



Genelleştirilmiş Lineer Karma Modellerde Tahmin Yöntemlerinin Uygulamalı Karşılaştırılması

Comparisons of Estimation Methods in Generalized Linear Mixed Models with an Application

Tuba Koç*, Mehmet Ali Cengiz

Ondokuz Mayıs Üniversitesi, Fen Edebiyat Fakültesi, İstatistik Bölümü, Atakum, Samsun

Özet

İstatistikte, Genelleştirilmiş Lineer Karma Modeller (GLKM), lineer karma modellerin özel bir halidir. Genelleştirilmiş lineer karma modeller, lineer tahmin edicilerin sahip olduğu sabit etkilere ilave olarak rastgele etkilerin olduğu genelleştirilmiş lineer modellerin bir durumudur. En çok olabilirlik yöntemi kullanılarak böyle modelleri uydurma, bu rastgele etkiler üzerinden integral alma işlemlerini içerir. Genelde bu integraller açık bir şekilde analitik formda ifade edilemezler. En çok olabilirlik yöntemi ile farklı yaklaşım yöntemleri geliştirilmiştir. Ancak bu yöntemlerin hiçbiri olası modeller ve veri setleri için iyi özelliklere sahip değildir. Bundan dolayı, Laplace, Nümerik Quadrature veya MCMC gibi yöntemler bilgisayar kullanımının artmasıyla birlikte gelişmiştir ve bu ileri yöntemler uygulanabilir hale gelmişlerdir. Bu çalışmada bağımlı değişkenin Binary veya Poisson olduğu durumlarda parametre tahminlerinin analizi için farklı yaklaşımlar karşılaştırılmış ve "Genç Nüfusun Sorun Algılaması: Trabzon Örneği" adlı çalışma verileri kullanılmıştır.

Anahtar Sözcükler: Sabit etki, Rastgele etki, Genelleştirilmiş lineer karma modeller, SAS

Abstract

In statistics, a Generalized Linear Mixed Model (GLMM) is a particular type of mixed model. It is an extension to the generalized linear model in which the linear predictor contains random effects in addition to the usual fixed effects. Fitting such models by maximum likelihood involves integrating over these random effects. In general, these integrals cannot be expressed in analytical form. Various approximate methods have been developed, but none has good properties for all possible models and data sets. For this reason, methods involving Laplace, Numerical Quadrature or Markov Chain Monte Carlo have increased in use as increasing computing power and advances in methods have made them more practical. This study compares different approaches that are applied for interpreting the parameters in mixture experiments and measuring the effects of the components in the case of the response, which has a Binary or a Poisson distribution, with application to Trabzon Youth survey data.

Keywords: Fixed effect, Random effect, Generalized linear mixed models, SAS

1. Giriş

Verilerin analizinde yaygın olarak kullanılan istatistik değerlendirme araçlarından biri lineer regresyon modelidir. Lineer regresyon modelinde normal dağılım varsayımı önemli bir rol oynamaktadır. Bağımlı değişkeninin sayı gibi kesikli değişken olduğu durumlarda normallik varsayımı sağlanmayabilir. Bir başka durum olarak da bağımlı değişkenin ikili olduğu durumları düşünebiliriz. Bu durumlarda bağımlı değişken sürekli değildir. Ayrıca bağımlı değişkenin sürekli olduğu fakat normal dağılım göstermediği durumlar da olabilir. Bu tür verilerin analizine imkan sağlayacak modeller Genelleştirilmiş Lineer Modellerdir (GLM).

Genelleştirilmiş lineer model kavramı ilk olarak Nelder ve Wedderburn (1972) tarafından geliştirilmiştir. Daha sonraki yıllarda McCullagh ve Nelder (1989), Aitken vd. (1989), Lindsey (1997), Uusipaikka (2000), McCulloch ve Searle (2001), Dobson (2002), Myers vd. (2001), Dunteman ve Ho (2006), GLM teorisi hakkında çalışmaya devam etmişlerdir. Ayrıca Cengiz (1997, 2005), Cengiz ve Percy (2001) genelleştirilmiş lineer modellerin genel bir özetini vermektedir.

Genelleştirilmiş Lineer Karma Modeller (GLKM), bir genelleştirilmiş lineer model ile lineer tahmin edicilerin rastgele etkilerinin birleştirilmesiyle elde edilmişlerdir ve uygulamada geniş bir alana sahiptirler. "Genelleştirme" kelimesi ile bağımlı değişkenin sadece normal dağılmadığı, "Karma" kelimesiyle de modeldeki genel

*Sorumlu yazarın e-posta adresi: tuba.koc@omu.edu.tr

sabit etkilere rastgele etkilerinde eklenmesi anlatılmak istenmiştir (Işık 2011).

GLKM yöntemi Breslow ve Clayton (1993), McGilchrist (1994) ve Lee ve Nelder (1996) tarafından kullanılmıştır. Genelleştirilmiş lineer karma modeller çeşitli istatistik model sınıflarını içine alır. GLKM, lineer karma modellerin varsayımlardan kaynaklanan eksiklerinin giderilmesi için geliştirilmiş bir yöntemdir. GLKM ayrıca McCulloch ve Searle (2001), Verbeke ve Molenberghs (2005), Littel ve vd. (2005) ve Jiang (2007) tarafından detaylı incelenmiştir.

İstatistiksel bir modelin parametrelerini tahmin etme birçok istatistiksel analiz için anahtar bir adımdır. GLKM'ler için bu parametreler sabit etki ve rastgele etki parametreleridir. GLKM'lerin parametrelerini tahmin etmede iki temel yaklaşım vardır. Bunlar; olabilirlik fonksiyonuna yaklaşım ve modele yaklaşımdır. Modele yaklaşımda algoritmalar genelde Taylor serileri ile ifade edilir. Bu yaklaşımlar aynı zamanda lineerleştirme yöntemleri olarak da bilinirler. Lineerleştirme yöntemiyle ilerleyen süreç durdurma kriteri sağlanıncaya kadar devam eder. Genelleştirilmiş lineer karma modellerde lineerleştirme yöntemi, rastgele etkileri içeren modeller için Pseudo ve Penalized yarıolabilirlik (PQL) (Wolfinger ve O'Connell 1993) yöntemi kullanılır. Bu yöntemle bulunan parametre tahminleri, rastgele etkilerin en iyi yansız tahmin edicilerinin tahminleri olacaktır. GLKM parametrelerini tahmin etmek için olabilirlik fonksiyonuna çeşitli integral yaklaşım yöntemleri vardır. Bunlar, Laplace yaklaşımı, Gauss-Hermite Quadrature (GHQ) ve Markov Zinciri Monte Carlo algoritmasıdır (MCMC). Bu çalışmada, 2009 yılında "Genç Nüfusun Sorun Algı-

laması; Trabzon örneği" (Murat vd. 2009) adlı çalışmada kullanılan Trabzon merkezde yaşayan 15-24 yaş arası 1286 gençle yüz yüze görüşülerek elde edilen veri seti ele alınmıştır. Farklı bağımlı değişken yapıları için sabit ve rastgele etkiler ayrı ele alınarak yukarıda bahsedilen farklı parametre tahmin yöntemleriyle karşılaştırmalar yapılmıştır.

2. Gereç ve Yöntemler

2.1 Genelleştirilmiş Lineer Karma Model (GLKM)

GLKM genelleştirilmiş lineer modellerde sabit etkileri içeren lineer tahmin edicilere rastgele etkilerin eklenmesiyle oluşturulur. Bir genelleştirilmiş lineer karma modelde rastgele etkiler lineer tahmin edicisinin bir parçasıdır ve verinin şartlı ortalaması ile lineer tahmin edicileri, bir lineer form ile bağlıdır.

y , $(n \times 1)$ boyutlu gözlem değerleri vektörü ve Y , $(r \times 1)$ boyutlu rastgele etkiler vektörü olmak üzere; $E[y|y] = g^{-1}(X\beta + Zy)$ şeklinde yazılır. Burada $g^{-1}(\cdot)$ türevlenebilen monoton link fonksiyonunun tersi, X , $(n \times p)$ boyutlu modeldeki sabit etkilere ilişkin tasarımı matrisi, Z , $(n \times r)$ boyutlu modeldeki rastgele etkilere ilişkin tasarımı matrisidir ve rastgele etkiler için ortalaması sıfır, varyansı G olan normal dağılım gösterdikleri varsayımı söz konusudur.

GLKM parametre tahmin yöntemlerinin avantaj ve dezavantajları ve genel yazılım isimleri Çizelge 1'de özetlenmiştir (Bolker vd. 2008).

GLKM'de kullanılan çeşitli uyum istatistikleri; Akaike'nin Bilgi Kriteri (AIC), (Akaike, 1974), Akaike Bilgi Kriterinin Küçük Örneklem Yanlı Düzeltilmiş Hali

Çizelge 1: GLKM yöntemlerinin avantaj ve dezavantajları özet tablosu

Yöntem	Avantajı	Dezavantajı	Yazılım
Penalized Yarıolabilirlik	Esnektir ve yaygın olarak kullanılır.	Küçük ortalamalar veya büyük varyanslar için yanlı olduğundan olabilirlik çıkarımları için uygun değildir.	PROC GLIMMIX (SAS), GLMM(Genstat), glmmPQL, glmer®
Laplace Yaklaşımı	PQL yönteminden daha doğru sonuçlar verir.	PQL yöntemine göre daha az esnek ve daha yavaştır.	PROC GLIMMIX, glmer®, AD Model Builder HLM
Gauss-Hermite Quadrature	Laplace yönteminden daha doğru sonuçlar verir	Laplace yönteminden daha yavaş çalışır ve rastgele etkiler 2-3 ile sınırlıdır.	PROC GLIMMIX, PROC NLMIXED(SAS), glmer®, glmmML®
Markov Chain Monte Carlo	Oldukça esnektir ve keyfi sayıda rastgele etkiler için doğru sonuçlar verir.	Çok yavaştır, Bayesci yaklaşım içersinde teknik açıdan hesaplanması zordur.	WinBUGS, JAGS, MCMCpack®, AD Model Builder

(AICC) dir. AICC, Hurvich ve Tsai (1989) ve Burnham ve Anderson (1998) tarafından geliştirilmiştir. Bayesci Bilgi Kriteri (BIC), Schwarz (1978), Tutarlı Akaike Bilgi Kriteri (CAIC), Bozdoğan (1987) ve Hannan ve Quinn Bilgi Kriteri (HQIC), Hannan ve Quinn tarafından 1979 yılında geliştirilmiş bilgi kriterlerindedir. En küçük bilgi kriterlerini veren modeller en uygun modellerdir. Çizelge 2’de bilgi kriterleri için formüller verilmiştir (SAS Institute Inc. 2011).

Burada l , log olabilirlik, log sözde olabilirlik veya log yarı (quasi) olabilirliğin maksimum değerini, d , modelin boyutunu ve n , n^* verinin boyutunu gösterir. d , n ve n^* modele bağlı büyüklüklendir.

Çizelge 2: Bilgi kriterleri

Kriter	Formül
AIC	$-2l + 2d$
AICC	$-2l + 2d^{n^*} / (n^* - d - 1) - 2l$
HQIC	$-2l + 2d \log n$
BIC	$-2l + d \log n$
CAIC	$-2l + d(\log n + 1)$

Çizelge 3: Tanımlanan değişkenlerin özellikleri

Değişken	Açıklama	Değişken Türü	Değişken Düzey Değerleri
y1	“Kendinizi Trabzon da nasıl hissediyorsunuz?”	Bağımlı	1: Mutlu, 2: Mutsuz
y2	“Günlük ortalama kaç saat internet kullanıyorsunuz?”	Bağımlı	0-15
x1	“Ailenizin aylık gelir ortalaması”	Bağımsız	1: 500-den az 2: 500-750 3: 751-1000 4: 1001-1500 5: 1501-2000 6: 2001-2500 7: 2500-...
x2	“Yaş grubu”	Bağımsız	1: 1983-85 2: 1986-88 3: 1989-91
x3	“Cinsiyet”	Bağımsız	1: E 2: B
x4	“Eğitim durumu”	Bağımsız	1: İlkokul mezunu 2: Ortaokul mezunu 3: Lise mezunu 4: Üniv. mezunu 5: Lise öğrencisi 6: Üniv., yüksek lisans, Doktora ögr.

2.2 Uygulama

Uygulamada, 2009 yılında “ Genç Nüfusun Sorun Algılaması; Trabzon” örneği adlı çalışmada (Murat vd. 2009) Trabzon merkezde yaşayan 15-24 yaş arası 1286 gençle yüz yüze görüşülerek elde edilen veriler kullanılmıştır. Yapılan bu çalışmada Trabzon’da yaşayan gençlerin kendilerini nasıl tanımladıkları, aile, akraba ve arkadaş çevreleriyle olan ilişkilerinin tespit edilmesi, kötü alışkanlıklarının olup olmadığı, yakın çevreye ve dünyevi meselelere duyarlılıklarının ölçülmesi, beklentilerinin ve sorunlarının neler olduğunun tespit edilmesi, mesleki becerilerinin olup olmadığının araştırılması ve boş zamanlarını nasıl değerlendirdiklerinin belirlenmesi, ideal anne ve babadan alınması gereken önemli unsurların tespit edilmesi, kendilerini Trabzon’da nasıl hissettikleri ve gelecekte Trabzon’da yaşamayı düşünüp düşünmediklerinin tespiti, gelecekte Trabzon’un nasıl bir şehir olacağını gençler tarafından tahmin edilmesi gibi benzer hususlarda gençlerin görüşlerinin tespit edilmesi ve Trabzon belediyesine gençlere yönelik neler yapılabileceğinin sunulması amaçlanmıştır. Çalışmada, cinsiyet, yaş ve eğitim kotalarına dikkat edilirken gelir rastgele etki olarak alınmıştır. Bu çalışma için Çizelge 3’ de tanımlanan değişkenler seçilmiştir.

İlk olarak “Kendinizi Trabzon da nasıl hissediyorsunuz?” soruyla ifade edilen bağımlı değişkene (y1 ikili (binary)), açıklayıcı değişkenlerin etkisini modellemek için GLM ve GLKM modeller, daha sonrada “Günlük ortalama kaç saat internet kullanıyorsunuz?” soruyla ifade edilen bağımlı değişkene (y2 poisson), açıklayıcı değişkenlerin etkisini modellemek için GLM ve GLKM modeller uygulanmıştır. GLM yöntemi uygulanırken tüm bağımsız değişkenler sabit etkili olarak alınmış ve GLKM yöntemi uygulanırken çalışmada cinsiyet, yaş ve eğitim kotalarına dikkat edildiğinden bu bağımsız değişkenler sabit etkili olarak aylık gelir ortalaması bağımsız değişkeni de rastgele etkili olarak alınmıştır. Her iki modelin kullanılmasındaki amaç GLKM'nin gücünü ortaya koymaktır. SAS9.3 makroları yazılarak yapılan analizlerden sonra, iki farklı model için GLM ve GLKM parametre tahminleri Çizelge 4’de verilmektedir.

Çizelge 4.’den anlaşıldığı gibi “ Kendinizi Trabzon da nasıl hissediyorsunuz?”, y1 bağımlı değişkeni için x2, yaş grubu açıklayıcı değişkeninin tüm seviyeleri için %5 anlamlılık seviyesinde p değerleri 0.05’ den büyük olduğundan GLM ve GLKM yöntemleri için model de anlamlı bir katkıya sahip değillerdir. x3, cinsiyet açık-

layıcı değişkeninin birinci seviyesi GLM yöntemi için p değeri 0.05’den küçük olduğundan modelde anlamlı bir katkıya sahip iken GLKM yöntemleri için anlamlı bir katkıya sahip değildir ve x4, eğitim durumu açıklayıcı değişkenine bakıldığında bir, iki ve üçüncü seviyeleri için p değerlerinin 0.05’ den küçük, dolayısıyla GLM ve GLKM yöntemleri için modelde anlamlı bir katkıya sahip, dördüncü seviyenin GLKM yöntemlerinden Laplace ve Quadrature için anlamlı olduğu ve beşinci seviyesinde hiçbir yöntem için modele anlamlı katkı sağlamadığı görülmektedir.

“Günlük ortalama kaç saat internet kullanıyorsunuz?” , y2 bağımlı değişkeni için x2, yaş grubu açıklayıcı değişkeninin birinci seviyesi için %5 anlamlılık seviyesinde p değerleri 0.05’ den küçük olduğundan GLM ve GLKM yöntemleri için model de anlamlı bir katkıya sahip ve ikinci seviyesi tüm yöntemler için modele anlamlı bir katkıya sahip değildir. x3, cinsiyet açıklayıcı değişkeninin birinci seviyesi %5 anlamlılık seviyesinde tüm yöntemler için modele katkı sağlar. Ve son olarak x4, eğitim durumu açıklayıcı değişkenine bakıldığında bir, iki üç ve beşinci seviyelerinin tüm yöntemler için modele anlamlı katkı sağladığı dördüncü seviyesinin ise %5 anlamlılık

Çizelge 4: İki farklı model için bağımsız değişkenlerin her bir model yöntemine göre p değeri

			YÖNTEMLER				
			Bağımsız Değişkenler	GLM	LAPLACE	QUADRATURE	RSPL
				P Değerleri	P Değerleri	P Değerleri	P Değerleri
Bağımlı Değişkenler	y1	x2	1	0.8861	0.0264	0.0264	0.8805
			2	0.9936	0.8827	0.8827	0.9745
		x3	1	0.0053	0.9784	0.9784	0.0052
			x4	1	0.0019	0.0052	0.0052
		2		0.0001	0.0019	0.0019	0.002
		3		<.0001	0.0002	0.0002	<.0001
		4	0.6070	<.0001	<.0001	0.6317	
	5	0.1268	0.6265	0.6265	0.0985		
	y2	x2	1	0.0013	0.0006	0.0006	0.0006
			2	0.2111	0.1450	0.1450	0.1442
		x3	1	<.0001	<.0001	<.0001	<.0001
			x4	1	<.0001	<.0001	<.0001
		2		<.0001	<.0001	<.0001	<.0001
		3		<.0001	<.0001	<.0001	<.0001
4		0.3893		0.1895	0.1895	0.1880	
5	0.0013	0.0005		0.0005	0.0005		

4. Kaynaklar

- Aitken, M., Anderson, D., Francis, B., Hinde, J. 1989.** Statistical modelling in GLIM, Oxford.
- Bolker, BM., Brooks, ME., Clark, CJ., Geange, SW., Poulsen, JR., Stevens, MHH., White, JS. 2008.** Generalized linear mixed models: a practical guide for ecology and evolution. *Trends in Eco. & E.*, 24(3): 127-135.
- Bozdogan, H. 1987.** Model selection and Akaike's Information Criterion (AIC): the general theory and its analytical extensions. *Psych.*, 52: 345-370.
- Breslow, N.E., Clayton, D.G., 1993.** Approximate Inference in Generalized Linear Mixed Models. *J.A.S.A.*, 88: 9-25.
- Burnham, K.P., Anderson, D.R., 1998.** Model Selection and Multimodel Inference. A Practical Information-Theoretic Approach, Springer, Newyork.
- Cengiz, M.A., 1997.** Bivariate Logistic Regression Analysis. Technical report, the University of Salford, MCS-97-11.
- Cengiz, M.A., Percy, D.F., 2001.** Mixed multivariate generalized linear models for assessing lower- limb arterial stenoses. *Statistic. Med.*, 20: 1663-1679.
- Cengiz, M.A., 2005.** Bayesian inference for bivariate generalized linear models in diagnosing renal arterial obstruction. *Statistic. Meth.*, 2: 168-174.
- Dobson, A.J., 2002.** An Introduction to Generalized Linear Models, 2nd ed. London: Ch.&H.
- Dunteman, G.H., Ho, Moon-Ho, R. 2006.** An Introduction to Generalized Linear Models, Sage Publication 145, USA.
- Hurvich, CM., Tsai, CL. 1989.** Regression and Time Series Model Selection in Small Samples. *Biol.* 76: 297-307.
- Işık, F. 2011.** Generalized Linear Mixed Models :An Introduction for Tree Breeders and Pathologists, Fourth International Workshop on the Genetics of Host-Parasite Interactions in Forestry, USA.
- Jiang, J. 2007.** Linear and Generalized Linear Mixed Models and Their Applications, Sprinder S., 257s, Newyork.
- Lee, Y., Nelder, J.A., 1996.** Hierarchical generalized linear models (with discussion) *J. Roy. Statist. Soc., B*, 58: 619-678.
- Lindsey, JK. 1997.** Applying Generalized Linear Models, Springer Verlag Newyork, 256s.
- Littell, RC., Milliken, GA., Stroup, WW., Wolfinger, RD. 2005.** SAS System for Mixed Models, SAS Institute Inc., Cary, NC, USA.
- McCullagh, P., Nelder, JA. 1989.** Generalized Linear Models, Ch. & H., 2nd ed.
- McCulloch, CE., Searle, SR. 2001.** Generalized Linear and Mixed Effects, J. Wiley& Sons, Inc., America, 325s.
- Mcgilchrist, CA. 1994.** Estimation in Generalised Mixed Models. *J. Roy. Statist. Soc. B.*, 56: 61-69.
- Murat, N., Cengiz, MA., Terzi, Y. 2009.** Genç Nüfusun Sorun Algılaması: Trabzon Örneği. *J. Inter. Soc. Res.*, 2(7): 175-184.
- Myers, RH., Montgomery, DC., Vining, GG. 2001.** Generalized Linear Models with Applications in Engineering and the Sciences J. Willey & Sons, NewYork.
- Nelder, JA., Wedderburn, RWM. 1972.** Generalized Linear Models, *J. Roy. Statist. Soc. A*, 135, 370-384.
- Sas Institute Inc., 2011.** SAS/STAT 9.3 User's Guide, Cary, NC,USA.
- Uusipaikka, E. 2000.** Confidence intervals in generalized regressions models, CRC Press.
- Verbeke, G., Molenberghs, G. 2000.** Linear Mixed Models for Longitudinal Data, Springer S., 568s, Newyork.
- Wolfinger, R., O'Connell, M. 1993.** Generalized linear mixed models: A pseudo-likelihood approach. *J. Statist. Computat. S.*, 48: 233-243.