$$

şeklinde tanımlanır. Bu tanım, Huber-tipi sıçrayan (skipped) ortalama ile aynıdır.

6. LTS yönteminin asimptotik etkinliği, eğer $\alpha = 0,5$ ($\varepsilon^* = \% 50$) ise sadece %7,2 olarak bulunur. Ancak, h değeri n değerine yaklaştıkça, LTS'nin etkinliği EKK'nın etkinliğine yaklaşacaktır.

8. UYGULAMA

Uygulamada, hata terimlerinin binom dağılıma uyduğu durumda EKK ve robust LTS regresyon tahmincilerinin simülasyon yöntemiyle karşılaştırılması yapılacaktır. Burada amaç; EKK ve robust LTS yöntemlerinden elde edilen parametre tahminlerinin simülasyon sürecinde başlangıç değeri olarak verilen parametre değerlerinden ne kadar sapma gösterdiklerinin ölçülmesidir. Bir başka ifadeyle amaç, Binom dağılımlı hata terimleri durumunda söz konusu regresyon yöntemlerinin hangisinin daha az sapmalı tahminler ürettiğinin gösterilmesidir. Bu amaca yönelik olarak, simülasyon sürecinde kullanılacak olan veriler hipotetik olarak türetilmiştir. Elde edilen veriler ve bu verilerin türetilme şekli aşağıda gösterilmiştir:

1. Düzgün dağılımdan 0 ile 1 aralığında iki tane (r_1 ve r_2) rassal sayının elde edilmesi

2. Elde edilen bu rassal sayılar Z olarak tanımladığımız aşağıdaki eşitlikte yerine koyularak bir $N(0,1)$ dağılmış değişken değeri elde edilir.

$$Z = \sqrt{-2 \ln(r_1)} \cos(2\pi r_2) \quad (11)$$

3. $N(0,1)$ dağılmış söz konusu Z değeri; $U = \mu + \sigma Z$ eşitliğinde yerine konularak U rassal değişkeni elde edilebilir.²³

Bu açıklamalar ışığında elde edilen $N(10, 6)$ dağılmış veriler aşağıdaki gibidir.

²³ An Introduction to Random Number Generation and Simulation, <http://dmawww.epfl.ch/benarous/Pmmi/interactive/rng7.htm>, 20.04.2004.

Tablo 1: Hipotetik olarak elde edilen yaklaşık olarak N(10,6) dağılmış veriler.

12,16774	11,68551	8,67988	10,70425	7,054214	13,56055	13,16472
13,93235	7,90377	10,30552	10,14960	9,86481	15,85245	11,27039
7,95474	10,92810	9,52312	5,24065	11,01156	5,81505	8,91277
10,47066	3,98880	13,67944	10,46852	7,994942	9,34246	11,91694
12,54917	9,75686	9,84980	9,04269	11,05708	6,54892	11,62940

Min: 3,9888 Ortalama: 10,113642 Medyan: 10,305520
Std, Sapma: 2,575998

8.1. Binom Dağılıma Uyan Hata Teriminin Elde Edilmesi

X, p parametresi ile bir Bernoulli rassal değişkeni olduğunda; X'in dağılım fonksiyonu,

$$F(x) = \begin{cases} 0, & x < 0 \\ 1 - p, & 0 \leq x < 1 \\ 1, & 1 \leq x \end{cases} \quad (12)$$

şeklinde verilmektedir. Dolayısıyla,

$$Fx^{-1}(u) = \begin{cases} 1, & 1 - p < u \\ 0, & dd \end{cases} \quad (13)$$

şeklinde alınır ve u'nun yerine 1-u koyulursa, p parametresi ile Bernoulli rassal değişkeninin U(0,1) rassal değişkeninden eğer u < p ise 1 ve diğer durumlarda 0 alınarak elde edilir. Bernoulli rassal değişkenini elde edildikten sonra bu rassal değişken kullanılarak Binom rassal değişkeni aşağıdaki gibi elde edilir:²⁴

1. Başlangıç değeri olarak X=0 alınır.
2. N tane döngü yapılarak yukarıda elde ettiğimiz B(1,p) rassal değişkeni X değişkenine eklenerek (X=X+B(1,p)) Binom rassal değişkeni elde edilir.

8.2. Simulasyon Uygulamasının Algoritması

1. Anakütle regresyon parametre değerleri simulasyonun başlangıcında $\beta_1 = 1$ ve $\beta_2 = 2$ olarak kabul edilmiştir.
2. Örneklem büyüklüğü n=35 seçilmiş ve her gözlem için normal dağılımdan elde edilen $X \sim [N(10,6)]$ değerleri sabitlenmiştir.
3. Hata terimleri (u_i değerleri) rassal olarak binom dağılımından elde edilmiştir. p=0.1, p=0.5 ve p=0.9 için 100 ve 200 döngülü simulasyon uygulaması yapılmıştır.
4. Anakütle katsayı değerleri ve hata terimleri bilindiğine göre,
$$Y_i = \beta_1 + \beta_2 X_i + u_i \quad (14)$$

²⁴ An Introduction to Random Number Generation and Simulation ,
<http://dmawww.epfl.ch/benarous/Pmmi/interactive/rng6.htm> , 20.04.2004.

- eşitliğinin kullanımıyla \bar{Y}_i değerleri bulunmaktadır.
5. Bulunan bu Y_i değerleri daha önce sabitlenmiş X değerlerine dayandırılarak, tahmin değerleri $\hat{\beta}_1$ ve $\hat{\beta}_2$ elde edilir.
 6. Bu tahminciler her defasında aynı $\beta_1 = 1, \beta_2 = 2$, X değerlerinin ve farklı u_i değerlerinin kullanımıyla çok sayıda deneme yapılarak çok sayıda $\hat{\beta}_1$ ve $\hat{\beta}_2$ değeri elde edilir.
 7. Çok sayıda elde edilen bu $\hat{\beta}_1$ ve $\hat{\beta}_2$ değerlerinin toplanıp deneme sayısına bölünmesiyle $\bar{\hat{\beta}}_1$ ve $\bar{\hat{\beta}}_2$ değerleri elde edilir. Bu $\bar{\hat{\beta}}_1$ ve $\bar{\hat{\beta}}_2$ değerleri birinci adımda verdiğimiz $\beta_1 = 1, \beta_2 = 2$ değerlerine eşitse tahmin edicinin “sapmasız” olduğunu göstermektedir.
 8. Tablolarda verilen sapma miktarları, tahmin değerleriyle $\beta_1 = 1, \beta_2 = 2$ değerleri arasındaki farkın mutlak değeridir.
($|\bar{\hat{\beta}}_i - \beta_i|$)

Bu uygulama, normal dağılmış veriler ile bu verilere bir ve iki aşırı değer (leverage point, $X_6=11.68551$ yerine $X_6=X_6+25.00$ ve $X_{11}=8.679885$ yerine $X_{11}=X_{11}+(-20)$) eklendiği durumdaki şekliyle tekrar ele alınarak sonuçları karşılaştırılmıştır. Orijinal veriler ve aşırı değer bulunan veriler karşısında EKK ve LTS regresyon tahmincilerinden elde edilen parametre tahminlerinin sapma miktarları karşılaştırılmıştır. LTS regresyon tahmincisinde h değeri $\alpha=0,25$ alınarak,

$$h = [n(1 - \alpha)] + 1 \quad (15)$$

eşitliğinden elde edilmiştir.

Tablo 2: 100 Deneme Sonucunda EKK ve LTS Yöntemlerinden Elde Edilen Parametre Tahminlerindeki Sapma Miktarları B(N,p)

	EKK Sapma Miktarı		LTS Sapma Miktarı		Hata Terimlerinin Dağılımı B(N,p)
	Sabit ($\beta_1=1$)	Eğim ($\beta_2=2$)	Sabit ($\beta_1=1$)	Eğim ($\beta_2=2$)	
Orijinal Değerler İçin	3,65956	0,017248	3,28937	0,025903	B (35, 0.1)
	17,35927	0,010299	17,38000	0,001261	B (35, 0.5)
	31,54559	0,004457	31,36990	0,037611	B (35, 0.9)
Tek Aşırı Değer Varlığında	15,42959	1,272161	3,15577	0,026130	B (35, 0.1)
	29,42862	1,273484	17,46000	0,000630	B (35, 0.5)
	43,40972	1,271358	31,63000	0,010965	B (35, 0.9)
İki Aşırı Değer Varlığında	17,82749	1,475443	3,51336	0,000843	B (35, 0.1)
	31,97494	1,480539	17,62350	0,005900	B (35, 0.5)
	45,91887	1,478116	31,5800	0,013247	B (35, 0.9)

Tablo 2’de görüldüğü gibi 100 denemelik simulasyon sürecinde, orijinal veri kümesinde (herhangi bir aykırı değer bulunmadığı durumda) EKK ve LTS yöntemleri benzer sonuçlar vermiştir. Söz konusu iki yöntemden elde edilen sabit parametre tahminleri başlangıç değerinden büyük sapmalar gözlenirken eğim parametresi tahminlerinde başlangıç değerine oldukça yakın, küçük sapmalar görülmektedir. Eğim parametresi tahminlerinde $p=0.1$ ve 0.9 için EKK yöntemi daha küçük sapma gösterirken $p=0.5$ için LTS yöntemi daha küçük sapma göstermiştir. Ancak, tek veya iki aşırı değer bulunduğu durumlarda, EKK yönteminden elde edilen sonuçlar başlangıç değerinden oldukça uzaktır, yani EKK yönteminden elde edilen hem sabit parametre hem de eğim parametresi tahminleri büyük sapmalar göstermektedir. Buna karşılık, LTS yöntemiyle elde edilen tahminler sabit parametre için EKK yöntemine benzerlik gösterirken eğim parametresi tahmininde orijinal veri kümesinde olduğu gibi küçük sapmalar göstermektedir. Bu durum LTS yönteminin Binom dağılımlı hata terimlerinde de aykırı değerlerden fazlaca etkilenmediğinin ve EKK yönteminin ise tek bir aşırı değer bile etkisinde kalarak daha fazla sapmalı sonuçlar verdiğinin bir göstergesidir.

Tablo 3: 200 Deneme Sonucunda EKK ve LTS Yöntemlerinden Elde Edilen Parametre Tahminlerindeki Sapma Miktarları $B(N,p)$

	EKK Sapma Miktarı		LTS Sapma Miktarı		Hata Terimlerinin Dağılımı $B(N,p)$
	Sabit ($\beta_1=1$)	Eğim ($\beta_2=2$)	Sabit ($\beta_1=1$)	Eğim ($\beta_2=2$)	
Orijinal Değerler İçin	3,60578	0,012348	3,60438	0,009423	B (35, 0.1)
	17,34241	0,010109	17,89660	0,029253	B (35, 0.5)
	31,42134	0,007487	31,48000	0,000000	B (35, 0.9)
Tek Aykırı Değer Varlığında	15,49504	1,280112	3,38500	0,004794	B (35, 0.1)
	29,43881	1,281681	17,84290	0,024552	B (35, 0.5)
	43,33189	1,263859	31,64890	0,016021	B (35, 0.9)
İki Aykırı Değer Varlığında	17,79300	1,469606	3,73168	0,022984	B (35, 0.1)
	31,70959	1,460223	17,89500	0,049518	B (35, 0.5)
	45,77887	1,465779	31,53000	0,001462	B (35, 0.9)

Simulasyonda deneme sayısının 200 çıkarıldığı durumda EKK ve LTS regresyon tahmincilerinden elde edilen parametre tahminlerine ilişkin Tablo 3’de gösterilen sonuçlar incelendiğinde 100 denemeli duruma benzerlik gösterdiği gözlemlenmektedir. Yani, aykırı değer içermeyen orijinal veri kümesinde EKK ve LTS tahminleri birbirine yakınken tek veya iki aykırı değer içeren durumda eğim parametresi için LTS yönteminin EKK yöntemine göre daha sapsız tahminler verdiği görülmektedir.

9. SONUÇ

Hata terimlerinin binom dağılıma uyduğu durumda, EKK ve LTS regresyon tahmincilerinin karşılaştırıldığı simülasyon çalışması sonucunda; orijinal veri kümesinde EKK ve LTS yöntemlerinin benzer sonuçlar verdiği, aşırı değerlerle verinin kirletildiği durumda ise LTS tahmincisinin bariz bir şekilde daha az sapmalı parametre tahminleri ürettiği görülmektedir. Buna göre; hata terimlerinin binom dağılıma uyduğu durumda, eğer veri kümesi aşırı değerler içeriyorsa, regresyon parametrelerinin LTS yöntemiyle tahmini daha sağlıklı sonuçlar verecektir. Bu durum dikkate alınarak; çeşitli teşhis araçları kullanılarak bir veya daha fazla aşırı değer tespit edilirse, doğrudan LTS yöntemiyle regresyon parametrelerinin tahmin edilmesi daha uygun bir yaklaşım olacaktır. EKK yönteminden elde edilen kalıntı değerlerinin aşırı değerlerin tespiti konusunda yetersizliği de göz önüne alınırsa; hem aşırı değerlerin tespiti, hem de parametre tahminlerinin sapmasızlığı için LTS yönteminin gerekliliği ortaya çıkmaktadır.

KAYNAKÇA

1. BİLLOR, N.; HADİ, S.A.; VELLEMAN, P.F. (2000), “BACON : Blocked Adaptive Computational Efficient Outlier Nominators”, **Computational Statistics and Data Analysis**, Vol.34, s.279-298.
2. HAMPEL, F. R. (1974), “The Influence Curve and Its Role in Robust Estimation”, *Journal of the American Statistical Association*, 69.
3. HAMPEL, F.R.; RONCHETTI E.M.; ROUSSEEUW P.J.; STAHEL W.A.(1986), **Robust Statistics**, New York: John Wiley.
4. HEİKKİLÄ,J.(2003),Robust Regression, Erişim: <http://www.ee.oulu.fi/~jth/robust.pdf>, Erişim Tarihi: 05.01.2003.
5. HU, F.; HU, J. (2000), “A Note On Breakdown Theory For Bootstrap Methods”, *Statistics and Probability Letters*, 50, s.49-53.
6. HUBER, Peter J. : **Robust Statistics**, New York , John Wiley, 1981.
7. HUBERT, M.; ROUSSEEUW, P.J.; VAN AELST, S.(2003), “Robustness”, *Encyclopedia of Actuarial Sciences*, Edited by B. SUNDT, J. TEUGELS, New York: John Wiley.
8. MATOUŠEK, J.; MOUNT, D.M.; NATENYAHU, N.S.(1998), “Efficient Randomized Algorithms For The Repeated Median Line Estimator”, *Algorithmica*, 20.

9. ÖZTÜRK, L.(2003), “Doğrusal Regresyonda Sağlam Kestirim Yöntemleri ve Karşılaştırılmaları”, Yayınlanmamış Doktora Tezi, Mimar Sinan Ün.
10. ROUSSEEUW, P.J.; LEROY, A.M.(1987), **Robust Regression and Outlier Detection**, New York: John Wiley.
11. ROUSSEEUW, P.J.; VAN AELST, S. (1999), “Positive-Breakdown Robust Methods in Computer Vision” Computing Science and Statistics, Vol:31, s. 451-460.
12. ROUSSEEUW, P.J.; VAN DRIËSSEN, K.(1999), “Computing LTS Regression for Large Data Sets” (Düzenlenmiş), Technical Report, University of Antwerp.
13. SCHUMACKER, R.E.; MONAHAN, M.P.; MOUNT R.E.(2002), “A Comparison of OLS to LTS and MM Robust Regression in S-PLUS”, **Southwest Educational Research Association 25th Annual Meeting**, Austin, Texas,2002.
14. STAUDTE, R.G.; SHEATHER S.J.(1990), **Robust Estimation and Testing**, New York : John Wiley.
15. TÜRKAY, H. (2004), “Doğrusal Regresyon Modellerinin Robust (Dayanıklı) Yöntemlerle Tahmini ve Karşılaştırmalı Uygulamaları”, Yayınlanmamış Doktora Tezi, İstanbul Üniversitesi.
16. ZUO,Y.(2001), “Some Qantitative Relationships Between Two Types of Finite Sample Breakdown Point”, Statistics and Probability Letters, vol.51, s.369-375.
17. _____ An Introduction to Random Number Generation and Simulation, Erişim:<http://dmawww.epfl.ch/benarous/Pmmi/interactive/rng6.htm> , Erişim Tarihi : 20.04.2004.
18. _____ An Introduction to Random Number Generation and Simulation, Erişim:<http://dmawww.epfl.ch/benarous/Pmmi/interactive/rng7.htm> , Erişim Tarihi : 20.04.2004.