

II. ARAŐTIRMALAR

Handwritten notes, possibly a list or index, with some illegible text.

Handwritten notes, possibly a list or index, with some illegible text.

Handwritten notes, possibly a list or index, with some illegible text.

Handwritten notes, possibly a list or index, with some illegible text.

DİSKRİMİNANT ANALİZİ VE BUNUNLA İLGİLİ BAZI PROBLEMLER ÜZERİNDE BİR ARAŞTIRMA (1)

Aydın ÖZTÜRK (2)

ÖZET :

Hangi popülasyona ait oldukları kesinlikle bilinmeyen fertlerin ayırım ve sınıflandırılmalarını konu edinen diskriminant analizinin genel teorisi incelenmiştir. Optimum sınıflandırma kaidesi çok değişkenli popülasyonlar için izah edilerek teori kesikli dağılış gösteren popülasyonlar için de genelleştirilmiştir.

Araştırmada diskriminant analizi ile ilgili yanlış sınıflandırma ihtimalleri, kalitatif verilerle yapılan sınıflandırma, değişkenlerin seçimi ve kayıp müşahadeler gibi konular ele alınarak bunlara ait çeşitli metodların mükayesesi yapılmıştır.

Giriş :

Popülasyonun parametreleri hakkında bilgi edinmek amacıyla alınan örneklere giren fertler tek bir vasıf bakımından ölçüye vurulabilecekleri gibi, birçok hallerde aynı şey iki ve daha fazla haller için de yapılabilir. Meselâ, aralarındaki farkın önemini tesbit etmek için iki ayrı muamelenin tatbik edildiği bir bitki çeşidinde yaprak genişliği, gövde uzunluğu ve kalınlığı gibi ölçülerin alınması istenebilir. Böyle ölçülere konu olan değişkenler arasında genellikle bir ilgi mevcut olup, bunların

ayrı ayrı analiz edilmeleri yanlış bir sonuca yol açar.

Çok değişkenli istatistikte, birden fazla ölçüye vurulan fertlerin ayırımı ve sınıflandırılmaları bahis konusu olduğu zaman bunun muayyen bir kriterle göre yapılması gerekir. Bu tip problemler çok-değişkenli istatistikte diskriminant analizi adı altında topanmıştır. Teorisi kısa zamanda geliştirilen diskriminant analizi tıp, ziraat, ve taksonomi gibi konularda geniş çapta uygulama alanı bulmuştur.

(1) Prof. Dr. Şaban Karataş, Prof. Dr. Fahrettin Tosun, ve Prof. Dr. İbrahim Aksöz'den müteşekkil bir jüri tarafından kabul edilen doktora tezinin özetidir.

(2) Atatürk Üniversitesi Ziraat Fakültesi, Zootečni Bölümü Doktor Asistanı.

Dergi Komisyonuna geliş tarihi: 28.7.1973.

Literatür Özeti:

Anderson (1951), diskriminant analizi konusunda geniş açıklamalarda bulunmuştur. Fertlerin mensup ol-

madıkları popülasyonlara yanlışlıkla sınıflandırılması ihtimallerinin minimum yapılmasını kriter olarak aşağıdaki istatistiği bulmuştur.

$$W(\underline{X}) = [\underline{u}(1) \quad \underline{u}(2)]' \underline{V}^{-1} \underline{X} - \frac{1}{2} [\underline{u}(1) + \underline{u}(2)]' \underline{V}^{-1} [\underline{u}(1) \underline{u}(2)]$$

$$> \ln \frac{q_2 c_{12}}{q_1 c_{21}}$$

Formüldeki $\underline{u}(1)$, $\underline{u}(2)$ birinci ve ikinci popülasyonlara ait ortalama vektörlerini \underline{V} varyans kovaryans matrisini q_1 , q_2 popülasyonlara ait a priori ihtimalleri, c_{12} ve c_{21} ise birinci ve ikinci popülasyonlara yanlışlıkla yapılan sınıflandırmanın maliyetini göstermektedir. Bu istatistiğe göre \underline{X} vektörünün müşahade edildiği feide ait $W(\underline{X})$ değeri eşitsizliğin sağındaki değerden büyükse ilgili fert birinci, aksi halde ikinci popülasyona dahil edilir.

$\underline{u}(1)$, $\underline{u}(2)$ ve \underline{V} parametleri bilinmedikleri takdirde bunların yerine örnek tahminleri olan $\bar{\underline{x}}(1)$, $\bar{\underline{x}}(2)$ ve S değerleri kullanılır. Böylece bulunan yeni istatistik $w(\underline{X})$, sapmasız değildir. Wald (1944) ve Sitgreaves (1952) $w(\underline{X})$ in örnek dağılımını bulmuşlardır. Anderson (1958) ikiden fazla olan çok-değişkenli normal popülasyonlarda ayırım ve sınıflandırmanın nasıl yapılacağını göstermiştir.

Sürekli değişkenler için geliştirilen teori kesikli değişkenler içinde genelleştirilebilir. Eldeki verilerin multinomiyal dağılım gösterdiği farzedilirse sürekli değişkenlerdeki gibi bulunan op-

timum kriterle sınıflandırma yapılabilir (Cochran ve Hopkins, 1962). İki sıklı değişkenlerin (dichotomus) tarif ettikleri popülasyonlar için yapılan sınıflandırmada Martin ve Bradley (1972) değişik bir ihtimal modeli kullanmıştır.

Materyal ve Metod :

Araştırmada iki grup veri kullanılmıştır. Bunlardan birincisi şans sayıları kullanılarak elde edilmiştir. Ortalaması ve varyans bilinen bağımsız normal değişkenler elde edildikten sonra çok - değişkenli dağılıma geçilmiştir. Diğer gruptaki veriler ise Atatürk Üniversitesi, Öğrenci İşlerindeki kayıtlardan alınmıştır.

Bu verileri kullanmak suretiyle disriminant analizi ile ilgili bazı problemler için şimdiye kadar ileri sürülmüş olan çeşitli metodların mukayesesi yapılmıştır.

Araştırmalar ve Sonuçları :

Literatür özetinde verilen teorinin ışığı altında yapılan çalışmalarda sırasıyla, yanlış sınıflandırma ihtimalleri, kalitatif verilerle yapılan sınıflandırma,

değişkenleri seçimi ve kayıp müşahadeler gibi konular ele alınmıştır.

Hata nisbetlerini tahmin için geliştirilen metodlardan beş tanesi alınarak mukayese edilmiştir. Bunlardan $D^2 = [\bar{x}(1) - \bar{x}(2)]' S^{-1} [\bar{x}(1) - \bar{x}(2)]$ si alınarak istatistiğine dayanan iki tanesinin en iyi olduğu sonucuna varılmıştır.

Öğrenci kayıtlarından alınan kalitatif verilerle yapılan sınıflandırmada Cochran ve Hopkins metoduyla Martin ve Bradley metodunun uygulaması gösterilmiştir. Ayrıca bir normal değişkenin kalitatif şekle getirilmesinden dolayı ayırım gücünde meydana gelen değişiklik incelenmiştir. İki popülasyon arasındaki mesafe çok küçük olduğu

zaman bu değişikliğin takriben % 63,6 olduğu gösterilmiştir.

Kayıp müşahadelerin tahmini için varyans - kovaryans matrisinin izinin ve determinatının minimum yapılmasını kriter olarak alan iki metodun mukayese yapılarak bunların birbirlerine olan üstünlükleri izah edilmiştir.

Son olarak diskriminant fonksiyonu için kullanılan değişkenlerin seçiminde bazı metodlar incelenmiştir. Diskriminant fonksiyonu ile çoklu regresyon arasındaki ilişkide gösterilerek değişkenlerin seçimi için regresyon analizindeki tekniklerden birinin kullanılmasının mümkün olabileceği belirtilmiştir.

A RESEARCH ON DISCRIMINANT ANALYSIS AND SOME RELATED TOPICS

In this reseacrh the general theory of discriminant analysis and some related topics such as probality of misclassification, choosing the best variables and the estimation of the missing observations were considered.

The poblem of classification aries when an investigator makes a number of measurements on an individual and wishes to classify the individual into one of several populations on the basis

of these measurements . In costructing a precedure of calassification it is desired to minimise the probality of misclassification. For this purpose a calassification procedure was obtained by minimising the expected cost of misclassification.

Anderson proposed the following statistics for the classifictaion of the individuals into one of two multivariate normal populations.

$$W(\underline{X}) = [\underline{u}(1) - \underline{u}(2)]' \underline{V}^{-1} \underline{X} - \frac{1}{2} [\underline{u}(1) + \underline{u}(2)]' \underline{V}^{-1} [\underline{u}(1) - \underline{u}(2)]$$

where $\underline{u}(1)$, $\underline{u}(2)$ are the mean vectors and \underline{V} is the variance - covariance matrix. When the parameters $\underline{u}(1)$, $\underline{u}(2)$ and \underline{V} are unknown, the sample estimates of these are used. The above mentioned statistical formula is the same as Fisher's linear discriminant function, except the constant term inc-

luded in $W(\underline{X})$. In the case where there are more than two populations, the allocation rule was generalised. The same rule was used to show how to classify individuals with the multivariate qualitative data.

The data used in the research was obtained by inkoving the simula-

tion techniques and gathering observations on the student records.

Comparing five of the methods developed to estimate the probability of misclassification, two of them were found to be the best.

Application of two different was illustrated for the classification with the multivariate qualitative data. The effect of replacing a normal variate by a qualitative variate was also mentioned for a special case.

The advantages and disadvantages of the two methods proposed to

estimate the missing observations were discussed.

In the last part of the research some which are used to select the best variables for discriminant function were considered. The relation between the multiple regression equation and the discriminant function was illustrated, and the possibility of the use of one of the regression techniques was pointed out for selecting the best variables of the discriminant function.

LITERATÜR

Anderson, T. W. (1951), Classification By Multivariate Analysis. *Psychometrika*, 16, 31.

Anderson, T. W. (1958), An introduction to Multivariate Statistical Analysis. John Wiley and Sons Inc. USA.

Cochran, W. G. ve Hopkins, C. E. (1961). Some Classification Problems With multivariate Qualitative Data. *Biometrics*. 17, 10.

Martin, D. C. ve Bradley, R. A. (1972), Probability Models. Estimation, and

Classification For Multivariate Dichotomus Populations. *Biometrics*, 28, 203.

Sitgreaves, R. (1952), On the Distribution of Two Random Matrices Used in Classification Procedures. *Annals of Math. Statist.* 23,263.

Wald, A. (1944) On a Statistical Problem Arising in the Classification of an Individual Into One of Two Groups. *Annals of Math. Statist.* 15,145.