

**T.C.  
MARMARA ÜNİVERSİTESİ  
FEN BİLİMLERİ ENSTİTÜSÜ**

**NONLİNEER LOJİSTİK REGRESYON VE  
UYGULAMASI**

**Esra Zeynep ŞENSOY**

**YÜKSEK LİSANS TEZİ  
MATEMATİK ANABİLİM DALI  
UYGULAMALI MATEMATİK PROGRAMI**

**DANIŞMAN  
Prof. Dr. Müjgan TEZ**

**İKİNCİ DANIŞMAN  
Prof. Dr. Aydın ERAR**

**İSTANBUL 2009**

**T.C.  
MARMARA ÜNİVERSİTESİ  
FEN BİLİMLERİ ENSTİTÜSÜ**

**NONLİNEER LOJİSTİK REGRESYON VE  
UYGULAMASI**

**Esra Zeynep ŞENSOY**

**YÜKSEK LİSANS TEZİ  
MATEMATİK ANABİLİM DALI  
UYGULAMALI MATEMATİK PROGRAMI**

**DANIŞMAN  
Prof. Dr. Müjgan TEZ**


**İKİNCİ DANIŞMAN  
Prof. Dr. Aydın ERAR**


**İSTANBUL 2009**


T.C.  
MARMARA ÜNİVERSİTESİ  
FEN BİLİMLERİ ENSTİTÜSÜ

## KABUL ve ONAY BELGESİ

Esra Zeynep ŞENSOY'un **Nonlinear Lojistik Regresyon ve Uygulaması** başlıklı Lisansüstü tez çalışması, M.Ü. Fen Bilimleri Enstitüsü Yönetim Kurulu'nun 29.06.2009 tarih ve 2009/15-09 sayılı kararı ile oluşturulan jüri tarafından Matematik Anabilim Dalı Uygulamalı Matematik Programında YÜKSEK LİSANS Tezi olarak kabul edilmiştir.

Danışman: Prof.Dr. Müjgan TEZ (Marmara Üniv.).....

1. Üye : Yard.Doç.Dr. Birsen E. ERDOĞAN (Marmara Üniv.).....

2. Üye : Yard.Doç.Dr. Dursun ÜSTÜNDAĞ (Marmara Üniv.).....

Tezin Savunulduğu Tarih : 22/07/2009

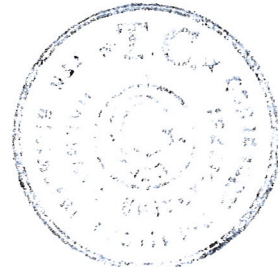
## ONAY

M.Ü. Fen Bilimleri Enstitüsü Yönetim Kurulu'nun 28.09.2009 tarih ve 2009/23-32 sayılı kararı ile Esra Zeynep ŞENSOY'un Matematik Anabilim Dalı Uygulamalı Matematik Programında Y.Lisans (MSc.) derecesi alması onanmıştır.

Fen Bilimleri Enstitüsü Müdürü

*Prof. Dr. Sevil ÜNAL*





# TEŐEKKÜR

Çalıőmalarımı yakından izleyerek beni yönlendiren ve her türlü yardım ve desteęi saęlayan danıőman hocalarım sayın Prof. Dr. Müjgan TEZ ve sayın Prof. Dr. Aydın ERAR'a, uygulama çalıőmasında istatistiksel program kullanımında her daim yardımcı olan Araő. Gör. Barıő AŐIKGİL'e, benden hiçbir zaman desteęini esirgemeyen sevgili aileme ve yüksek lisans arkadaőım Bürya ÖZKARTAL'a en içten dileklerle teşekkür etmeyi bir borç bilirim.

# İÇİNDEKİLER

	SAYFA
TEŞEKKÜR.....	i
İÇİNDEKİLER .....	i
ÖZET .....	iv
ABSTRACT .....	vi
SEMBOLLER .....	viii
KISALTMALAR.....	x
ŞEKİL LİSTESİ.....	xi
TABLO LİSTESİ .....	xii
<b>BÖLÜM I. GİRİŞ.....</b>	<b>1</b>
I.1 LOJİSTİK REGRESYON MODELİ .....	1
I.2 LOJİSTİK REGRESYONUN LİNEER REGRESYONDAN FARKI ve KULLANIM AMAÇLARI.....	5
<b>BÖLÜM II. LOJİSTİK REGRESYON MODELİNİN KURULMASI VE ANALİZİ .....</b>	<b>7</b>
II.1 PARAMETRE TAHMİN YÖNTEMLERİ .....	7
II.1.1 En Çok Olabilirlik Yöntemi.....	7
II.1.2 Yeniden Ağırlıklandırılmış İteratif En Küçük Kareler Yöntemi ....	10
II.1.3 Minimum Logit Ki-Kare Yöntemi.....	11
II.2 PARAMETRELERİN ÖNEM TESTİ .....	11
II.3 ÇOKLU LOJİSTİK REGRESYON MODELİ .....	13
II.4 LOJİSTİK REGRESYON MODELİNDE KATSAYILARIN YORUMLANMASI .....	17
III.4.1 İki Sonuçlu (Dichotomous) Bağımsız Değişken .....	17
III.4.2 Çok Sonuçlu (Polychotomous) Bağımsız Değişken .....	20
III.4.3 Sürekli Bağımsız Değişken .....	21
III.4.4 Etkileşim Varlığında Odds Oranlarının Kestirimi.....	22
II.5 MODEL UYUMUNUN BELİRLENMESİ.....	24
II.5.1 Ki-Kare İstatistiği ( $X^2_{B0}$ ) .....	25
II.5.2 Pearson Ki-Kare İstatistiği ve Sapma Ölçütü.....	25
II.5.3 Hosmer-Lemeshow (G) İstatistiği .....	27

II.5.4 Blok Ki-Kare İstatistiđi .....	28
II.5.5 Sınıflandırma Tablosu .....	28
II.5.6 ROC Eğrisi Altındaki Alan .....	29
II.6 MODEL DEĞERLENDİRMESİ.....	30
<b>BÖLÜM III. DOĞRUSAL OLMAYAN LOJİSTİK REGRESYON</b>	
<b>MODELİ .....</b>	<b>34</b>
III.1 SPLAYNLAR VE PARÇALI POLİNOM UYDURMASI.....	35
<b>BÖLÜM IV. UYGULAMA .....</b>	<b>46</b>
IV.1 UYGULAMA-1 .....	46
IV.1 UYGULAMA-2.....	55
<b>BÖLÜM V. SONUÇLAR VE TARTIŞMA .....</b>	<b>66</b>
<b>KAYNAKLAR.....</b>	<b>68</b>
<b>EKLER.....</b>	<b>71</b>
<b>ÖZGEÇMİŞ.....</b>	<b>77</b>

# ÖZET

## NONLİNEER LOJİSTİK REGRESYON VE UYGULAMASI

Lojistik Regresyon, bağımsız değişkenlerin bağımlı değişkenler üzerindeki etkisini olasılık olarak ortaya koyar. Risk faktörlerinin olasılık olarak belirlenmesini sağlar. Lojistik Regresyonun amacı, en az değişkeni kullanarak en iyi uyuma sahip olacak şekilde yanıt değişkeni ile bağımsız değişkenler arasındaki ilişkiyi tanımlayabilmek ve amaca yönelik kabul edilebilir model kurmaktır. Lojistik Regresyonda neden sonuç ilişkisinin ortaya konulması amacıyla, bağımlı değişken olumlu-olumsuz, başarılı-başarısız gibi kategorik olarak ikili (binary) kodlanmıştır. Bu yöntem, sayısal verilerle kolay yorumlanabildiği için popülerdir. Bu nedenlerle son zamanlarda epidemiyolojik çalışmalarda, biyolojide, ziraatte, taşımacılıkta, ekonomi gibi alanlarda yaygın bir şekilde kullanılmaktadır.

Lojistik Regresyon modelinde lojistik fonksiyonun doğal logaritmasının alınmasıyla elde edilen lojit fonksiyon her zaman doğrusal olmayabilir. Bazen verinin doğrusallığa uygun olmayışından bazen de model denkleminin istenilen sonucu iyi karşılayamamasından kaynaklanan bu durum model denkleminin karesel ya da kübik olarak biçimlendirilmesiyle çözümlenebilir. Bu durumda splayn fonsiyonlarla lojistik regresyon çalışılacaktır.

Bu çalışmada, lojistik regresyon modelinin kurulması ve analizi başlığı altında parametre tahmin yöntemleri, parametrelerin önem testi açıklanmıştır. Modelde bağımsız değişken sayısının birden fazla olması durumunda ‘Çoklu Lojistik Regresyon Model incelenmiştir. Lojistik regresyon için katsayıların yorumlanması stratejileri açıklanmıştır.

Lojistik regresyon fonksiyonu doğrusal olmadığında kullanılacak yöntemler hakkında bilgiler verilmiştir. Buna göre çalışmamızda yer alan segmentli fonsiyonlardan bahsedilerek splayn fonsiyonlar açıklanmıştır.

Uygulama çalışmamızda, 243 kişilik veri kümesinde hemoglobin kan değeri incelenerek ilgili risk durumu üzerine Doğrusal Olmayan Lojistik Regresyon Modeli

kurulmuştur. SAS version 9.1 kullanılarak model analiz edilmiştir. Ayrıca segmentli lojistik regresyona bir başka örnek olacak Mulla'nın albümin çalışması Matlab programına uyarlanarak incelenmiştir.

Son olarak her iki uygulama çalışmasının sonuçları, farklı istatistiksel programlarda değerlendirilmiştir.

**Haziran, 2009**

**Esra Zeynep ŞENSOY**

# **ABSTRACT**

## **NONLINEAR LOGISTIC REGRESSION AND ITS APPLICATION**

Logistic Regression produces effects of independent variables on dependent variables as probability. Risk factors could be determined as probability by Logistic Regression. The aim of Logistic Regression Analysis is to establish the best acceptable model with least variable, which gives the relationship between outcome and independent variables. In Logistic Regression, dependent variable is coded binary as positive-negative, successful-unsuccessful for showing dose-response relationship. This method is popular because of easy interpretation by numerical data. Thus it is commonly used in many fields including business and finance, ecology, health policy, agriculture, biology and transportation.

Logit function obtained by computing the natural log of the logistic function in Logistic Regression model can not always be linear. Sometimes data is not available as linear. Sometimes model equation can not satisfy the probable cause. This situation can be analyzed as quadratic or cubic shape. At this stage, logistic regression is studied with spline functions.

At this study, the title under building logistic regression models and analysis, parameter estimation method and test of significance of parameters are explained. When model has more than one variable, ‘Multiple Logistic Regression Model’ is investigated. For logistic regression, strategies for interpretation coefficients are explained.

The information is given about that the method can be used, when Logistic Regression isn’t linear. According to that, spline functions are explained to talk about segmented functions in our study.

In our application work, blood hemoglobin values in the data set of 243 people were examined. The Nonlinear Logistic Regression Model was established based on the risk status. The model is analyzed to use SAS version 9.1.

Additionally, Mulla's albumin study was adapted to Matlab programme as an example of segmented logistic regression.

Finally, the results of both application studies are evaluated by different statistical programmes.

**June, 2009**

**Esra Zeynep ŞENSOY**

$y$	: Bağımlı Değişken (Yanıt Değişkeni )
$y_i$	: i.Deneğin Bağımlı Değişkeninin Değeri
$\pi(x_i)$	: Lojistik Fonksiyon ( i. gözleme ait lojistik fonksiyon değeri )
$\hat{\pi}(x_i)$	: Lojistik Fonksiyon Kestirimi ( i. gözlem için )
$\beta_0$	: Herhangi Bir Parametre
$\beta_1$	: Eğim Parametresi
$\beta^l$	: Parametreler vektörü
$\hat{\beta}_0$	: $\beta_0$ Parametresinin Kestirimi
$\chi^2$	: Ki-kare İstatistiği
<b>100(1-<math>\alpha</math>)%CI</b>	: Güven Aralığı

## KISALTMALAR

<b>LR</b>	: Lojistik Regresyon
<b>Var</b>	: Varyans
<b>HGB</b>	: Hemoglobin
<b>EKK</b>	: En Küçük Kareler Kestirimi
<b>ML</b>	: En Çok Olabilirlik ( Maksimum Likelihood )
<b>IRLS</b>	: İteratif Olarak Ağırlıklandırılmış En Küçük Kareler

## SEMBOLLER

<b>D</b>	: Sapma ( Deviance ) İstatistiđi
<b>E(y/x) = <math>\pi</math> (x)</b>	: Koşullu Ortalama ( x değeri için y'nin beklenen değeri )
<b>f(z)</b>	: Lojistik Fonksiyon ( basit form )
<b>g(x)</b>	: Lojit Fonksiyon
<b>G</b>	: Hosmer-Lemeshow İstatistiđi
<b>h<sub>i</sub></b>	:Leverage Uzaklık Ölçütü
<b>I(<math>\beta</math>)</b>	: Fisher Bilgi Matrisi
<b>J</b>	: Jacobiyen Matris
<b>J<sup>T</sup></b>	: Jacobiyen Matrisinin Transpozu
<b>L(.)</b>	: Olabilirlik Fonksiyonunun Deđeri
<b>OR</b>	: Odds Oranı
<b><math>\widehat{OR}</math></b>	: Odds Oran Kestirimi
<b>P</b>	: Olasılık Fonksiyonu
<b>Pr</b>	:Anlamlılık Test Ölçütü
<b>ROC</b>	: Receiver Operating Characteristic
<b>SE</b>	: Standart Hata
<b><math>\widehat{SE}</math></b>	: Standart Hata Kestirimi
<b>SSE</b>	: Artık Kareler Toplamı
<b>SSR</b>	: Regresyon Kareler Toplamı
<b>ST</b>	: Skor Testi
<b>n</b>	: Gözlem Sayısı
<b>V(.)</b>	: Varyans
<b>Var(<math>\beta</math>)</b>	: $\beta$ Parametresinin Varyansı
<b>w<sub>i</sub></b>	: i için Ağırlık Fonksiyonu
<b>W</b>	: Wald İstatistiđi
<b><math>\varepsilon</math></b>	: Hata
<b>x</b>	: Bağımsız Deđişken
<b>x<sub>i</sub></b>	: i.Deneđin Bağımsız Deđişkeninin Deđeri

# ŞEKİL LİSTESİ

	<u>SAYFA NO</u>
Şekil I.1 Lojistik Fonksiyonun Grafikte Gösterilmesi.....	2
Şekil I.2 Lojistik Fonksiyondaki z Değerinin Tanım Aralığı.....	2
Şekil I.3 Lojistik Fonksiyon S-eğrisi ve Risk Değerlendirme Şekli.....	3
Şekil II.1 ROC Eğrisi Altında Kalan Alanın Grafıksel Gösterimi.....	29
Şekil II.2 Leverage Değerine(h) Karşı Kestirilen Olasılık Değerinin ( $\hat{\pi}$ ) Çizimi.....	32
Şekil III.1 Adım Fonksiyonu ve Kırık Çizgi. 0. ve 1.Dereceden Splayn Yaklaşımıyla İki Düzgün Eğri.....	35
Şekil III.2 Bir Parabolik Splayn. 0.Dereceden Bir Splayn Histogramının Her Bir Parçasındaki Alana Eşit Olan Parabolik Splayn.....	36
Şekil III.3 Kesilmiş Üslü İfade. $t-2 \leq x \leq t+2$ için $6(x-t)_+^0$ , $6(x-t)_+$ , $3(x-t)_+^2$ ve $(x-t)_+^3$ fonksiyon çizimleri.....	38
Şekil III.4 Lineer- Karesel Lojistik Regresyon Modeli.....	41
Şekil III.5 Lineer- Karesel Lojistik Regresyon Modelin Türevi .....	41
Şekil IV.1 Serum Albüminin Risk Grafığı.....	62

# TABLO LİSTESİ

## SAYFA NO

<b>Tablo II.1</b> Bağımsız Değişkenin İkili Olması Durumunda Lojistik Regresyon Modelinin Değerleri.....	17
<b>Tablo IV.1</b> HGB Modelleri için Response Profili.....	50
<b>Tablo IV. 2</b> HGB Model_1 için Model Uygunluk İstatistiği.....	50
<b>Tablo IV. 3</b> HGB Model_2 için Model Uygunluk İstatistiği.....	50
<b>Tablo IV. 4</b> HGB Model_1 için Sıfır Hipotez Testi.....	51
<b>Tablo IV. 5</b> HGB Model_2 için Sıfır Hipotez Testi.....	51
<b>Tablo IV. 6</b> HGB Model_1 için En Çok Olabilirlik Kestirimi.....	52
<b>Tablo IV. 7</b> HGB Model_2 için En Çok Olabilirlik Kestirimi.....	52
<b>Tablo IV. 8</b> HGB Model_1 için Odds Ratio Kestirimi.....	52
<b>Tablo IV. 9</b> HGB Model_2 için Odds Ratio Kestirimi.....	53

# BÖLÜM I

## GİRİŞ

Regresyon analizi bir bağımlı değişken ile bir veya birden fazla bağımsız değişken arasındaki bağıntının açıklanması sürecidir.

Regresyon analizinde bağımsız değişkenler (X), nicel veya nitel olabilir. Bağımlı değişkenin nicel veya nitel olması, kullanılan yöntemin çözümü ve yorumlanmasını etkiler. Bağımlı değişkenin nitel olması durumunda Lojistik Regresyon Analizi uygun bir yöntemdir.

Lojistik Regresyon Analizinin kullanım amacı, İstatistik'teki diğer model yapılandırma teknikleri ile aynıdır: En az değişkeni kullanarak en iyi uyuma sahip olacak şekilde bağımlı ile bağımsız değişkenler arasındaki ilişkiyi tanımlamaktır [3].

### I.1 LOJİSTİK REGRESYON MODELİ

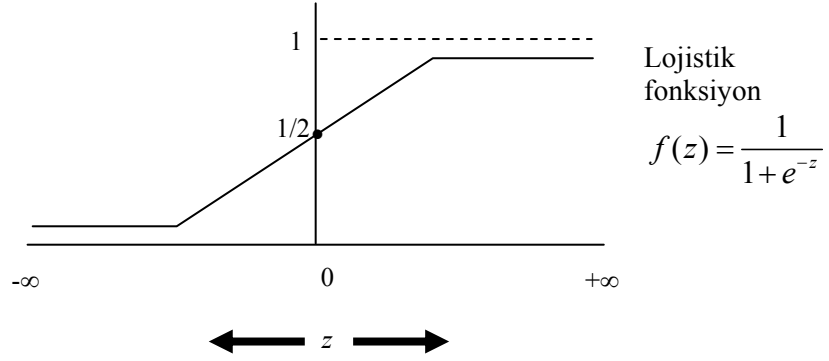
Lojistik Regresyon yöntemi bir veya birden fazla açıklayıcı değişken ile yanıt değişkeni arasındaki ilişkileri inceler. Yanıt değişkeni kesikli olup, iki veya daha fazla olası değere bağlıdır. Bu yöntemde, açıklayıcı değişkenlerin bağımlı değişkenler üzerindeki etkileri olasılık olarak hesaplanarak risk faktörlerinin olasılık olarak belirlenmesi sağlanır [15, 18].

Alışlagelen doğrusal regresyon analizinde bağımlı değişkenin değeri kestirilirken, lojistik regresyon analizinde bağımlı değişkenin alacağı değerlerden birinin gerçekleşme olasılığı kestirilir.

Lojistik modelin matematiksel formu,

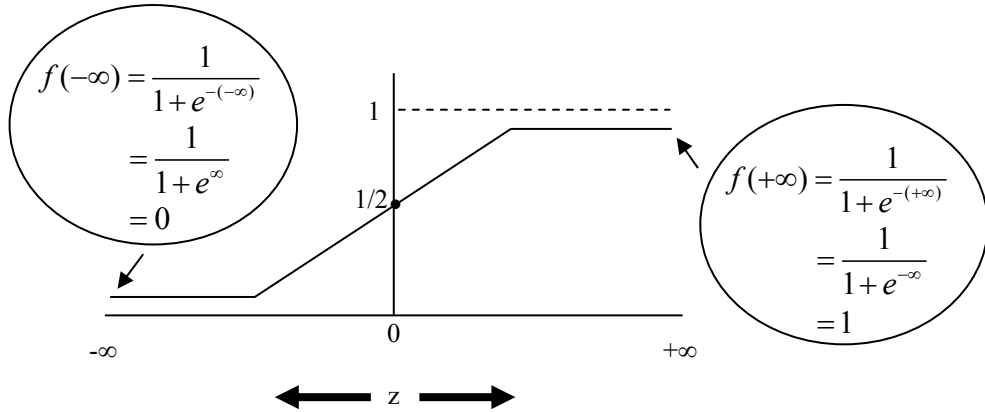
$$f(z) = \frac{1}{1 + e^{-z}}$$

şeklinde tanımlanır.



Şekil I.1 Lojistik Fonksiyonun Grafikte Gösterilmesi

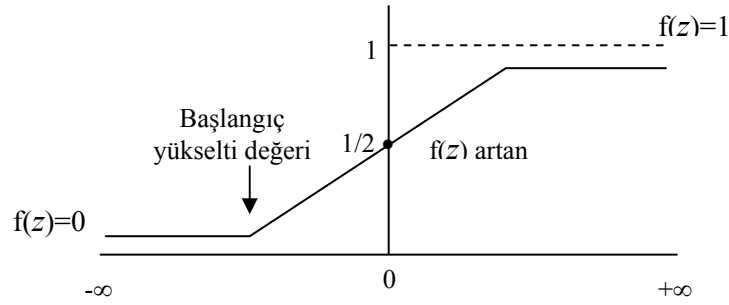
Şekil I.1 ve Şekil I.2'den görüldüğü gibi  $z$ 'nin tanım aralığı  $(-\infty, +\infty)$  dir.  $z$ 'nin değeri ne olursa olsun  $f(z)$  fonksiyonunun değeri 0 ile 1 arasında kalır.



Şekil I.2 Lojistik Fonksiyondaki  $z$  Değerinin Tanım Aralığı

$f(z)$  lojistik fonksiyonunun 0 ile 1 arasında bir değişime sahip olması lojistik fonksiyonun tercih edilmesindeki önemli bir nedendir; çünkü model 0 ile 1 arasında yer alan olasılıklar üzerine kurulmuştur. Böylece 0-1 aralığı dışında lojistik model için bir risk kestiricisi ortaya çıkmaz.

Lojistik modelin tercih edilmesinin diğ er bir nedeni de lojistik fonksiyonun biçimidir. Ş ekil I.3'te görü ldü ğ ü gibi lojistik fonksiyon S-ş ekinde bir sigmoid eğ ri meydana getirir. S ş ekilli lojistik fonksiyonda  $z$ , ç eş itli risk faktörlerinin katılımını gösteren bir indeks olarak kabul edilirse  $f(z)$  de  $z$  de ğ erindeki riski gösterir. Ş ekil I.3'den de anlaş ılaca ğ ı gibi yükselti de ğ erine kadar bireyin riski minimumdur. Sonra risk ortadaki  $z$  de ğ erlerinde hızla artmakta ve  $z$  yeteri kadar arttı ğ ında 1 civarında kalmaktadır.



**Ş ekil I.3** Lojistik Fonksiyon S-eğ risi ve Risk De ğ erlendirme Ş ekli

Yükselti de ğ eri, epidemiologlar tarafından örne ğ in, hastalık koş ullarının de ğ iş ikli ğ ini belirtmek için kullanılmış tır [11].

Lojistik modelde yanıt de ğ iş keninin kesiklili ğ i nedeni ile yanıt de ğ iş keninin açıklayıcı de ğ iş ken ile iliş kisi X-Y çiziminden açıkça görülemez. Bunun için yanıt de ğ iş kenini yerine  $f(z)$  olasılık de ğ erine karşı çizim yapılır;  $f(z)$  fonksiyonu sürekli olup Ş ekil I.3'teki gibidir [2].

$k$  sayıda bağımsız (açıklayıcı) de ğ iş ken varlı ğ ında lojistik fonksiyon,

$$z = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$$

ve  $\beta_0$  ve  $\beta_i$  bilinmeyen parametreler olmak üzere,

$$f(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + e^{-(\beta_0 + \sum \beta_i X_i)}}$$

biçiminde yazılır.

Öte yandan lojistik regresyon fonksiyonu  $x$ 'deki her birim de ğ iş me sonucunda  $E(Y/x)$ 'de oluşan de ğ iş ikli ğ i gösterir. Böylece lojistik regresyon modeli  $\pi(x) = E(Y/x)$  iken,

$$\pi(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} \quad (I.1)$$

$$\pi(x)/(1 - \pi(x))$$

biçiminde verilir [8].

Koşullu ortalama "lojit dönüşüm" adı verilen,

$$g(x) = \ln \left[ \frac{\pi(x)}{1 - \pi(x)} \right] = \ln \left[ \frac{P(y=1|x)}{P(y=0|x)} \right] = \beta_0 + \beta_1 x \quad (I.2)$$

dönüşümüyle doğrusal biçime getirilebilir. Eşitlikteki  $\pi(x)/(1 - \pi(x))$  oranı odds olarak adlandırılır.

Bu dönüşümün önemi,  $g(x)$ 'in diğer regresyon modellerinde istenen çoğu özelliği taşımasıdır. Lojit  $g(x)$ , parametreleri bakımından doğrusaldır. Lojit  $g(x)$  süreklidir ve  $x$ 'in tanım aralığına göre 0 ve 1 arasında değer alabilir.

Lojistik modelin yanıt değişkenini oluşturan lojit dönüşüm  $g(x)$ 'in özellikleri aşağıdaki şekilde sıralanabilir:

- i)  $\pi(x)$  arttıkça  $g(x)$  de artar.
- ii)  $\pi(x)$  yani lojistik fonksiyon 0 ve 1 arasında iken  $g(x)$  tüm gerçel doğru üzerinde değerler alır.
- iii) Eğer  $\pi(x) < 0,5$  ise  $g(x) < 0$

Eğer  $\pi(x) > 0,5$  ise  $g(x) > 0$ 'dır [2].

Lojistik regresyon modelinde yanıt değişkeninin gözlem değeri  $y = E(Y/x) + \varepsilon$  ile gösterildiğinde,  $\varepsilon$  hata terimi olarak adlandırılır ve gözlemin koşullu ortalamadan sapma miktarını ifade eder.  $\varepsilon$ 'nin klasik modelde sıfır ortalama ve  $\sigma^2$  sabit varyansı ile normal dağılıma sahip olması önemli bir varsayımdır. Fakat iki sonuçlu bağımlı değişken için,  $x$  verildiğinde yanıt değişkeni  $y = \pi(x) + \varepsilon$  ile gösterilir ve  $\varepsilon$ 'un bir veya iki olası değeri vardır. Eğer,

$$y=1 \text{ ise } \pi(x) \text{ olasılığıyla, } \varepsilon = 1 - \pi(x),$$

$$y=0 \text{ ise } \pi(x) \text{ olasılığıyla, } \varepsilon = -\pi(x) \text{ olur. Böylece } \varepsilon, \text{ sıfır ortalama ve}$$

$\pi(x)[1 - \pi(x)]$  varyansına sahip bir dağılım gösterir:

$$E(\varepsilon) = [(1 - \pi(x))\pi(x)] + [-\pi(x)(1 - \pi(x))] = 0$$

$$V(\varepsilon) = E(\varepsilon^2) - [E(\varepsilon)]^2 = E(\varepsilon^2)$$

$$= [(1 - \pi(x))^2 \pi(x)] + [(-\pi(x))^2 (1 - \pi(x))]$$

$$= [(1 - 2\pi(x) + \pi(x)^2)\pi(x)] + [(\pi(x))^2 (1 - \pi(x))]$$

$$= [\pi(x) - \pi(x)^2]$$

$$= \pi(x)[1 - \pi(x)]$$

Böylece yanıt değişkeninin koşullu dağılımı,  $\pi(x) = E(y/x)$  koşullu dağılımına göre bir binom dağılımıdır.

Özet olarak, yanıt değişkeninin iki düzeyli olması durumunda regresyon analizinde:

- 1) Regresyon eşitliğindeki koşullu ortalama 0 ve 1 arasında bir değer olmalıdır.
- 2) Normal dağılımlar değil de binom dağılımları hatanın dağılımını tanımlar ve analiz bunun üzerine kuruludur.
- 3) Doğrusal regresyonda kullanılan ilkeler, lojistik regresyon analizinde de yol göstericidir [ 8].

## **1.2 LOJİSTİK REGRESYONUN LİNEER REGRESYONDAN FARKI ve KULLANIM AMAÇLARI**

Lojistik regresyonda, doğrusal regresyon analizinde olduğu gibi bazı değişken değerleri göz önüne alınarak tahmin yapılmaya çalışılır. Fakat bu iki yöntem arasında üç önemli fark vardır:

a) Doğrusal regresyon analizinde tahmin edilecek bağımlı değişken sürekli iken, Lojistik Regresyon analizinde bağımlı değişkenler kesikli bir değer almaktadır.

b) Doğrusal regresyon analizinde bağımlı değişkenin değeri tahmin edilirken, lojistik regresyon analizinde ise bağımlı değişkenin alabileceği değerlerden birinin gerçekleşme olasılığı tahmin edilir.

c) Doğrusal regresyon analizinde bağımsız değişkenin çoklu normal dağılım göstermesi koşulu aranırken lojistik regresyon analizinde böyle bir şart yoktur [3].

Öte yandan Lojistik Regresyon Analizi, diskriminant analizi ve çapraz tablo uygulamalarına alternatif olarak da uygulanmaktadır.

Lojistik regresyon analizinin güncel olmasının nedenleri şöyle özetlenebilmektedir:

a) Yanıt değişkeni kesikli iken açıklayıcı değişkenlerin hem sürekli hem de kesikli olma durumlarında uygulanabilmektedir.

b) Lojistik modelin parametreleri epidemioloji de yapılan ölçümlere benzediği için yorumları kolay olmaktadır. Epidemiolojide katsayıların exponansiyeli hastalık riski olarak yorumlanmaktadır.

c) Lojistik modelin parametre sayısı karşılık gelen doğrusal regresyon modeli ve diskriminant fonksiyonu ile aynı olmaktadır.

d) Lojistik modele dayalı analizler için standart paket programlar vardır

e) Açıklayıcı değişkenlerin olasılık fonksiyonlarının dağılımı üzerinde kısıt olmaması (yarı parametrik) nedeni ile çeşitli testler uygulanabilmektedir.

Epidemioloji ve diğer medikal uygulamaların yanı sıra deneysel verilerin analizinde, askeri konularda, meteorolojide, ziraatte, taşımacılıkta, ekonomi v.b. alanlarda sıkça kullanılan lojistik regresyon analizi farklı varsayımlar durumunda aynı lojistik formülasyona götürdüğü için varsayım bozulmalarına karşı daha güçlü bir yöntemdir [2,7].

Bu çalışmada, Lojistik Regresyon Modeli'nin kurulması, parametrelerin kestirim yöntemleriyle elde edilmesi ve kestirilen bu parametrelerin anlamlılığının incelenmesi açıklanacaktır. Kestirilen katsayılara ilişkin standart hatalar oldukça yüksek çıkıyorsa, önemli olduğu bilinen değişkenler önemsiz gözüküyorsa ve regresyon modelinin veriyi iyi karşılayamaması durumunda Doğrusal Olmayan Lojistik Regresyon Modeli'nin kullanılması önerilir. Bu amaçla Doğrusal Olmayan Lojistik Regresyon ve karşılaşılan model yapıları anlatılacaktır.

## BÖLÜM II

### LOJİSTİK REGRESYON MODELİNİN KURULMASI VE ANALİZİ

$(x_i, y_i)$ ,  $i=1, 2, \dots, n$ ,  $n$  tane bağımsız gözlem olsun.  $y_i$  iki düzeyli yanıt değişkeni ve  $x_i$  de  $i$ 'ninci denek için bağımsız değişken değerini temsil etsin. Yanıt değişkeninin 0 ve 1 kodlanmasının, belirli bir riskin yokluğunu ve varlığını temsil ettiği varsayalım.  $\pi(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$  lojistik regresyon denklemi için bilinmeyen  $\beta_0$  ve  $\beta_1$  parametrelerinin kestirilmesi gerekir. Bu bölümde bu parametrelerin kestirimi, modele katkıları, modelin anlamlılığı analiz incelenecektir.

#### II.1 PARAMETRE KESTİRİM YÖNTEMLERİ

##### II.1.1 En Çok Olabilirlik Yöntemi

Regresyon analizlerinde, hataların normal dağılması, ortalamasının sıfır olması ve bağımsız değişkenin her bir seviyesi için varyansın sabit kalması önemli varsayımlardır. Fakat iki yanıtli model olan lojistik regresyon analizinde bu varsayımlar geçerli değildir. Bu durumda  $Y = \pi(x) + \varepsilon$  modelinde yer alan hata terimi,  $Y=1$  ve  $Y=0$  durumunda iki olasılık değerinden birini alır:  $Y=1$  ise  $\pi(x)$  olasılığıyla  $\varepsilon = 1 - \pi(x)$ ;  $Y=0$  ise  $1 - \pi(x)$  olasılığıyla  $\varepsilon = -\pi(x)$ 'dir. Böylece  $\varepsilon$ 'nin ortalamasının sıfır ve varyansının  $\pi(x)[1 - \pi(x)]$  olduğu; bağımlı değişkenin( $y$ ) binomiyal dağılım gösterdiği görülür [3].

En çok olabilirlik yöntemi, olabilirlik fonksiyonu maksimum olacak biçimde bilinmeyen parametrelerin kestirimidir.

Y, 0 ve 1 değerlerini almış iken,

- $\pi(x)$ ,  $x$  verildiğinde Y'nin 1'e eşit olma koşullu olasılığıdır ve  $P(Y=1/x) = \pi(x)$  ile gösterilir.

- $1 - \pi(x)$ ,  $x$  verildiğinde Y'nin 0'a eşit olma koşullu olasılığıdır ve  $P(Y=0/x) = 1 - \pi(x)$  ile gösterilir.

$\pi(x_i)$ ,  $x_i$ 'deki  $\pi(x)$  değerini vermek üzere,

$y_i=1$  ise  $(x_i, y_i)$  gözlem çiftinin olabilirlik fonksiyonuna katkısı  $\pi(x_i)$ ,

$y_i=0$  ise  $(x_i, y_i)$  gözlem çiftinin olabilirlik fonksiyonuna katkısı  $1 - \pi(x_i)$

olasılığı kadardır.  $\beta = (\beta_0 \beta_1)'$  parametre vektörü iken,  $(x_i, y_i)$  ikililerinin olabilirlik fonksiyonuna katkısını göstermek için aşağıdaki yol izlenir:

$$\pi(x_i)^{y_i} [1 - \pi(x_i)]^{1-y_i}$$

Gözlemler birbirinden bağımsız olduğundan olabilirlik fonksiyonu,

$$l(\beta) = \prod_{i=1}^n \pi(x_i)^{y_i} [1 - \pi(x_i)]^{1-y_i} \quad (\text{II.1})$$

olarak bulunur. Log-olabilirlik (log likelihood) fonksiyonu ise,

$$L(\beta) = \ln[l(\beta)] = \sum_{i=1}^n \{y_i \ln[\pi(x_i)] + (1 - y_i) \ln[1 - \pi(x_i)]\} \quad (\text{II.2})$$

olur.  $L(\beta)$ 'yı maksimum yapan  $\beta$  değerini bulmak üzere,  $L(\beta)$ 'nin  $\beta_0$  ve  $\beta_1$ 'e göre türevleri alınırsa,

$$\sum_{i=1}^n [y_i - \pi(x_i)] = 0 \quad (\text{II.3})$$

$$\sum_{i=1}^n x_i [y_i - \pi(x_i)] = 0 \quad (\text{II.4})$$

"olabilirlik eşitlikleri" elde edilir.

Lojistik regresyonda bu eşitlikler  $\beta_0$  ve  $\beta_1$ 'e göre doğrusal olmayan, üstel denklem olduklarından bu denklemlerin çözümü için iteratif yöntemler kullanılır.

Lojistik modelin en çok olabilirlik kestirimlerini bulmak için iterasyona başlarken başlangıç değerlerini vermenin çeşitli yolları vardır. Bunlardan biri

grafiksel gösterimlerden göz ile kestirimde bulunmaktadır. Başlangıç değerlerinin doğruluğu iterasyon sayısı ve kestirimlerin etkinliği üzerinde önemli etkiye sahiptir. İyi bir başlangıç değeri ile az sayıda iterasyon sonucu optimal çözüme ulaşılabilmektedir. Ençok olabilirlik yönteminde her adımda yapılacak düzeltme miktarı( $\delta$ ), tek açıklayıcı değişken durumunda lojit tablolarından elde edilmektedir [2]. Veriler birbirinden çok ayrıık olduğunda en çok olabilirlik kestirimlerinde yakınsama elde edilememektedir [14].

(II.3) ve (II.4) eşitliklerinden elde edilen  $\beta$  değeri, en çok olabilirlik kestirimidir ve  $\hat{\beta}$  ile gösterilir.  $\hat{\pi}(x_i)$  de  $\pi(x_i)$ 'nin en çok olabilirlik kestirimidir ve verilen  $x = x_i$  değeri için Y'nin 1'e eşit olma koşullu olasılığının kestirimini verir.

(II.4) eşitliği  $\sum_{i=1}^n y_i = \sum_{i=1}^n \hat{\pi}(x_i)$  şeklinde yazılabilir. Bu eşitlik, y'nin gözlenen değerlerinin toplamının, y'nin beklenen değerleri toplamına eşit olduğunu gösterir. Bu özellik sonraki bölümlerde modelin uyumu tartışılırken yararlı olacaktır [8].

En çok olabilirlik ve Newton yöntemi yardımıyla parametre kestirimi aşağıdaki gibi bulunur [14]:

$$\pi(x_i) = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \text{ lojistik fonksiyon olsun.}$$

(II.2) nolu denklemin parametrelere göre daha açık biçimde yazılmasıyla,  $L(\beta)$ ,

$$L(\beta) = \sum y_i \cdot \beta_0 + \sum y_i \cdot \beta_1 x_i - \sum \ln(1 + e^{\beta_0 + \beta_1 x_i})$$

olur ve parametrelere göre türevleri alınıp sifıra eşitlenirse,

$$\frac{\partial L(\beta)}{\partial \beta_0} = \sum y_i - \sum \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} = \sum y_i - \sum \pi(x_i) = 0$$

$$\frac{\partial L(\beta)}{\partial \beta_1} = \sum y_i x_i - \sum x_i \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} = \sum y_i x_i - \sum x_i \pi(x_i) = 0$$

elde edilir. Bu iki denklemin Newton yöntemine göre iteratif olarak çözümünüyle kestiriciler bulunur:

$$f_1(\beta_0, \beta_1) = \sum y_i - \sum \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} = \sum y_i - \sum \pi(x_i)$$

$$f_2(\beta_0, \beta_1) = \sum y_i x_i - \sum x_i \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} = \sum y_i x_i - \sum x_i \pi(x_i)$$

$$J(\boldsymbol{\beta}) = \begin{bmatrix} \frac{\partial f_1}{\partial \beta_0} & \frac{\partial f_1}{\partial \beta_1} \\ \frac{\partial f_2}{\partial \beta_0} & \frac{\partial f_2}{\partial \beta_1} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial f_1}{\partial \beta_0} &= -\sum \left( \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right) \left( \frac{1}{1 + e^{\beta_0 + \beta_1 x_i}} \right) & \frac{\partial f_1}{\partial \beta_1} &= -\sum x_i \left( \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right) \left( \frac{1}{1 + e^{\beta_0 + \beta_1 x_i}} \right) \\ &= -\sum \pi(x_i)(1 - \pi(x_i)) & &= -\sum x_i \pi(x_i)(1 - \pi(x_i)) \end{aligned}$$

$$\begin{aligned} \frac{\partial f_2}{\partial \beta_0} &= -\sum x_i \left( \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right) \left( \frac{1}{1 + e^{\beta_0 + \beta_1 x_i}} \right) & \frac{\partial f_2}{\partial \beta_1} &= -\sum x_i^2 \left( \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right) \left( \frac{1}{1 + e^{\beta_0 + \beta_1 x_i}} \right) \\ &= -\sum x_i \pi(x_i)(1 - \pi(x_i)) & &= -\sum x_i^2 \pi(x_i)(1 - \pi(x_i)) \end{aligned}$$

$$J = - \begin{bmatrix} \sum \pi_i(1 - \pi_i) & \sum x_i \pi_i(1 - \pi_i) \\ \sum x_i \pi_i(1 - \pi_i) & \sum x_i^2 \pi_i(1 - \pi_i) \end{bmatrix}$$

Jakobiyen matrisi oluşturulduktan sonra iteratif çözüm vektörüyle iterasyon denklemleri çözülür:

$$\boldsymbol{\beta}^{(k+1)} = \boldsymbol{\beta}^{(k)} - J^{-1}(\boldsymbol{\beta}^{(k)}) \cdot F(\boldsymbol{\beta}^{(k)}) \quad (\text{II.5})$$

### II.1.2 Yeniden Ağırlıklandırılmış İteratif En Küçük Kareler Yöntemi

Lojistik regresyonda parametre tahmininde kullanılan bir başka yöntem de iteratif olarak ağırlıklandırılmış en küçük kareler yöntemidir. Bu yöntem, hata terimlerine ilişkin varyansların eşit olmadığı zamanlarda, değişen varyanslılık durumu söz konusu olduğunda doğru tahmin sağlar.

Gruplandırılmış verilerde J grubun her birinde  $n_j$  denemeden  $r_j$  başarı elde edildiğinde, başarı oranı  $P_j = \frac{r_j}{n_j}$  şeklinde gösterilsin. Varyansı

$$\sigma_j^2 = \text{Var}(r_j / P_j) = \text{Var}(P_j) = P_j \cdot \left( \frac{1 - P_j}{n_j} \right)$$

olduğundan her binom dağılımlı gözlem için varyans değişmektedir. Bu durumda lojit( $r_j/n_j$ )'nin açıklayıcı değişkenler üzerinde,

$$w_j = \frac{n_j}{P_j(1 - P_j)}$$

ağırlığı ile ağırlıklandırılmış regresyonu uygulanmalıdır. Fakat  $w_j$  ağırlık değerleri de  $P_j$ 'nin bir fonksiyonu olduğu için en küçük kareler yöntemi

iteratif olarak uygulanarak ağırlık değerleri her adımda yeniden elde edilmek suretiyle çözüme ulaşılır [4 , 23].

### II.1.3 Minimum Logit Ki-Kare Yöntemi

Ağırlıklı en küçük kareler kestirim yönteminin özel bir biçimidir. Berkson'un (1955) geliştirdiği bu yöntemde,  $2 \times j$  çapraz tablolarındaki beklenen ve gözlenen lojit değerler arasındaki farktan yararlanılmaktadır. Bu yöntem tekrarlı veriler olması durumunda kullanılmaktadır.

Veriler J grupta tekrar edildiğinde ve her grupta tekrar sayısı çok olduğunda katsayı kestirimleri ağırlıklı en küçük kareler yöntemleri elde edilebilmektedir.

II.1.2 kısmında değinilen yöntemde,  $P_j$  başarı olasılığı, lojistik fonksiyon eşitliğinde tanımlandığı gibidir. Bu olasılık üzerinde yapılan lojit dönüşüm, bu yöntemde yanıt değişkenini oluşturmaktadır. Kestirimde kullanılan ağırlık değerleri  $n_j P_j (1-P_j)$  olarak elde edilmektedir. Yöntem, lojit değeri olarak tanımlanan yanıt değişkeninin açıklayıcı değişkenler üzerindeki ağırlık değeri ile ağırlıklandırılmış regresyonundan EKK kestirimlerini elde etmeye dayanmaktadır. Buradan tek adımda bulunan ağırlıklı EKK kestirimleri minimum lojit Ki-Kare kestirimleri adını almaktadır.

Olasılık değerinin 0 veya 1 olduğu durumda lojit değeri tanımlı olmayacağı için  $P_j$  yerine  $P_j + 1/2n_j$  değerinin konulduğu ayarlanmış lojit ki-kare yöntemi kullanılmaktadır [2].

## II.2 PARAMETRELERİN ÖNEM TESTİ

Lojistik regresyonda katsayıların anlamlılık testleri bağımsız değişkeni içeren ve içermeyen modellerden bulunacak log-olabilirlik fonksiyonu ile yapılmaktadır.

Olabilirlik fonksiyonu kullanarak gözlenen değerlerle kestirilen değerlerin karşılaştırılması,

$$D = -2 \ln \left[ \frac{\text{Mevcut modelin olabilirliği}}{\text{Doymuş modelin olabilirliği}} \right] \quad (\text{II.6})$$

biçiminde verilir. Burada doymuş model veri noktaları kadar parametre içeren modeldir.

Eşitlik (II.6)'te parantez içerisinde verilen ifade olabirlik oranı olarak bilinir;  $(-2\ln)$  katının alınması, matematiksel olduğu kadar dağılımı bilinen bir değer elde etmektedir. Bu değer hipotez testi amacıyla kullanılmaktadır. Böyle bir test “Olabilirlik Oran Testi ” adını alır. Log olabilirlik eşitliği kullanılarak (II.6) eşitliği,

$$D = -2 \cdot \sum_{i=1}^n \left[ y_i \ln \left( \frac{\hat{\pi}_i}{y_i} \right) + (1 - y_i) \ln \left( \frac{1 - \hat{\pi}_i}{1 - y_i} \right) \right] \quad (\text{II.7})$$

biçiminde elde edilir. Burada  $\hat{\pi}_i = \hat{\pi}(x_i)$  'dir.

(II.7)'deki D istatistiği bazı yazarlar tarafından sapma (deviance) istatistiği olarak adlandırılır ve uyum iyiliğine karar verirken kullanılır [14]. Lojistik regresyon için sapma, doğrusal regresyondaki artık kareler toplamı ile aynı anlamı taşır ve (II.7) eşitliğindeki sapma, doğrusal regresyon için hesaplanırsa

$$\text{SSE} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \text{ 'ye eşit olur.}$$

Bağımsız bir değişkenin anlamlılığına karar vermek için, denklemde bağımsız değişkenin olduğu ve olmadığı durumlardaki D değerleri karşılaştırılır. Modeldeki bağımsız değişkeni kapsamasından dolayı ortaya çıkan D'deki değişim aşağıdaki gibidir:

$$G = D (\text{değişkensiz model için}) - D (\text{değişkenli model için}) \quad (\text{II.8})$$

Bu istatistik lojistik regresyonda, doğrusal regresyonda kullanılan F testinin pay kısmı olan regresyon kareler toplamı ile aynıdır. G'yi hesaplamak için farkı alınacak D değerlerinin her ikisi için de doymuş modelin olabilirlikleri ortak olduğu için G istatistiği şu şekilde ifade edilir:

$$G = -2 \cdot \ln \left[ \frac{\text{Değişkensiz modelin olabilirliği}}{\text{Değişkenli modelin olabilirliği}} \right] \quad (\text{II.9})$$

Tek bağımsız değişkenli ve değişkenin modelde olmadığı durumda  $\beta_0$ 'ın en çok olabilirlik kestirimi  $\ln(n_1/n_0)$ 'dır. Burada  $n_1 = \sum y_i$  ve  $n_0 = \sum (1 - y_i)$  iken, kestirim değeri de  $n_1/n$  sabit değeridir. Bu durumda G,

$$G = -2 \ln \left[ \frac{\left( \frac{n_1}{n} \right)^{n_1} \left( \frac{n_0}{n} \right)^{n_0}}{\sum_{i=1}^n \hat{\pi}_i^{y_i} (1 - \hat{\pi}_i)^{(1-y_i)}} \right] \quad (\text{II.10})$$

ya da

$$G = 2 \left\{ \sum_{i=1}^n [y_i \ln(\hat{\pi}_i) + (1 - y_i) \ln(1 - \hat{\pi}_i)] - [n_1 \ln(n_1) + n_0 \ln(n_0) - n \ln(n)] \right\} \quad (\text{II.11})$$

olur.  $\beta_1=0$  hipotezi altında, G istatistiği 1 serbestlik derecesiyle ki-kare ( $\chi^2$ ) dağılımına sahiptir.

Bağımsız değişkenin modele katkısının önemini test etmek için Wald ve Score testleri de mevcuttur. Bu testler için gerekli varsayımlar, olabilirlik oran testindeki varsayımlarla aynıdır.

Wald testi, eğim parametresinin en çok olabilirlik kestiriminin, onun standart hatasına oranı ile yapılır. Elde edilen oran,  $\beta_1=0$  hipotezi altında standart normal dağılımı gösterir:

$$W = \frac{\hat{\beta}_1}{\widehat{SE}(\hat{\beta}_1)} \sim N(0,1)$$

Bazı yazarlar Wald testi ile bağımsız değişkenin katsayısı gerçekte önemli olduğu halde test sonucunda katsayının önemsiz çıkabileceğini belirtmişlerdir [5]. Hem G olabilirlik oran testi, hem de W wald testinde  $\beta_1$  için en çok olabilirlik kestirimlerini hesaplamak gerekir.

Bir başka test de Score testidir. Hesaplama işlemlerini azaltması bu testin en büyük avantajı olmakla beraber, birçok paket programda bulunmaması ise kullanılmasını kısıtlamaktadır [6]. Score testi, log olabilirliğinden türetilen dağılım teorisine dayanmaktadır. Aslında matris hesaplamaları gerektiren çok değişkenli bir testtir.

Tek değişkenli durumda, bu test, (II.3) eşitliğindeki türevlerin verilmesiyle (II.4) eşitliğindeki türevlerin koşullu dağılımına dayanır. Burada  $\beta_0$  ve  $\beta_1$  katsayılarının hesaplanmasında (II.4) eşitliği kullanır:  $\beta_0 = \ln(n_1/n_0)$  ve  $\beta_1 = 0$  dır. Bu parametre değerlerine bağlı olarak  $\hat{\pi} = n_1/n = \bar{y}$  olur. Böylece (II.4) eşitliğinin sol tarafı  $\sum x_i (y_i - \bar{y})$ 'ye dönüşür. Kestirilmiş varyans  $\bar{y}(1 - \bar{y}) \sum (x_i - \bar{x})^2$  dir. Score testi (ST) için test istatistiği,

$$ST = \frac{\sum_{i=1}^n x_i (y_i - \bar{y})}{\sqrt{\bar{y}(1 - \bar{y}) \sum_{i=1}^n (x_i - \bar{x})^2}}$$

olarak verilir [8].

### II.3 ÇOKLU LOJİSTİK REGRESYON MODELİ

Bu bölümde, lojistik regresyon modelini birden çok bağımsız değişken olması durumu için kısaca ele alacağız.

$\mathbf{X}'=(x_1, x_2, \dots, x_k)$  vektörü ile gösterilen k tane bağımsız değişken kümesi olsun. Yanıt değişkeni için  $P(Y=1/x) = \pi(x)$  yazılsın. Çoklu lojistik regresyon modelinin lojiti aşağıdaki denklemle ifade edilir:

$$g(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k \quad (\text{II.12})$$

Buradaki lojistik regresyon modeli,

$$\pi(x) = \frac{e^{g(x)}}{1 + e^{g(x)}} \quad (\text{II.13})$$

şeklindedir.

Bağımsız değişkenlerden bazıları nominal ise bunları sürekli değişkenlermiş (aralık ölçekli) gibi kabul ederek modele dahil etmek uygun değildir. Değişkenin farklı düzeylerini göstermek için göstermelik (dummy, kukla) ya da tasarım (desing) değişkenleri kullanılır [ 8].

Genelde, nominal bir değişken d kategoriye sahipse, d-1 tane göstermelik ( $D_{ju}$ ,  $u=1, \dots, d-1$ ) tanımlanır. Böylece nominal ölçekli değişken içeren modelin lojiti aşağıdaki gibi olur:

$$g(x) = \beta_0 + \beta_1 x_1 + \dots + \sum_{i=1}^{d_j-1} \beta_{ju} D_{ju} + \beta_k x_k \quad (\text{II.14})$$

$\beta'=(\beta_0, \beta_1, \dots, \beta_k)$  iken çok değişkenli durumda, parametre kestirimleri için en çok olabilirlik yöntemi kullanılır. Olabilirlik fonksiyonu yaklaşık olarak (II.1) nolu denklemde verildiği gibi olmalıdır. Fakat tek değişiklik olarak  $\pi(x)$  denklem (II.13) de tanımlandığı gibidir. Log olabilirlik fonksiyonunun k+1 tane katsayıya göre türevi alınarak k+1 tane olabilirlik eşitlikleri,  $j=1, 2, \dots, k$  olmak üzere,

$$\sum_{i=1}^n [y_i - \pi(x_i)] = 0 \quad (\text{II.15})$$

ve

$$\sum_{i=1}^n x_i [y_i - \pi(x_i)] = 0 \quad (\text{II.16})$$

biçiminde elde edilir [11].

Kestirilen katsayıların varyans ve kovaryansları, log olabilirlik fonksiyonunun ikinci dereceden kısmî türevler matrisinden elde edilir. Bu kısmi türevler,  $j,u=0,1,2,\dots,k$ ,  $\pi_i = \pi(x_i)$  olmak üzere aşağıdaki biçimdedir:

$$\frac{\partial L(\beta)}{\partial \beta_j^2} = -\sum_{i=1}^n x_{ij}^2 \pi_i (1 - \pi_i) \quad (\text{II.17})$$

ve

$$\frac{\partial^2 L(\beta)}{\partial \beta_j \partial \beta_u} = -\sum_{i=1}^n x_{ij} x_{iu} \pi_i (1 - \pi_i) \quad (\text{II.18})$$

Eşitlik (II.17) ve (II.18)'de verilen terimlerin negatiflerini kapsayan  $(k+1) \times (k+1)$  boyutlu matris  $\mathbf{I}(\beta)$  ile gösterilir ve "bilgi (information) matrisi" olarak adlandırılır. Kestirilen katsayıların varyans ve kovaryansları bilgi matrisin tersinden elde edilir.  $\mathbf{Var}(\beta) = \mathbf{I}^{-1}(\beta)$ . Çok özel durumların dışında bu matrisin elemanlarını açık bir şekilde yazmak mümkün değildir. Matrisin  $j$ . köşegen elemanı  $\text{Var}(\beta_j)$ 'dir; köşegen dışındaki matris elemanları ise  $\text{Cov}(\beta_j, \beta_u)$  ile gösterilir. Varyans ve kovaryansların kestirimlerini  $\text{Var}(\beta)$ 'yi  $\hat{\beta}$  kullanarak buluruz. Kestirilmiş katsayıların kestirilmiş standart hataları, ( $j=0,1,2,\dots,k$ ) için,

$$\widehat{\text{SE}}(\hat{\beta}_j) = \left[ \widehat{\text{Var}}(\hat{\beta}_j) \right]^{1/2}$$

ile verilir [8].

Kestirilen modelin incelenmesinde ve uyumun değerlendirilmesinde faydalı olacak bilgi matrisi  $\hat{\mathbf{I}}(\hat{\beta}) = \mathbf{X}' \mathbf{V} \mathbf{X}$  şeklindedir. Burada  $\mathbf{X}$ ,  $n \times (k+1)$  boyutundaki bir matris olup her bir birey ve denek için verileri kapsar:

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix}$$

$\mathbf{V}$  ise  $n \times n$  boyutlu ve köşegen elemanları  $\hat{\pi}_i(1 - \hat{\pi}_i)$ , olan ( $i=1,2,\dots,n$ ) bir matristir [8, 11]:

$$\mathbf{V} = \begin{bmatrix} \hat{\pi}_1(1 - \hat{\pi}_1) & 0 & \dots & 0 \\ 0 & \hat{\pi}_2(1 - \hat{\pi}_2) & \dots & 0 \\ \vdots & 0 & \dots & \vdots \\ 0 & \dots & 0 & \hat{\pi}_n(1 - \hat{\pi}_n) \end{bmatrix}$$

Öte yandan anlamlılık testleri için, tek bağımsız değişkenli durumda verilen olabilirlik oran testinin çok değişkenli duruma genelleştirilmiş biçimi kullanılır. Olabilirlik oran testi (II.5) de verildiği gibi G istatistiğine bağlıdır. Aralarındaki tek fark  $\hat{\beta}$  vektöründe k+1 parametreyi kapsayan modelde  $\hat{\pi}$  değeridir. Modelde birlikte değişenler için k tane "eğim" katsayısının sıfıra eşit olması hipotezi (null hipotezi) altında, G istatistiği p serbestlik derecesiyle ( $\chi^2$ ) ki-kare dağılımı gösterir.  $H_0$  hipotezinin reddedilmesi doğrusal regresyondaki yoruma benzer yoruma sahiptir: En az bir veya tüm k katsayılarının sıfırdan farklı olması yorumudur.

Genel olarak, test edilecek hipotez,

$$H_0: \beta_q = \beta_{q+1} = \dots = \beta_{k-1} = 0$$

$H_1$ : En az bir  $\beta_p$  değeri sıfırdan farklıdır.

iken tüm parametreler üzerinde bulunacak olabilirlik fonksiyonu L(F); hipoteze göre kısaltılmış model için ise olabilirlik fonksiyonu L(R) olsun.

Test istatistiği,

$$\chi^2 = 2 \ln \left[ \frac{L(R)}{L(F)} \right] = -2 [\ln L(R) - \ln L(F)]$$

olur. Örnek hacmi yeterince büyük olduğunda  $H_0$  hipotezi doğru ise  $\chi^2$  istatistiği yaklaşık olarak  $\chi^2_{(1-\alpha; k-q)}$  şeklinde dağılım gösterir. Serbestlik derecesi,  $V = (n-q) - (n-k) = k-q$  şeklindedir. Böylece  $\chi^2 \leq \chi^2_{(1-\alpha; k-q)}$  olduğunda  $H_0$  kabul edilirken  $\chi^2 > \chi^2_{(1-\alpha; k-q)}$  olduğunda  $H_1$  kabul edilir [3].

Katsayıların hepsinin veya bir kısmının sıfırdan farklı olduğu sonucuna varmadan önce, değişkenler tek tek Wald test istatistiği ile test edilebilir. Wald testi için kullanılan eşitlik,  $W_j = \hat{\beta}_j / \widehat{SE}(\hat{\beta}_j)$ 'dir. "Bir katsayının ( $\hat{\beta}_j$ ) sıfıra eşit olması" hipotezi altında Wald istatistiği standart normal dağılım gösterir.

Öte yandan, Score testinin çok değişkenli durumlardaki ifadesi,  $L(\beta)$ 'nin  $\beta$ 'ya göre k tane türevinin koşullu dağılımı dikkate alınarak hesaplanır. Bu testin hesaplanması Wald testiyle aynı zorluklara sahiptir ve Score testinin de olabilirlik oran testine karşı bir üstünlüğü yoktur [8].

## II.4. LOJİSTİK REGRESYON MODELİNDE KATSAYILARIN YORUMLANMASI

Önceki bölümlerde model kurma ve modeldeki değişkenlerin önem testi üzerinde durulmuştu. Burada ise katsayıların yorumlanması incelenecektir. Lojistik regresyonda model uygunluğunun değerlendirilmesi ile ilgili yöntemler doğası gereği daha karışıktır.

### II.4.1 İki Sonuçlu (Dichotomous) Bağımsız Değişken

$x$  bağımsız değişkeni iki düzeyli olsun ve 0 , 1 biçiminde kodlandığını varsayalım.  $x=1$  ve  $x=0$  için lojitteki fark  $g(1)-g(0)=[\beta_0+\beta_1]-[\beta_0]=\beta_1$  şeklindedir. Bu durumda konuyu açıklamak için “Odds Ratio” larına gerek duyulur.

Lojistik regresyonda  $x$  bağımsız değişkeni 0 ve 1 ile kodlandığında mümkün olan olasılıklar  $\pi(x)$  ve  $1-\pi(x)$  değerleridir. Bu durum Tablo II.1’de gösterilmiştir.  $x=1$  iken riskin bireyler arasındaki olma olasılığı  $\pi(1)/[1-\pi(1)]$  olarak tanımlanır. Benzer şekilde  $x=0$  iken sonucun bireyler arasında olma olasılığı da  $\pi(0)/[1-\pi(0)]$  olarak tanımlanır. Odds Oranı (Odds Ratio) OR ile gösterilir ve aşağıdaki denklem ile verilir:

$$OR = \frac{\pi(1)/[1-\pi(1)]}{\pi(0)/[1-\pi(0)]} \quad (II.19)$$

**Tablo II.1** Bağımsız Değişkenin İkili Olması Durumunda Lojistik Regresyon Modelinin Değerleri

Sonuç Değişkenleri (Y)	Bağımsız Değişken (X)	
	$x=1$	$x=0$
$y=1$	$\pi(1) = \frac{e^{\beta_0+\beta_1}}{1+e^{\beta_0+\beta_1}}$	$\pi(0) = \frac{e^{\beta_0}}{1+e^{\beta_0}}$
$y=0$	$1-\pi(1) = \frac{1}{1+e^{\beta_0+\beta_1}}$	$1-\pi(0) = \frac{1}{1+e^{\beta_0}}$
Toplam	1.0	1.0

Tablo II.1'deki deęerler (II.19) denkleminde yerine konulursa odds oranı ařaęıdaki gibi bulunur:

$$\begin{aligned}
 OR &= \frac{\left(\frac{e^{\beta_0+\beta_1}}{1+e^{\beta_0+\beta_1}}\right) / \left(\frac{1}{1+e^{\beta_0+\beta_1}}\right)}{\left(\frac{e^{\beta_0}}{1+e^{\beta_0}}\right) / \left(\frac{1}{1+e^{\beta_0}}\right)} \\
 &= \frac{e^{\beta_0+\beta_1}}{e^{\beta_0}} \\
 &= e^{(\beta_0+\beta_1)-\beta_0} \\
 &= e^{\beta_1}
 \end{aligned} \tag{II.20}$$

Dolayısıyla iki düzeyli baęımsız deęiřkenin lojistik regresyonu için odds oranı  $OR = e^{\beta_1}$  dir. Odds deęerlerinin logaritması lojit olarak adlandırılır. Bu deęerlere iliřkin lojit deęerler,

$$g(1) = \ln \left\{ \pi(1) / [1 - \pi(1)] \right\}$$

$$g(0) = \ln \left\{ \pi(0) / [1 - \pi(0)] \right\}$$

olmak üzere log-odds (odds) logaritması,

$$\begin{aligned}
 \ln(OR) &= \ln \left[ \frac{\pi(1) / [1 - \pi(1)]}{\pi(0) / [1 - \pi(0)]} \right] = g(1) - g(0) \\
 &= \ln(e^{\beta_1}) \\
 &= \beta_1
 \end{aligned}$$

řeklinde ifade edilir; lojit fark (logit difference) olarak da adlandırılır. Yani  $x$  deęiřkeninin iki farklı deęeri için bulunan log-odds oranlarının farkı, bu iki kořulun karřılařtırılmasında kullanılacak log-odds oranını verir. Yukarıdaki örnekte log-odds oranı  $\ln(OR) = \ln(e^{\beta_1}) = \beta_1$  olup, regresyon modelindeki eęim katsayısını vermektedir. Bir bařka yorum ise,  $g(1)-g(0)$  olmasıyla  $\beta_1 = g(x+1)-g(x)$  řeklinde yazılarak baęımsız deęiřkendeki bir birimlik deęiřmeye karřılık lojitteki deęiřmeyi göstermesidir [11].

Odds Oranı geniř kullanıma sahip bir ölçüm birimidir. Epidemiolojik arařtırmalarda güçlü bir analitik metot olduęu kanıtlanan lojistik regresyonun sečilmesinin temel nedeni katsayılarının yorumlanabilir olmasıdır. Lojistik regresyonda katsayıların yerine odds oranlarının yorumlanması oldukça büyük önem

taşıır.  $x=1$  olan bireyler arasında sonuç deęişkeninin gerçekleşme olasılığı,  $x=0$  olan bireylere göre ne kadar çok veya az olduęu açısından özellikle epidemiolojide yaygın olarak kullanılan bir ilişki ölçüsüdür. Örneğın  $y$  yanıt deęişkeni bireyin akcięer kanseri olup olmadığını,  $x$  deęişkeni de bireyin sigara kullanım durumunu gösterir. Eđer  $\widehat{OR} = 2$  ise, akcięer kanserinin sigara içen kişilerde içmeyenlere göre iki kat daha sık görüldüğünü belirtir.

Odds oranı yardımı ile yorumlama birçok durumda görelı risk (relative risk) adı verilen bir nicelięe dayanır. Bu parametre  $\pi(1)/\pi(0)$  oranına eşittir. Eşitlik (II.18)den  $[1-\pi(0)]/[1-\pi(1)] \approx 1$  ise odds oranı görelı riske yaklaşır. Bu yaklaşım ancak  $\pi(x)$ ,  $x=0$  ve  $x=1$  için yeterince küçük ise geçerli olur [8].

Odds oranı, OR, lojistik regresyonda önemli bir parametredir. Teorik olarak örneklem genişlięi yeterince büyük olduęu zaman  $\widehat{OR}$  'nin dağılımı normaldir. Ancak örneklem büyüklüğü çoęu çalışmada sağlanamamaktadır. Bundan dolayı çıkarsamalar genellikle  $\ln(\widehat{OR}) = \hat{\beta}_1$  'in örneklem dağılımına dayanır. Odds oranının  $100(1-\alpha)\%$  güven aralıęı kestirimi, ilk olarak  $\beta_1$  katsayısı üzerinden elde edilir:

$$\exp\left[\hat{\beta}_1 \pm Z_{1-\alpha/2} \widehat{SE}(\hat{\beta}_1)\right]$$

Güven aralıęının sınırları 1'den büyük çıktıęı zaman dağılım sağa yatık olur ve aralıkta 1 deęerinin olmaması durumunda odds oranının istatistiksel olarak önemli olduęuna karar verilir. Bu durum sadece bağımsız deęişkenin 0 ve 1 deęerleriyle kodlanmasında geçerlidir [8, 11].

İlişki ölçüsü olarak odds oranının öneminden dolayı, nokta ve aralık kestirimleri genellikle lojistik regresyon analizi sonucunda tablolarda bulunur.

Bağımsız deęişkenin dięer kodlama türlerinde ise lojit farkın deęeri yeni kodlar üzerinden yeniden hesaplanır; sonra da  $e$ 'nin üssü şeklinde yazılır. İki farklı düzeydeki  $x=a$ 'ya karşı  $x=b$  için herhangi bir bağımsız deęişkenin odds oranının logaritması, bu iki deęer için hesaplanan lojitlerin farkına eşittir:

$$\begin{aligned} \ln\left[\widehat{OR}(a,b)\right] &= \hat{g}(x=a) - \hat{g}(x=b) & (II.21) \\ &= (\hat{\beta}_0 + \hat{\beta}_1 a) - (\hat{\beta}_0 + \hat{\beta}_1 b) \\ &= \hat{\beta}_1 (a-b) \end{aligned}$$

Kestirilen odds oranı, aşağıdaki gibidir:

$$\widehat{OR}(a,b) = \exp\left[\hat{\beta}_1(a-b)\right]. \quad (\text{II.22})$$

Bu ifade sadece  $(a-b)=1$  olduğu zaman  $\exp(\hat{\beta}_1)$  'e eşit olur. Eşitlik (II.21) ve (II.22)'deki  $\widehat{OR}(a,b)$  Odds notasyonunu gösterelim:

$$\widehat{OR}(a,b) = \frac{\hat{\pi}(x=a)/[1-\hat{\pi}(x=a)]}{\hat{\pi}(x=b)/[1-\hat{\pi}(x=b)]} \quad (\text{II.23})$$

a ve b, marjinal yöntem olarak da bilinen, eğer  $a=1$ ,  $b=-1$  (örneğin; kadın:-1, erkek:+1) biçiminde kodlanmış ise, (II.21) eşitliği,

$$\begin{aligned} \ln\left[\widehat{OR}(\text{kadın}, \text{erkek})\right] &= \hat{g}(D=1) - \hat{g}(D=-1) \\ &= [\beta_0 + \beta_1(D=1)] - [\beta_0 + \beta_1(D=-1)] \\ &= 2\hat{\beta}_1 \end{aligned}$$

ve kestirilen odds oranı  $\widehat{OR} = \exp(2\hat{\beta}_1)$  olur.

Kodlama metodu, aynı zamanda odds oranının güven aralığının alt ve üst sınırlarının hesaplanmasını da etkiler. Genel olarak denklem (II.23) de verilen odds oranı için güven aralığının alt ve üst sınırları,

$$\exp\left[\hat{\beta}_1(a-b) \pm Z_{1-\alpha/2} |a-b| \widehat{SE}(\hat{\beta}_1)\right] \quad (\text{II.24})$$

şeklindedir;  $|a-b|$ ,  $(a-b)$  değerinin mutlak değeridir. Genelde tüm iki düzeyli değişkenlerin 0 veya 1 diye kodlanması ve bu değişkenlerin aralık ölçekli değişken olarak varsayılması önerilir [ 8, 11].

#### II.4.2 Çok Sonuçlu (Polychotomous) Bağımsız Değişken

Nominal bir bağımsız değişkenin ikiden fazla düzey içerdiği durumda, bu değişkenin kategorilerini gösteren göstermelik/tasarım değişkenlerinin oluşturulması gerekir.

Göstermelik değişkenlerin seçimi, bir düzey referans seçilerek yapılabilir. Referans grup olarak seçilen düzey için 0 ve göstermelik değişkenler 1 değerini alır. İki düzeyin birbirine göre karşılaştırılması, örneğin,

$$\begin{aligned} \ln \left[ \widehat{OR}(\text{siyah}, \text{beyaz}) \right] &= \left[ \hat{\beta}_0 + \hat{\beta}_1(D_1 = 1) + \hat{\beta}_2(D_2 = 0) + \hat{\beta}_3(D_3 = 0) \right] \\ &\quad - \left[ \hat{\beta}_0 + \hat{\beta}_1(D_1 = 0) + \hat{\beta}_2(D_2 = 0) + \hat{\beta}_3(D_3 = 0) \right] \\ &= \hat{\beta}_1 \end{aligned}$$

biçiminde uygulanır [8].

Odds oranları için güven aralıkları, iki düzeyli bağımsız değişkenlerde kullanılan yaklaşımdan aynı şekilde elde edilir. İlk olarak Log-odds (lojistik regresyon katsayısı) için güven limitleri bulunur ve daha sonra bu limitleri ve üssü alınarak odds oranı için güven limitleri elde edilir. Genel olarak, katsayı için %100(1- $\alpha$ ) güven limitleri  $\hat{\beta}_j \pm Z_{1-\alpha/2} \widehat{SE}(\hat{\beta}_j)$  dir. Odds oranı için güven aralığı bu sınırların e-üssü alınarak,

$$\exp \left[ \hat{\beta}_j \pm Z_{1-\alpha/2} \widehat{SE}(\hat{\beta}_j) \right] \quad (\text{II.25})$$

elde edilir.

Tasarım değişkenlerinin kodlanmasının bir başka yolu ise ortalamadan sapma metodudur. Bu kodlama yöntemi genel ortalamadan grup ortalamasının sapmasının etkisini açıklar. Lojistik regresyonda, "grup ortalaması" grubun lojitidir ve "genel ortalama" ise tüm grupların ortalama lojitidir. Bu yöntemle göre, tasarım değişkeninin tüm değerleri -1 alınıp, geri kalan diğer değişkenler için 0 ve 1 kodlaması kullanılır.

Ortalamadan sapma kodlaması metoduyla kestirilen katsayıların yorumlanması, referans grup metodu kadar kolay ve açık değildir. Kestirilen katsayıların üssü, belirli bir grup için odds değerinin, oddsların geometrik ortalamasına oranını verir. Fakat buradan çıkan sonuç gerçeği yansıtmaz. Çünkü pay ve paydada iki farklı kategori için oddslar mevcuttur. Kestirilen katsayının üssü ortalama odds değerine göre ifade edilmektedir.

### II.4.3 Sürekli Bağımsız Değişken

Lojistik regresyon modeli bir sürekli bağımsız değişkeni içerdiği zaman, kestirilen katsayıların yorumlanması değişkenin modele nasıl girdiğine ve değişkenin birimine bağlıdır. Doğrusallık varsayımı altında, lojit,  $g(x) = \beta_0 + \beta_1 x$  'dir.  $\beta_1$  eğim katsayısı,  $x$ 'deki "1" birimlik artışın log-odds değerinde meydana getirdiği değişimi verir. Herhangi bir  $x$  değeri için  $\beta_1 = g(x+1) - g(x)$  dir. "1" değeri çoğu zaman biyolojik olarak pek ilginç değildir. Örneğin, yaş değişkenindeki "1" birim artış veya sistolik

kan basıncındaki 1 mm Hg artış önemli sayılmayacak kadar küçük bir değişimdir. Yaştaki 10 yıllık veya 10 mm Hg değişimi daha yararlı kabul edilebilir. Diğer yandan  $x$ 'in tanım aralığı 0'dan 1'e doğru ise, "1" birimlik değişim çok büyük olacaktır; 0,01 birimlik artış da daha gerçekçi olacaktır. Bu nedenden dolayı sürekli ölçekli kovaryantlar hakkında faydalı bir yorum sağlamak için kovaryanttaki " $c$ " değişimi için nokta ve aralık kestirimine bakmak gerekir.

$x$ 'deki  $c$  birimlik bir değişim için log-odds oranı  $g(x+c)-g(x)=c\beta_1$  lojit farkından elde edilir ve karşılık gelen odds oranı  $OR(c)=OR(x+c,x)=\exp(c\beta_1)$  olarak bulunur.  $\beta_1$ 'in en çok olabilirlik kestirimi  $\hat{\beta}_1$ , güven aralığı kestiriminde aşağıdaki gibi kullanılır:  $OR(c)$ 'nin %100  $(1-\alpha)$  güven aralığı kestirimi,

$$\exp\left[c\hat{\beta}_1 \pm Z_{1-\alpha/2}c\widehat{SE}(\hat{\beta}_1)\right]$$

şeklindedir.

Güven aralığının başlangıç ve bitiş noktalarının  $c$ 'nin seçimine bağlı olmasından dolayı  $c$ 'nin belirli bir değeri tüm tablo ve hesaplamalarda açıkça belirtilmelidir.  $c$ 'nin daha rastgele, keyfi seçimi bazı durumlarda zorluklar çıkarabilir.

Lojitin kovaryantta doğrusal olmadığına inanılıyorsa, göstermelik değişkenlerin kullanılması ve gruplandırılması düşünülmelidir. Alternatif olarak, yüksek dereceli terimler (örn:  $x^2$ ,  $x^3$ , ...) veya kovaryantta doğrusal olmayan ölçekleme (örn.  $\log(x)$ ) düşünülmelidir. Dolayısıyla, sürekli kovaryantları modellemenin önemli noktası lojitteki ölçeklendirmedir [8].

Özet olarak, sürekli bir değişkenin kestirilmiş katsayısının yorumlanması, nominal ölçekli değişkenlerinkine benzer. En önemli fark ise sürekli değişken için anlamlı bir değişimin tanımlanmasıdır.

#### II.4.4 Etkileşim Varlığında Odds Oranlarının Kestirimi

Zaman zaman iki ya da daha çok değişken söz konusu olduğunda etkileşim terimleri de modele girebilir. Risk faktörü ve başka bir değişken arasında etkileşim olduğunda risk faktörünün odds oranı kestirimi, onunla etkileşime giren değişkenin değerine bağlıdır. Bu durumda önceki bölümlerde verilen odds oranlarını kestirmek için kullanılan formülde bir değişiklik yaparak etkileşim içinde olan değişkenler arasındaki lojit farkı da hesaba katılır [8].

Bir modelde iki bağımsız değişken ve bunların etkileşimi yer alsın. F risk faktörü, x kovaryant ve onların etkileşimi de F.X şeklinde ifade edilsin. F=f ve X=x için bu modelin lojiti,

$$g(f.x) = \beta_0 + \beta_1 f + \beta_2 x + \beta_3 fx \quad (\text{II.26})$$

şeklinde olur. X=x sabit tutulduğunda F=f'e karşı F=f<sub>i</sub> düzeyleri için log-odds oranı,

$$g(f_1.x) = \beta_0 + \beta_1 f_1 + \beta_2 x + \beta_3 f_1 x$$

ve

$$g(f_0.x) = \beta_0 + \beta_1 f_0 + \beta_2 x + \beta_3 f_0 x$$

ifadelerinin farkının lojitinin alınmasıyla,

$$\begin{aligned} \ln[OR(F = f_1, F = f_0, X = x)] &= g(f_1, x) - g(f_0, x) \quad (\text{II.27}) \\ &= (\beta_0 + \beta_1 f_1 + \beta_2 x + \beta_3 f_1 x) - (\beta_0 + \beta_1 f_0 + \beta_2 x + \beta_3 f_0 x) \\ &= \beta_1 (f_1 - f_0) + x \beta_3 (f_1 - f_0) \end{aligned}$$

şeklinde elde edilir. Buradaki parametreler yerine kestirilmiş değerlerini kullanarak kestirilmiş log-odds değeri bulunmuş olur. Bu eşitliğin varyans kestirimi,

$$\begin{aligned} \widehat{Var} \left\{ \ln \left[ \widehat{OR}(F = f_1, F = f_0, X = x) \right] \right\} &= (f_1 - f_0)^2 \cdot \widehat{Var}(\hat{\beta}_1) + [x(f_1 - f_0)]^2 \cdot \widehat{Var}(\hat{\beta}_3) \\ &\quad + 2x(f_1 - f_0) \cdot \widehat{Cov}(\hat{\beta}_1, \hat{\beta}_3) \quad (\text{II.28}) \end{aligned}$$

biçiminde verilir [8].

Lojistik regresyon programlarının birçoğu modeldeki kestirilen parametrelerin varyans ve kovaryansının kestirimini verir. Kestirimler elde edildikten sonra, eşitlik (II.28) de bu değerler yerine konularak kestirilen log-odds oranının varyansının bir kestirimi bulunur. Log-odds oranı için %100(1-α) güven aralığı kestiriminin sınır (başlangıç ve bitiş) noktaları,

$$\exp \left( \left[ \hat{\beta}_1 (f_1 - f_0) \right] + x \hat{\beta}_3 (f_1 - f_0) \right) \pm Z_{1-\alpha/2} \widehat{SE} \left\{ \ln \left[ \widehat{OR}(F = f_1, F = f_0, X = x) \right] \right\} \quad (\text{II.29})$$

olarak elde edilir.

F iki düzeyli bir risk faktörü ise bu ifade, varyansı ve log-odds kestiricileri daha basit şekil alır. Eğer f<sub>i</sub>=1, f<sub>0</sub>=0 alırsak log-odds oranının kestirimi,

$$\ln \left[ \widehat{OR}(F = 1, F = 0, X = x) \right] = \hat{\beta}_1 + \hat{\beta}_3 x$$

olur. Varyans kestirimi

$$\widehat{Var}\left\{\ln\left[\widehat{OR}(F = f_1, F = f_0, X = x)\right]\right\} = \widehat{Var}(\hat{\beta}_1) + x^2\widehat{Var}(\hat{\beta}_3) + 2x\widehat{Cov}(\hat{\beta}_1, \hat{\beta}_3)$$

şeklindedir. Odds oranı için kestirilen güven aralığının sınır noktaları

$$\exp\left(\left(\hat{\beta}_1 + \hat{\beta}_3x\right) \pm Z_{1-\alpha/2}\widehat{SE}\left\{\ln\left[\widehat{OR}(F = 1, F = 0, X = x)\right]\right\}\right)$$

şeklinde elde edilir [8].

## II.5 MODEL UYUMUNUN BELİRLENMESİ

Modele gerekli tüm değişkenler alındıktan sonra kestirilen lojistik regresyon modelinin sonuç değişkenini tanımlamakta ne kadar etkili olduğu, uyum iyiliği (goodness-of-fit) testiyle anlaşılır.

Bağımlı değişkeninin gözlenen değerleri  $y$  vektörü şeklinde gösterilsin ve  $y'=(y_1, y_2, \dots, y_n)$  olsun. Model tarafından kestirilen değerler  $\hat{y}$  ile gösterilsin ve  $\hat{y}'=(\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n)$  olsun. Eğer  $y$  ve  $\hat{y}$  arasındaki uzaklık, yani artık(residual), özet ölçüleri küçükse ve her bir  $(y_i, \hat{y}_i)$ ,  $i=1,2,3,\dots,n$  ikilisinin bu özet ölçülere katkısı sistematik değil ve modelin hata yapısına göre küçükse modelin uyumlu olduğuna karar verilir. Böylece, uygun bir modele tamamen karar vermek için hem  $y$  ve  $\hat{y}$  arasındaki uzaklığın özet ölçüsünün hem de bu ölçülerin her bir parçasının tek tek incelenmesi gereklidir.

Uyum iyiliği için şu aşamalar uygulanır: (a) tümel uyum ölçülerinin hesaplanması ve değerlendirilmesi, (b) genellikle grafiksel olarak özet istatistik parçalarının ayrı ayrı incelenmesi ve (c)  $y$  ve  $\hat{y}$  parçaları arasındaki uzaklık veya farkın diğer ölçümlerinin incelenmesi.

Öte yandan, uyum iyiliği kovaryantın hepsiyle değil, modeldeki kovaryantlarla belirlenir. Örneğin  $x'=(x_1, x_2, \dots, x_k)$  bağımsız değişken vektörü ve  $J$  gözlenen  $x$ 'in farklı değerlerinin sayısını gösterebilir.  $m_j$  denek sayısını gösterebilir,  $j=1,2,\dots,J$ .

- Eğer analize dahil edilen denekler arasında aynı kovaryant değerine sahip denekler varsa, kovaryant sayısı ( $J$ ) < toplam denek sayısı ( $n$ ) olur, ( $m_j$  değeri büyük olma eğiliminde).

- Model geliştirme aşamasında, bağımsız değişkenlerin değerleri her bir denek için farklı ölçüm değerine sahipse  $J=n$  dir [8].

Özellikle modelde sürekli değişkenler yer alıyorsa  $J \approx n$  olması beklenilir. Uyum iyiliği testinde kullanılan farklı istatistikler vardır. Bunlar aşağıda açıklanmıştır ve bu istatistiklerde  $J=n$  olduğu varsayılmıştır.

### II.5.1 Ki-Kare İstatistiği ( $\chi^2_{B0}$ )

Modelde sadece sabit terim varken sözkonusu hatayı gösterir. Diğer bir anlatımla ( $\chi^2_{B0}$ ) istatistiği, modelde sadece sabit terim olduğunda  $-2\text{LogL}$  istatistiğini vermektedir. Yani ilk ki-kare istatistiği modeldeki tüm  $\beta$  katsayılarının sıfır olduğu hipotezini kabul eden  $-2\text{LogL}$  istatistiğidir.

### II.5.2 Pearson Ki-Kare İstatistiği ve Sapma Ölçütü

Lojistik regresyonda gözlenen ve kestirilen değerler arasındaki farkın (doğrusal regresyondaki  $y - \hat{y}$  artık fonksiyonu gibi) birçok ölçüsü vardır. Lojistik regresyonda kestirilen değerler her bir kovaryantın deseni için hesaplanır ve o kovaryant deseni için kestirilen olasılığa bağlıdır; kestirilen değer  $\hat{y}_j$  ile gösterilir:

$$\hat{y}_j = m_j \hat{\pi}_j = m_j \frac{e^{\hat{g}(x_j)}}{1 + e^{\hat{g}(x_j)}}$$

Buradaki  $\hat{g}(x_j)$  kestirilen lojittir.  $m_j$  kovaryant değerleri birbirinden farklı olan denek sayısıdır.

Gözlenen ve kestirilen değerler arasındaki farkın iki ayrı ölçüsünü inceleyelim. Bunlar Pearson artığı ve deviance artığıdır. Kovaryant deseni için Pearson artığı

$$r(y_j, \hat{\pi}_j) = \frac{(y_j - m_j \hat{\pi}_j)}{\sqrt{m_j \hat{\pi}_j (1 - \hat{\pi}_j)}} \quad (\text{II.30})$$

dir. Bu artıklar üzerinden hesaplanan istatistik Pearson ki-kare istatistiği,

$$\chi^2 = \sum_{j=1}^J r(y_j, \hat{\pi}_j) \quad (\text{II.31})$$

ile verilir. Sapma artığı aşağıdaki gibi tanımlanır:

$$d(y_j, \hat{\pi}_j) = \pm \left\{ 2 \left[ y_j \cdot \ln \left( \frac{y_j}{m_j \hat{\pi}_j} \right) + (m_j - y_j) \cdot \ln \left( \frac{(m_j - y_j)}{m_j (1 - \hat{\pi}_j)} \right) \right] \right\}^{1/2} \quad (\text{II.32})$$

Burada denklemin önündeki ( $\pm$ ) işareti ( $y_j - m_j \hat{\pi}_j$ ) ifadesinin işareti ile aynıdır.  $m_j=1$  ve  $y_j=0$  kovaryant deseni için deviance artışı aşağıdaki şekle dönüşür:

$$d(y_j, \hat{\pi}_j) = -\sqrt{2m_j |\ln(1 - \hat{\pi}_j)|}$$

Başka bir örnek olarak  $m_j=1$  ve  $y_j=1$  için sapma artışı,

$$d(y_j, \hat{\pi}_j) = \sqrt{2 |\ln(\hat{\pi}_j)|}$$

olur. Sapma artıklarına dayanan özet istatistiği,

$$D = \sum_{j=1}^J d(y_j, \pi_j)^2 \quad (\text{II.33})$$

olarak yazılır [8].

Kurulan modelin tüm yönlerden doğru olduğu varsayımı altında  $\chi^2$  ve D istatistiklerinin dağılımının J-(k+1) serbestlik derecesiyle ki-kare dağılımına sahip olduğu varsayılır. Aynı zamanda D istatistiği (k+1) parametrelili kestirim modeline karşılık J parametrelili doymuş modelin olabilirlik oran testidir. Benzer durum  $H_0$  hipotezi geçerli iken  $\chi^2$  istatistiğinin gösterdiği dağılım için de geçerlidir:

$$D = -2 \ln \left[ \frac{\text{Uyarlanmış modelin olabilirliği}}{\text{Doymuş modelin olabilirliği}} \right]$$

Diğer bir deyişle olabilirlik oran testi,

$$D = -2 \sum_{i=1}^n \left[ y_i \ln \left( \frac{\hat{\pi}_i}{y_i} \right) + (1 - y_i) \ln \left( \frac{1 - \hat{\pi}_i}{1 - y_i} \right) \right]$$

ile verilir.  $j=n$  alınmasıyla bu istatistik D, sapma olarak da adlandırılır [12, 14].

Sapma, genelde analize bağımsız değişken ilave edildiğinde modelin hatasını gösterir. Bu nedenle -2LogL istatistiği bağımlı değişkendeki açıklanmayan varyansın anlamlılığını gösterir. Bu istatistiğin anlamlı olmaması lojistik regresyon analizinde istenen durumu göstermektedir. Log olabilirlik değeri 0 ve 1 arasında değerler almaktadır. Bu oran bağımlı değişkenin bağımsız değişkenler tarafından tahmin edilme olasılığını gösterir. -2LogL istatistiği yaklaşık olarak ki-kare dağılımına uyduğundan lojistik regresyon analizindeki bu istatistik regresyon analizindeki artık kareler toplamına benzer. Olabilirlik oranı 1 ise -2LogL sifıra eşit olur. Model, doymuş modelle uyumlu tam olarak temsil edilmesi durumunda bu sonuç gerçekleşir. Özetle -2LogL istatistiği ne kadar küçük değere sahipse daha iyi bir model oluşur [8, 14].

### II.5.3 Hosmer-Lemeshow (G) İstatistiği

Hosmer ve Lemeshow, kestirilen olasılık değerlerinin gruplandırılmasını önermişlerdir.  $J=n$  ise veri matrisindeki  $n$  tane sütun.  $n$  tane kestirilen olasılık değerine karşılık gelir. Burada gruplama yapmadaki amaç daha düzgün bir beklenen değer oluşturulmasıyla mevcut dağılımı ki-kare dağılımına yaklaştırarak anlamlı, yorumlanabilir model elde etmektir. Gruplama iki farklı tipte yapılabilir:

- (a) kestirilen olasılıkların yüzdesi dikkate alınarak
- (b) kestirilen olasılıkların sabit değerleri dikkate alınarak.

İlk yöntemde gözlenen ve kestirilmiş beklenen frekansları karşılaştıran  $\hat{C}$  istatistiğini inceleyeceğiz. Burada 10'lu risk grubu kullanılır. Yani tüm gözlemler 10 gruba ayrılır. Bu yöntem yeterli sayıda gözlem olduğunda geçerlidir. Olumsuz yönü ise gerçek değerlerin göz ardı edilmesidir.  $\hat{C} \sim \chi^2_{(g-2)}$ ,  $g$  grup sayısıdır.

İkinci yöntemde sabit kesim noktaları üzerine bir gruplandırma yapılır. Örneğin kestirilen olasılıklara göre sabit gruplar oluşturularak deneklerin ilgili grupta yer alması sağlanır.

Gruplama yoluyla verileri azaltma işleminde, gruplardaki veri sayısı azaldığından dolayı uyumdan sapmalar görülebilir. Hosmer-Lemeshow testi yorumlama ve geniş veri kümesini rahat çözmek amaçlı en yaygın tercih edilen bir uyum testidir. Bu istatistik SPSS'de "Hosmer-Lemeshow G" olarak bilinmesinin yanı sıra Model ki-kare istatistiği olarak adlandırılır. Bağımsız değişkenlerden hiçbirinin bağımlı üstünlük oranıyla anlamlı doğrusal bir ilişki göstermediğini ileri süren sıfır hipotezini test etmektedir. Diğer bir anlatımla bu istatistik, sabit terimin dışındaki tüm lojit katsayılarının sıfıra eşit olup olmadığını sınamaktadır.

Model ki-kare istatistiği bir olabilirlik oran testidir ve bu yüzden modelde bağımsız değişkenin olmadığı  $-2\text{LogL}$  istatistiği ile modelde bağımsız değişkenlerin olduğu  $-2\text{LogL}$  istatistiği arasındaki fark olarak hesaplanmaktadır. İncelenen modelin parametre sayısı ile yalnız sabit terimli modelin parametreleri arasındaki farka eşit bir serbestlik derecesi ile ki-kare dağılımına uymaktadır. Lojistik regresyon analizinde model ki-kare değerinin anlamlı olması arzu edilen durumu göstermektedir. Model ki-kare testi regresyon analizindeki F testine benzer [8, 10].

#### II.5.4 Blok Ki-Kare İstatistiği

Bir blok değişkeninin modele dahil edilmesiyle model ki-kare istatistiğine meydana gelen değişmeyi gösterir. Bu istatistik adimsal lojistik regresyon analizinde "Step ki-kare" adıyla hesaplanır. Her adımda modele tek bir değişken ilave edilip çıkartılıyorsa, doğal olarak blok ve step ki-kare istatistikleri eşit olmaktadır. Kategorik bağımsız değişken modele dahil edildiğinde blok ki-kare istatistiği ile test edilmektedir. Bu durumda kategorik değişken ile ilgili tüm göstermelik değişkenler blok olarak modele dahil edilmektedir [8, 10].

#### II.5.5 Sınıflandırma Tablosu

Kurulan bir lojistik modelin sonuçlarını özetlemek için sezgiye dayanan en iyi yol sınıflandırma tablosu oluşturmaktır. Bu tablo, sonuç değeri y'nin düzeyleri ile kestirilen lojistik olasılıklar tarafından türetilen ikili bir değişkenin çapraz sınıflandırılmasıyla elde edilir.

Türetilen bu ikili bağımsız değişkeni elde etmek için c kesim noktası belirlenir ve kestirilen her bir olasılık değeri c ile karşılaştırılır. Eğer kestirilen olasılık c değerini geçerse, türetilen ikili değişken 1'e eşit olur, diğer durumlarda 0'a eşittir. c'nin en yaygın kullanılan değeri 0.5'dir.

Bu yaklaşımda kestirilen olasılıklar grup üyeliğini tahmin etmek için kullanılır:

$$Y=0 \text{ grubunda } P(Y=1)=Q_1 \text{ ve } X \sim N(0,1)$$

Y=1 grubunda  $X \sim N(\mu,1)$  olduğu varsayılırsa lojistik regresyon için eğim katsayısı  $\beta_1 = \mu$  olur ve sabit değer,

$$\beta_0 = \ln \left[ \frac{Q_1}{(1-Q_1)} \right] - \frac{\mu^2}{2}$$

ile verilir. Lojistik regresyon modelinin Hatalı Sınıflandırma Olasılığı (PMC),

$$PMC = Q_1 \Phi \left\{ \frac{1}{\beta_1} \ln \left[ \frac{(1-Q_1)}{Q_1} \right] - \frac{\beta_1}{2} \right\} + (1-Q_1) \Phi \left\{ \frac{1}{\beta_1} \ln \left[ \frac{Q_1}{(1-Q_1)} \right] - \frac{\beta_1}{2} \right\}$$

dir.  $\Phi$ ,  $N(0,1)$  dağılımının kümülatif dağılım fonksiyonudur. Böylece beklenen hata oranı, eğimin büyüklüğünün bir fonksiyonu olduğundan, model uyumuyla ilişkisiz olduğundan doğru veya yanlış sınıflandırma uyum iyiliği için bir kriter olarak dikkate alınmaz. Sınıflandırma tablosu bağımlı değişken gruplarındaki denek,

gözlem sayısına duyarlıdır. Grup büyüklüğü arttıkça doğru sınıflandırma olasılığının güvenilirliği artar [8].

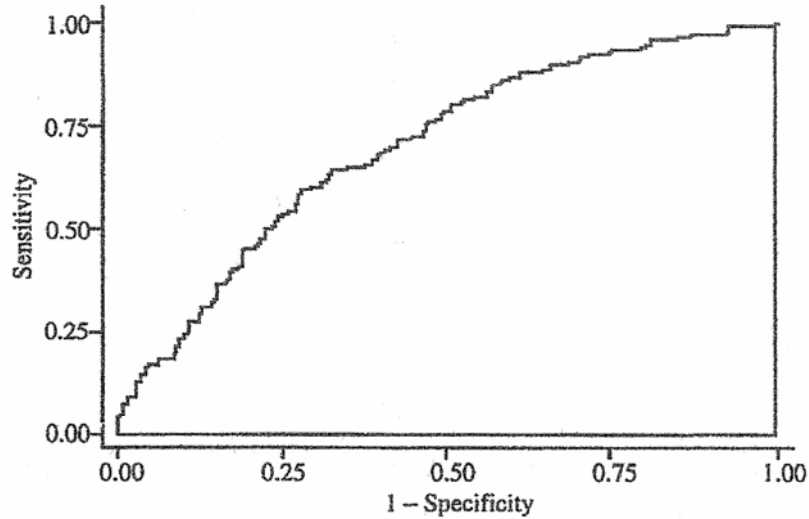
### II.5.6 ROC Eğrisi Altındaki Alan

Doğruluk yüzdeleri, bir test sonucunun doğru sınıflandırılmasındaki tek kesim noktasına dayanır. ROC (Receiver Operating Characteristic) eğrisi altındaki bölge, sınıflandırmada doğruluk tanımından daha fazlasını verir. ROC eğrisi altında kalan alan 0'dan 1'e kadar değişir. Gerçekleşmeyen duruma karşı gerçekleşen değerlerin farkını görmek bu ayrımı yapabilmek için modelin bir yeterlilik ölçüsü olarak kullanılır.

Sınıflandırma tablosunda duyarlılık ve özgüllük  $\left( \begin{matrix} \text{sensitivity, specificity} \\ (\hat{y}=1 \wedge y=0) & (\hat{y}=0 \wedge y=0) \end{matrix} \right)$

sadece ek bilgi verir, model uyumuna karar vermede doğrudan etkili değildir.

Hedefimiz, sınıflandırmanın amacı için optimal bir kesim noktasını seçmek ise (sensitivity ve specificity) doğruluk yüzdelerinin her ikisini de maksimum yapan bir kesim noktası seçmek olacaktır.



Şekil II.1 ROC Eğrisi Altında Kalan Alanın Grafikselleştirilmesi

Bu Şekil II.1'de ROC eğrisi altındaki alan 0.6989'dur. İstatistiksel paket programlarla ilgili değerler (kesim noktası gibi) girildikten sonra ROC eğri altındaki alana ulaşılır. Genel bir kural olarak: (ROC eğrisi altında kalan alanın doğru sınıflandırma ölçüsüne göre bir olasılık olduğu ve (0,1) arası değer aldığı unutulmamalıdır).

ROC = 0.5 alana sahipse ayırım yapılamaz, yorumlanamaz.

$0.7 \leq \text{ROC} < 0.8$  ise kabul edilebilir ayırım olduğu düşünülür.

$0.8 \leq \text{ROC} < 0.9$  ise mükemmel ayırımdır.

$\text{ROC} \geq 0.9$  ise çok iyi ayırımdır.\*

\* Pratikte ROC eğrisi altındaki alanın 0.9'dan büyük gözlenmesi olağan dışıdır. Gerçekte önceki konularda belirtildiği gibi tam bir ayırım olduğunda lojistik regresyon modelinin katsayılarını kestirmek imkansızdır, fakat tam ayırma, sınıflandırmaya yakın bir ihtimal ROC eğrisi altında kalan alan  $> \%90$  olduğunda kabul edilebilir [ 8].

## II.6 MODEL DEĞERLENDİRME

İstatistikte geliştirilen bir modelin geçerliliğinin değerlendirilmesi büyük önem taşır. Lojistik regresyon modelinin uygunluğunun değerlendirilmesinde genelde gerçek olasılıklar ile kestirilen olasılıklar arasındaki standart farka bakılır [10]. Gözlemin, yeni modelin hangi açıdan etkili olduğu, değişik hataların dağılımları, ilişki ölçümü gibi göstergeler incelenir. Lojistik regresyonda hesaplanabilen hata türleri:

### a) Standart Olmayan Hatalar

Standart olmayan hatalar ( $e_i$ ), fiili (gerçek) olasılıklar ile kestirilen olasılıklar arasındaki farka eşittir. Modelin hataları lojit ölçekte hesaplanmış ise bu hatalara lojit hatalar adı verilir. Lojit hatalar aşağıdaki formülle hesaplanır:

$$\text{Lojit Hata} = \frac{e_i}{P_i(1 - P_i)} \quad (\text{i. birimin standart olmayan hatası})$$

### b) Standart Hatalar

Standart hatalar, standart olmayan hataların ( $e_i$ ) kendi standart sapmalarına bölünmesiyle elde edilir. Standart hatalar, aşağıdaki formülle hesaplanır:

$$Z_i = \frac{e_i}{\sqrt{P_i(1 - P_i)}} \quad (\text{i. birimin standart hatası})$$

Her bir birimin standart hatası  $\chi^2$  uygunluk istatistiğinin bir bileşeni olarak görülebilir. Büyük örnekler için standart hatalar 0 ortalama ve 1 standart sapma ile normal dağılıma uyar.

### c) Sapma Değerleri

Modelin serbestlik derecesiyle de ilgili olan sapma değeri deviance artığından elde edilir.  $y=0$  riskli karşılaşılmayan durum için sapma değeri,

$$\text{Sapma} = \sqrt{-2 \cdot \ln(P_i)}$$

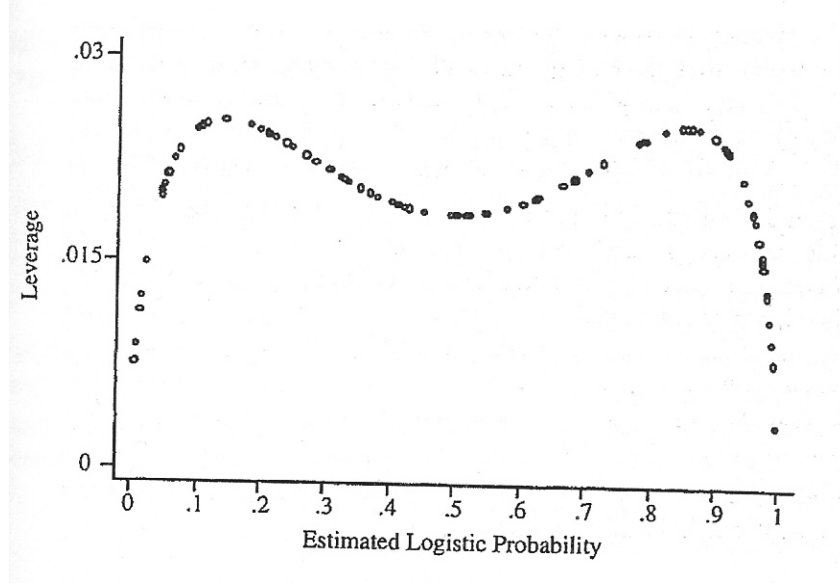
$y=1$  riskiyle karşılaşıldığında,

$$\text{Sapma} = \sqrt{-2 \cdot \ln(1 - P_i)}$$

formülünden hesaplanır. Büyük sapma değerleri modelin ilgili verileri iyi temsil etmediğini gösterir. Büyük örnekler için sapma değerleri yaklaşık olarak normal dağılıma uymaktadır.

### d) Uzaklık (Leverage) Değerleri

Uzaklık değerleri kestirilen değerler üzerinde büyük etkisi olan birimlerin belirlenmesi amacıyla kullanılır. Uzaklık değerleri 0 (tamamen etkisiz) ve 1 (tamamen etkili) aralığı içinde değerler alır. Uzaklık değerlerinin ortalaması  $P/n$  oranına eşittir.  $P$ , sabit terim dahil, modelde tahmin edilen parametre sayısını ve  $n$  örnek hacmini göstermektedir. Böylece uzaklık değerleri ortalama uzaklık değeri ile karşılaştırılmaktadır. Ayrıca herhangi birbirimin kestirilen olasılığı %10'dan küçük veya %90'dan büyük olması durumunda, söz konusu birim etkili bir birim olsa bile uzaklık değeri küçük hesaplanabilmektedir [10].



**Şekil II.2** Leverage Değerine (h) Karşı Kestirilen Olasılık Değerinin ( $\hat{\pi}$ ) Çizimi

Eğer kestirilen olasılık 0.1 ve 0.9 arasında ise, leverage uzaklık olarak düşünülebilecek bir değer verir. Kestirilen olasılık 0.1 ve 0.9 aralığı dışında olduğu zaman, leverage değeri uzaklığı tam anlamıyla ölçemeyebilir.

#### e) Cook Uzaklığı (Cook's Distance)

Cook uzaklığı değeri herhangi bir birimin model üzerindeki etkisini gösterir. Cook uzaklığı belirli bir birimin modelden çıkartılması durumunda lojistik regresyon katsayılarının ne kadar değişeceğini gösterir. Aşağıdaki formülle hesaplanır:

$$CU_i = Z_i^2 \left( \frac{h_i}{1-h_i} \right)$$

Formülde  $Z_i$  standartlaştırılmış hataları ve  $h_i$  ise uzaklık (leverage) değerini gösterir. Formülden kolayca görülebileceği gibi Cook uzaklığı hem standart hataya hem de uzaklık değerine bağlıdır.

#### f) Grafik Yöntemler

Yukarıda bahsedilen istatistiklerden uygun olanlar kullanılarak normal olasılık ve diğer grafikler elde edilebilmektedir. Örneğin Şekil II.7 grafik yöntem olarak incelenip, yorumlanabilir. Uzaklık değerleri ortalama uzaklık değeri ile karşılaştırıldığından, kestirilen olasılık 0.1 ve 0.9 arasında olduğunda ortalama

uzaklık deęerinden uzaklařan gözlemin etkililięine bakılır. İlgili gözlem noktasının artık deęeri çok fazla büyük çıkıyorsa etkili bir gözlemdir ve bu durumda, bu gözlem model yapısını etkileyeceęinden gözlem noktasının veriden çıkarılamayacaęına karar verilir.

Sonuç olarak bu modellerin deęerlendirilmesi, yorumlanmasıyla modellerin uyum iyilięi arařtırılmıř olur [10].

## BÖLÜM III

### DOĞRUSAL OLMAYAN LOJİSTİK REGRESYON MODELİ

Bağımlı değişkenin değerinde gözlemlenen değişim bağımsız değişkenin aldığı değere bağlı olarak sabit çıkmıyorsa, bağımlı değişken ile bağımsız değişken arasında doğrusal olmayan bir ilişki vardır denir. Burada doğrusal olmayan regresyon modelleri olarak “segmented” veya diğer bir deyişle “splayn” fonksiyonlar incelenecektir.

Bağımlı değişken ile bağımsız değişkenler arasındaki doğrusal olmayan ilişkiyi bulmak için lojistik regresyon modelinin doğrusal, karesel, kübik veya daha yüksek dereceli bir model olup olmadığı sınanmalıdır. Eğer türetilen bağımsız değişken sayısı çok fazla ise, kestirilen katsayılara ilişkin standart hatalar oldukça yüksek çıkma eğilimi gösterir ve modele dahil edilen doğrusal terimler istatistiksel anlamda önemli olduğu halde önemsiz çıkar. Bu durumda lojistik regresyon modelinin doğrusallık varsayımının bozulduğu görülür. Böylece doğrusal olmayan lojistik regresyon analizi kullanılarak verilerin modellenmesi tercih edilir.

Doğrusal olmayan lojistik regresyon modeli bilinmeyen parametre vektörünün bileşenlerinin en az birine göre doğrusal olmayan bir fonksiyona sahiptir. Doğrusal olmayan lojistik regresyon analizinin teorisi ve yöntemleri, doğrusal lojistik regresyon modelinin teorisi ve yöntemleriyle benzerdir.

Doğrusal olmayan regresyon modelleri oldukça yararlı olmasına karşın, verilerin normal olmaması halinde, kullanılan basit model uygun değildir. Bazen dönüşüm yapılarak bu sorun çözülebilir. Bununla birlikte, uygulamada veriler göz önüne alınarak en uygun fonksiyonel şekli tahmin etmek zordur. Parametrik modelin uygun olduğu düşünülebilen durumlarda bile, bağımsız değişkenlere göre bağımlı değişkenin fonksiyonel bağımlılığını ortaya çıkarmada, doğrusal olmayan regresyon yararlı bir tekniktir [22].

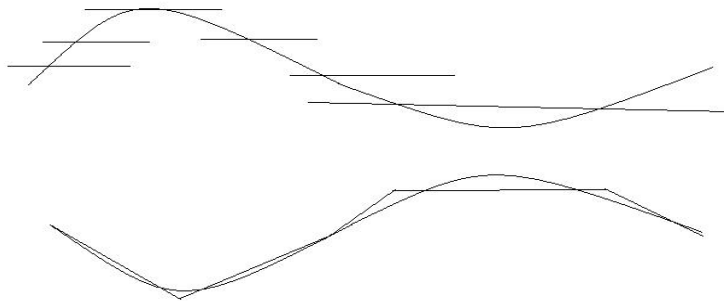
### III.1 SPLAYNLAR VE PARÇALI POLİNOM UYDURMA

“Splayn”, bir dizi veri noktalarına polinomial bir eğri uydurma veya bu noktalar arasından pürüzsüz olarak geçen ve birçok parçadan oluşan esnek bir eğridir. Splayn fonksiyonlar yeni yeni geliştirilen matematiksel araçlardır. Bu fonksiyonların temel düşüncesi, tanımlanan aralığı bağımsız değişkenlerin gözlem değerleri yardımıyla alt aralıklara bölerek, her bir alt aralıkta farklı bir polinomial fonksiyon ile bağımlı ve bağımsız değişkenler arasındaki ilişkiyi modellendirerek istenilen mertebeden türevi olan sürekli fonksiyon elde etmektir [ 1 ].

Bazen düşük dereceli polinomların veriye uydurulmasında zayıf kaldığı, polinom derecesinin yükseltilmesinin de durumu düzeltmediği görülür. Bunun nedeni artık kareler toplamındaki başarısızlık veya açıklanmamış yapıları gösteren artık çizimleridir. Bu problem, fonksiyon,  $x$  dizisinin farklı parçalarında farklı davrandığında ortaya çıkar. Bazen  $x$  ve/veya  $y$ 'deki dönüşümler bu problemi ortadan kaldırır. İdeal yaklaşım, verilerdeki  $x$  dizisini bölmek ve her bir bölüme uygun bir eğri oluşturmaktır. Splayn fonksiyonlar parçalı polinom uydurmayı gerçekleştirmek için kullanışlı bir yoldur [16].

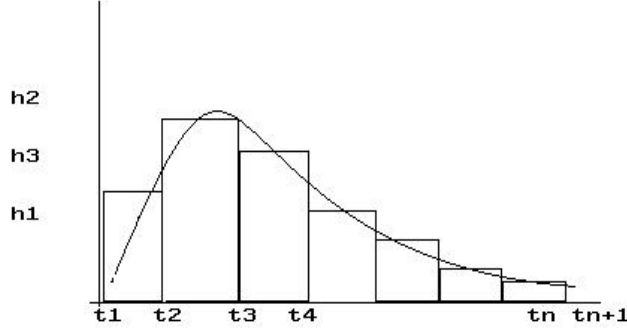
Parçalı polinomların sıradan polinomlara göre bükülgenliği zayıftır, iyi düzgünleştirilmeleri zordur. Splaynlar  $n-1$  sürekli türevlerle kırılma noktalarıyla  $n$ .dereceden parçalı polinomlardan oluşur. Splaynların kırılma noktaları, polinom parçalarının birbirine bağlandığı ortak noktalar olan “düğüm(knots)” olarak adlandırılır.  $n=2$  veya 3 veya daha fazla ise splayn düzgün, esnek şekildedir.

Splaynların iki klasik örneği vardır: 0.dereceden olan adım fonksiyonları ve 1.dereceden splaynlar olan kırık çizgiler. Bu örnekler düzgünleştirilmenin olmadığı durumlar şeklinde aşağıda gösterilmiştir.



**Şekil III.1** Adım Fonksiyonu ve Kırık Çizgi. 0. ve 1.Dereceden Splayn Yaklaşımıyla İki Düzgün Eğri

Bir histograma ait parabolik splayn yaklaşımı Şekil III.2' de oldukça düzgünleştirilmiş bir biçimde gösterilmiştir.



**Şekil III.2** Bir Parabolik Splayn. 0.Dereceden Bir Splayn Histogramının Her Bir Parçasındaki Alana Eşit Olan Parabolik Splayn

Histogram  $h_i$  yüksekliğine,  $t_i$  düğüm noktasına sahiptir ve parabolik splayn her parçasının altındaki bölge histogramdaki basamaklara karşılık gelir. Parabolik splaynın ayrıntılı incelenmesi eğrilerin 2.türevlerindeki sıçramaları gösterir.

Genelde parçalı polinomlardan çok splaynların tercihi daha fazladır; çünkü düzgünleştirme durumları bazı katsayıları yok sayar. Bu, düğümlerdeki değerler ve düğümlere ihtiyaç duyulan kırık çizgilerin oldukça kolay görülmesini sağlar. Parçalı polinom  $pp(x)$ ,  $(k+1)$  kırılma noktalı,  $n$ .dereceden  $k$  tane  $p_i(x)$  polinom kümesine sahip olan denklem,

$$pp(x) = p_i(x) \quad x \in [t_i, t_{i+1}], \quad i = 1, 2, \dots, k$$

biçiminde gösterebilir.  $p_i(x)$  fonksiyonu  $a_{i0}, a_{i1}, \dots, a_{in}$  katsayılarına sahiptir. Kırılma noktaları kümesi  $T = (t_1, t_2, \dots, t_{k+1})$  dir. Parçalı polinomun açık şekilde temsili aşağıdaki gibidir:

$$pp(x) = a_{i0} + a_{i1}(x-t_i) + a_{i2}(x-t_i)^2 + \dots + a_{in}(x-t_i)^n$$

$pp(x)$ 'i hesaplamak için doğru  $i$ .aralığı bulmamızı sağlayan eşitlik,

$$a_{0i} + \left[ \frac{a_{0i+1} - a_{0i}}{t_{i+1} - t_i} \right] (x-t_i)$$

ile verilir. Dolayısıyla değerlendirme boyunca karşılaştığımız  $(x-t_i)$ 'nin katsayısı olan  $a_{1i}$  yi hesaplamak zorunluluğu vardır.

Benzer bir gösterim yüksek dereceli splaynlar için mevcuttur. İlk  $[t_1, t_2]$  aralığında splayn, herhangi bir polinom olabilir; dolayısıyla  $a_{i0}, a_{i1}, \dots, a_{in}$  katsayılarına ihtiyaç vardır. Sonraki aralıkta  $p_2(x)$ 'in ilk  $(n-1)$  türevini  $t_2$  deki  $p_1(x)$ 'e eşitlemek için,  $t_2$  değerinde hesaplama yapılır. Bu polinomlar  $j=0, \dots, n$  için  $a_{1j}$ 'den hesaplanabilir.  $t_2$  değeri için  $p_2(x)$ 'in  $n$ .türevi,  $p_1(x)$ 'den  $p_2(x)$  polinomuna geçiş için bir sıçrama noktasıdır. Çünkü bu geçiş noktası bir polinomun diğeriyle birleştiği ortak bir yerdur. Bu şekilde 3. ve 4. ve sonraki düğümler hesaplanabilir. Dolayısıyla  $(i=2, 3, \dots, k)$   $b_i$  sıçramasının  $p_1(x)$  açılımında yer alan katsayıların bulunmasıyla  $pp(x)$  gösterimi elde edilmiş olur [20].

“Kesilmiş üslü ifade”(truncated power basis) kullanarak basit bir matematiksel biçim oluşturulmasıyla fonksiyon gösterimi geliştirilir.  $t, x$  fonksiyonunun bir parametresi olacak şekilde fonksiyon,

$$(x-t)_+^i = \begin{cases} (x-t)^i & , x \geq t \\ 0 & , x \leq t \end{cases}$$

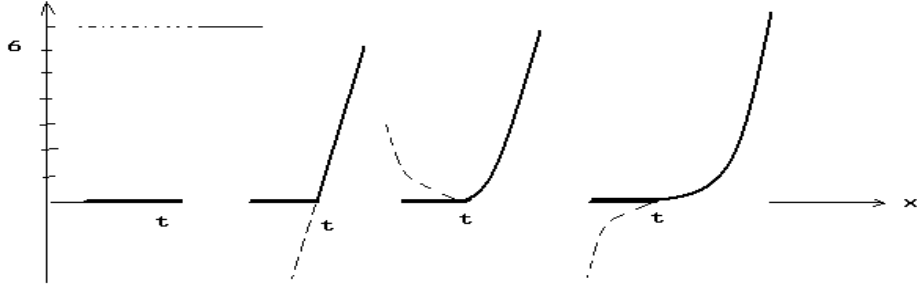
olarak yazılabilir.

Şekil III.3 'de  $i=0, 1, 2, 3, 4$  için fonksiyon çizimleri yer almaktadır.  $(x-t)_+^i$  nin türevleri ve integralleri kolayca hesaplanır. Özel olarak  $x^j$  fonksiyonuna benzer bir örnek aşağıda verilmiştir:

$$\frac{d^j (x-t)_+^i}{dx^j} = i(i-1)\dots(i-j+1)(x-t)_+^{i-j} \quad , j \leq i \quad \text{için.}$$

Böylece,  $(x-t)_+^i$  nin  $n$ .türevi adım fonksiyonudur.  $t$ ' nin sol tarafı 0, sağ tarafı  $n!$  olur.  $[t_2, t_3]$  ikinci aralığında  $p_1(x) + b_2(x-t_2)_+^n$  ,  $p_1(x)$ 'den  $p_2(x)$ 'e geçişi sağlayan bir fonksiyon adım fonksiyonudur. Bu fonksiyon,  $t_2$  de  $n$ .türevdeki  $n!.b_2$  nin bir sıçramasıyla  $p_1(x)$ 'den oluşturulur. Bu şekilde splaynın bir başka gösterimi de elde edilmiş olur:

$$pp(x) = \sum_{j=0}^n a_j (x-t_j)^j + \sum_{i=2}^k b_i (x-t_i)_+^n$$



**Şekil III.3** Kesilmiş Üslü İfade.  $t-2 \leq x \leq t+2$  için  $6(x-t)_+^0$ ,  $6(x-t)_+$ ,  $3(x-t)_+^2$  ve  $(x-t)_+^3$  fonksiyon çizimleri.

Şekil III.3'deki fonksiyonlar düzgünleştirilmiştir ve kübik olanın kırılma noktası gözle görülmeyecek şekildedir [20].

Özetle, splaynlar k.dereceden parçalı polinomlardır. Parçaların ortak noktaları olan düğümler önemli bir rol üstlenir. Çoğunlukla fonksiyon değerlerinin ve ilk (k-1) türevlerinin düğümlere uygunluk gösterdiği görülür; çünkü splayn, (k-1) sürekli türevlere sahip bir fonksiyondur.

İlk ve ikinci türevlerde sürekli kübik splayn (k=3), genellikle uygulamalı problemler için uygundur.  $h$  düğüm noktalı bir kübik splayn,  $t_1 < t_2 < \dots < t_h$ , aşağıdaki gibi yazılır:

$$E(y) = s(x) = \sum_{j=0}^3 \beta_{0j} x^j + \sum_{i=1}^h \beta_i (x-t_i)_+^3$$

Düğüm noktalarını açık bir biçimde gösteren splayn modelde süreklilik şartı aranmayabilir. Çünkü değişim noktalarına kadar zaten sürekli olmakla beraber ayrı parçalar şeklinde incelendiğinden bu koşulun aranmasına gerek yoktur. Yukarıda tanımlanan splayn fonksiyon buna örnek olarak verilebilir. Bu şekilde daha az süreklilik koşulu, daha uygun modelin elde edilmesini sağlar. Splayn fonksiyonların kesikli polinomlardan oluşması yani verinin bazı bölgelerde süreklilik şartının sağlamaması durumunda kullanılması ve bu noktalarda değişim noktalarının kullanılmasıyla oluşturulan farklı dereceli polinomların düğüm noktalarıyla birbirine bağlanması ilgili verinin kolayca çözülmesine yardımcı olur.

Düğüm noktalarının bilindiğini varsayalım. Eğer düğüm noktaları kestirilmiş parametreler ise problem bir doğrusal olmayan regresyon problemidir. Düğüm

noktalarının yeri bilindiğinde yukarıda tanımlanan eşitliğin oluşturulması, doğrusal en küçük karelerin açık bir uygulamasıyla tamamlanabilir [16].

Her bir parçadaki polinom derecesine, düğümlerin yeri ve sayısına karar vermek kolay değildir. Wold (1974), bölüm başına en az 4-5 veri noktasıyla mümkün birkaç düğümün var olabileceğini belirtmiştir. Buradaki uygulama önemlidir; çünkü splayn fonksiyonların büyük ölçüdeki önemli esnekliği, “overfit” (aşırı uyumu) veriyi çok kolay kılar. Ayrıca Wold, her bölümün bir büküm noktası ve bir ekstrem noktasından daha fazla olmayacağını belirtmiştir. Yani her bölümde sadece bir değişim söz konusudur. Ekstrem noktaları bölüm merkezinde yer almalı ve büküm noktaları da düğümlere yakın olmalıdır [16].

Pastor ve Guallar’ın makalelerindeki spline yönteminin ise anlaşılması daha kolaydır ve Matlab programında çözümlenmesi oldukça yalındır. Makalelerinde iki parçalı lojistik regresyon modellerini tanımlamışlardır. Epidemiyolojik çalışmalar hastalık riski ve maruz kalan değişken arasındaki ilişkiyi açıklamak için ideal olduğundan makalelerinde alkol alımı ile miyokard enfarktüsünü incelemişlerdir. Burada doz-tepki ilişkisinin, ilgili değişken  $x$ , bilinmeyen bir eşik düzeyine ulaştığı zaman beklenmedik şekilde değiştiği görülmüştür ve bu değişim noktası  $\lambda$  ile gösterilmiştir. Alışılmış doz-tepki analiz yöntemleri değişim noktalarının yerini veya güven aralığını tahmin etmede çıkarsama sağlamadığı için kategorik olarak ilgili  $x$  değişkenini alt aralıklara ayırmışlardır. Sınırlardaki değişim noktalarının kararsızlığından kaçınmak için splayn regresyon kullanarak her bir kategoride doğrusal veya karesel model bulmayı amaçlamışlardır. Splayn regresyon kategoriler içinde risk değişimine izin verdiğiinden doz-tepki analizi dediğimiz lojistik regresyon ilişkisini kestirmede bu yöntem kategorik analizden daha esnektir. Düğüm yeri araştırmacı tarafından keyfi olarak ayarlanmak zorundadır.

Pastor ve Guallar’ın yaptıkları çalışmada kullandıkları yöntem, bu çalışmadaki yöntemle paralel bir çalışmadır.

Buradaki lojistik regresyon modelindeki amaç, doz-tepki ilişkisinin belirgin bir biçimdeki değişimlerinde, maruz kalan  $x$  değişkeninde bir potansiyel değişim noktası kestirmektir. Lojit ve  $x$  değişkeni arasında değişim noktasının her iki yanında iki farklı polinomial fonksiyonun oluşturulması istenmektedir. İki segmentli lojistik regresyon lojit terimle açıklanabilir:

$$g(x, z_1, \dots, z_p) = \log \left( \frac{\pi(x, z_1, \dots, z_p)}{1 - \pi(x, z_1, \dots, z_p)} \right) \\ = f(x, \lambda) + \alpha_1 z_1 + \dots + \alpha_p z_p \quad (\text{III.1})$$

Burada  $\pi(x, z_1, \dots, z_p) = P(y=1 | x, z_1, \dots, z_p)$  verilen kovaryant kümesi için hastalık olasılığını belirtir.  $\lambda$  ise bilinmeyen değişim noktalarıdır. İki segmentli polinomial

fonksiyon  $f(x, \lambda) = \begin{cases} f_1(x) , & x < \lambda \\ f_2(x) , & \text{diğer durumlarda} \end{cases}$  ise farklı çok terimliler için

$f_1(x)$  ve  $f_2(x)$  vardır.  $\alpha_1, \dots, \alpha_p$  kovaryantların regresyon katsayılarıdır.

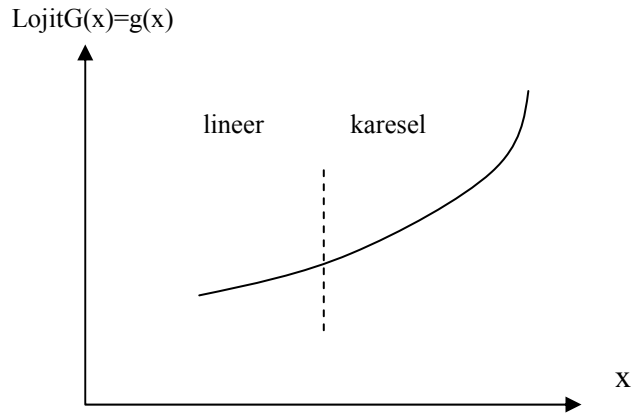
Pratik uygulamalar için doğrusal, karesel veya kübik gibi polinomlar sınıfı belirtilmelidir. Değişim noktasında düzgün geçişin istenilen özelliği en küçük karesel polinomlardan birini seçmektir. Diğer taraftan kareselden daha büyük dereceli polinomların katsayılarının epidemiyolojide yorumlanması zor olacaktır. İki segmentli lojistik regresyon için temel alternatifler doğrusal-karesel, karesel-doğrusal veya karesel-karesel gibi modellerdir.

Kovaryantların bulunduğu durumda doğrusal-karesel biçimdeki iki segmentli lojistik regresyon aşağıdaki gibi ifade edilebilir:

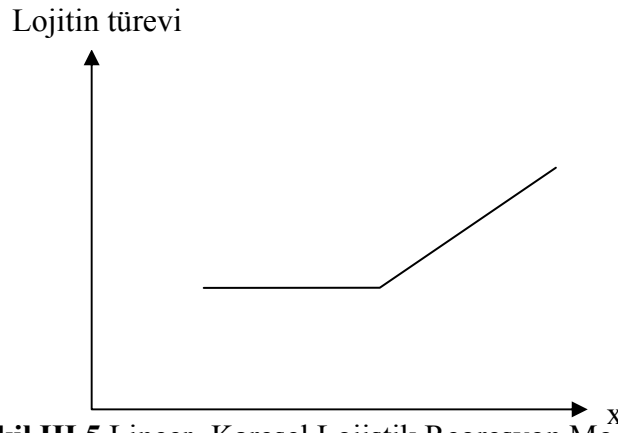
$$g(x, z_1, \dots, z_p) = \beta_0 + \beta_1 x + \beta_2 (x - \lambda)^2 I_{(x \geq \lambda)} + \alpha_1 z_1 + \dots + \alpha_p z_p \quad (\text{III.2})$$

$$I_{(x \geq \lambda)} = \begin{cases} 1 , & x \geq \lambda \\ 0 , & \text{diğer durumlarda} \end{cases}$$

$\lambda$  bilinmeyen değişim noktasını temsil eder.  $\beta_0, \beta_1, \beta_2$  etkilenen  $x$  değişkeninin bilinmeyen regresyon katsayılarıdır. Doğrusal-karesel lojitte iki segmentli modelde, etki,  $\lambda$  değişim noktasının alt parçasında kalan lojitte doğrusaldır. Bu parçada lojit  $\beta_1 + 2\beta_2(x - \lambda)$  şeklindedir.  $\lambda$  değişim noktasının üst tarafında, etki kareseldir ve  $\beta_0 + \beta_1 x + \beta_2 (x - \lambda)^2$  şeklinde verilir. Bu modelde lojitin türevi sürekli olduğundan iki parçaya karşı risk eğiliminin düzgün bir değişimi vardır. Eğim,  $\lambda$  değişim noktasının aşağısında ise  $\beta_1$  'dir. Değişim noktasının yukarısına doğru bir değişiklik başlar; burada eğim  $\beta_1 + 2\beta_2(x - \lambda)$  şeklinde doğrusal olacak şekilde sürekli bir değişim gözlenir.



Şekil III.4 Lineer- Karesel Lojistik Regresyon Modeli



Şekil III.5 Lineer- Karesel Lojistik Regresyon Modelin Türevi

Şekil III.4 'deki doğrusal-karesel lojistik regresyon modelini uydurmak için  $\alpha_1, \dots, \alpha_p$  kovaryantlarını ;  $\beta_0, \beta_1, \beta_2$  maruz kalan değişken katsayılarını ve  $\lambda$  değişim noktasını kestirmek gerekir. Bu bilinmeyen parametrelerin en çok olabilirlik kestiricileri  $\hat{\lambda}, \hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2$  ve  $\hat{\alpha}_1, \dots, \hat{\alpha}_p$  şeklindedir. Bunlar log-likelihood fonksiyonunu maksimize edecek şekilde kestirilir:

$$\begin{aligned}
 l &= \log(L) \\
 &= \sum_{i=1}^n [y_i \log(\pi_i) + (1 - y_i) \log(1 - \pi_i)] \\
 &= \sum_{i=1}^n [y_i g_i - \log(1 + e^{g_i})]
 \end{aligned}$$

$y_i$  i.kişi için hastalık durumunu belirtir.  $\pi_i = \pi(x_i, z_{i1}, \dots, z_{ip})$  ve  $g_i = g(x_i, z_{i1}, \dots, z_{ip})$  hastalık olasılığını ve i.kişi için bu olasılığın lojitini verir. (III.2) nolu eşitlik doğrusal olmadığından ve devamında sabit olmadığından standart iteratif olarak yeniden ağırlıklandırılmış en küçük kareler yöntemi bu modelin uydurulmasında kullanılamaz. Bunun yerine tüm model parametrelerini eşzamanlı olarak kestirmek için IRLS yönteminin başka bir versiyonu önerilmiştir [19].

Standart regresyondaki çıkarsamalar, en çok olabilirlik oran istatistiğinin asimptotik  $\chi^2$  dağılımından direkt olarak elde edilen olabilirlik temelli yöntemler ya da ML (maksimum likelihood) kestiricilerinin asimptotik normal dağılımı üzerine kurulu, Wald tipi yaklaşımlar kullanılmasına yönlendirilir. İki segmentli regresyonda, bu asimptotik özellikler belli bir değişim noktasının varlığına bağlıdır. Değişim noktasının varlığı için kavramsal bir test, doz-yanıtın doğrusal homojen bir örneği sayılan  $H_0: \beta_2 = 0$  hipotezinin test edilmesidir. Eğer  $\beta_2 = 0$  ise ikinci segmentin olmadığı sonucu ortaya çıkar; ama bu istenilen bir durum değildir. Çünkü segmente sahip bir veri üzerinde çalışılır. O halde  $\lambda$  iyi tanımlanmamıştır; parametre uzayının bir bozulması mevcuttur. Bu durumda ML kestiricileri asimptotik olarak normal değildir; Wald ve olabilirlik oran istatistiklerinin dağılımları sırasıyla standart normal ve  $\chi^2$  dağılımlarına yakınsamaz. Burada hipotez testleri veya güven aralıkları gibi tanımlı tüm çıkarsamalar,  $\lambda$ 'nın varlığına dayandırılır.  $\lambda$  için  $\%100(1 - \alpha)$  güven aralığında Wald modelinde  $\hat{\lambda} \pm Z_{1-\alpha/2} s(\hat{\lambda})$  kestirilir.  $s(\hat{\lambda})$ , modifiye edilmiş IRLS algoritmasının son iterasyonunda Fisher bilgi matrisinin tersinden elde edilen standart hatanın bir kestiricisidir.  $\lambda$  için daha düzgün güven aralığı olabilirlik oran istatistiğinin kullanılmasıyla elde edilebilir:

$$\hat{\lambda} = 2l(\hat{\lambda}, \hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2) - 2l(\lambda, \hat{\beta}_0(\lambda), \hat{\beta}_1(\lambda), \hat{\beta}_2(\lambda))$$

$\hat{\beta}_j(\lambda)$ ,  $\lambda$  verildiğinde ( $j=0, 1, 2$ )  $\beta_j$ 'nin ML kestirimidir.  $\chi^2_{(1; 1-\alpha)}$ , 1 serbestlik dereceli  $\%100(1 - \alpha)$  güven düzeyli ki-kare değeri olmak üzere,  $\lambda$  olabilirliğe dayalı güven aralığı  $\hat{\lambda} \leq \chi^2_{(1; 1-\alpha)}$  biçiminde gösterilir. Bu aralık genellikle Wald yaklaşımıyla oluşturulan terimlerde daha geçerlidir.

Spline fonksiyondaki gruplama farklı şekillerde gösterilebilir. Liu ve arkadaşlarının, doğum ile yumurtalık kanserindeki geçici azalmaya ilişkin çalışmalarında  $\lambda$  değişim noktası fark biçiminde gösterilmiş olup bu şekilde her bir değişim noktasından sonraki bölüm için değişken sayısının artmasına neden olmuşlardır. Makalelerinde kullandıkları spline model,

$$\beta_0 + \beta_p (\text{parity}) + \beta_a (\text{age}) + \beta_{\text{afd}} (\text{age at first delivery}) + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_1^3 + \beta_4 x_2^2 + \beta_5 x_3^2$$

dir. Burada yaş ve ilk doğumdaki yaş birer kategorik zaman değişkenleridir.  $x_1$  ise ilk doğumdan itibaren geçen zamanın sürekli değişkeni olmak üzere  $\lambda$  değişim noktaları 3 ve 38 olduğundan diğer değişkenler şöyledir [13]:

$$\begin{cases} x_2 = x_1 - 3 & , x_1 > 3 \\ x_2 = 0 & , \text{diğer durumlarda} \end{cases}$$

$$\begin{cases} x_3 = x_1 - 38 & , x_1 > 38 \\ x_3 = 0 & , \text{diğer durumlarda} \end{cases}$$

Pastor ve Guallar, tüm parametre uzayında kalan parametrelerin sürekli olacak şekilde ilk kısmi türevleriyle iki segmentli bir polinomial fonksiyon olan  $g(x, z_1, \dots, z_p) = f(x, \lambda) + \alpha_1 z_1 + \dots + \alpha_p z_p$  lojistik modeline uyan bir algoritma geliştirdiklerini belirtmektedirler.  $f(x, \lambda)$ 'nın özel parametrisasyonlarına örnek olarak,

$$f(x, \lambda) = \beta_0 + \beta_1 x + \beta_2 (x - \lambda)^2 I_{(x \geq \lambda)} \text{ doğrusal-karesel,}$$

$$f(x, \lambda) = \beta_0 + \beta_1 x + \beta_2 (x - \lambda)^2 I_{(x \leq \lambda)} \text{ karesel- doğrusal,}$$

$$f(x, \lambda) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 (x - \lambda)^2 I_{(x \geq \lambda)} \text{ karesel-karesel,}$$

modeller verilebilir. Buradaki başlıca amaç iki segmentli lojistik regresyon modelinin log-olabilirlik fonksiyonu,

$$l = \sum_{i=1}^n [y_i g_i - \log(1 + e^{g_i})] \quad (\text{III.3})$$

ifadesinin  $Q = (\lambda, \beta_0, \beta_1, \dots, \beta_q, \alpha_1, \dots, \alpha_p)$  bilinmeyen parametreler vektörüyle maksimize edilmesiyle en çok olabilirlik kestiricilerini  $\hat{Q} = (\hat{\lambda}, \hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_q, \hat{\alpha}_1, \dots, \hat{\alpha}_p)$  elde etmektir. Skor vektörü  $u(Q) = \partial l / \partial Q$ , aşağıdaki biçimde hesaplanabilir:

$$\frac{\partial l}{\partial Q} = \sum_{i=1}^n (y_i - \pi_i) \frac{\partial g}{\partial Q_r} \quad , \quad r = 1, \dots, p + q + 2 \quad (\text{III.4})$$

$g_i$  parametre uzayı olan  $Q$ 'ya göre birinci kısmi türevlere sahip sürekli bir ifadedir.  $E(y_i) = \pi_i$  ve  $\text{var}(y_i) = \pi_i(1 - \pi_i)$  olduğundan,  $u(Q)$  0 olasılığına sahiptir ve kovaryans matrisi  $I(Q)$  Fisher bilgi matrisine eşittir:

$$I_{r,s}(Q) = \left\{ E(u(Q)u^T(Q)) \right\}_{r,s} = \sum_{i=1}^n \pi_i(1 - \pi_i) \frac{\partial g_i}{\partial Q_r} \frac{\partial g_i}{\partial Q_s} \quad (\text{III.5})$$

Buradaki  $I_{r,s}(Q)$  Fisher matrisinin  $(r,s)$ ncinci elemanlarını gösterir  $(r, s = 1, \dots, p + q + 2)$ . En çok olabilirlik kestiricisi  $\hat{Q}$ ,  $u(Q) = 0$  olabilirlik denklemlerinin  $p + q + 2$  tanesinin çözümüyle elde edilir. Bu denklemler  $Q$  parametrelerinde doğrusal olmadıklarından iteratif yöntemlerle  $\hat{Q}$  elde edilmelidir. İki segmentli lojistik regresyon modeli Fisher skorlama yöntemi kullanılarak uyarlanabilir. Fisher bilgi matrisi  $I(Q^{(k)})$  ve skor vektörü  $u(Q^{(k)})$  olmak üzere  $Q$ 'nın mevcut kestirimi  $Q^{(k)}$  olsun.  $Q^{(k+1)}$  bir sonraki yeni parametre kestirimi olmak üzere  $I(Q^{(k)}) = (Q^{(k+1)} - Q^{(k)}) = u(Q^{(k)})$  denkleminde türetilir. Bu da (III.4) ve (III.5) denklemlerinin değerlendirilmesiyle aşağıdaki denklem sisteminin sonuçlarıdır  $(r = 1, \dots, p + q + 2)$ :

$$\sum_{i=1}^n \pi_i^{(k)} (1 - \pi_i^{(k)}) \frac{\partial g_i}{\partial Q_r} \Big|_{Q^{(k)}} \left[ \sum_{s=1}^{p+q+2} \frac{\partial g_i}{\partial Q_s} \Big|_{Q^{(k)}} (Q_s^{(k+1)} - Q_s^{(k)}) - \frac{y_i - \pi_i^{(k)}}{\pi_i^{(k)} (1 - \pi_i^{(k)})} \right] = 0 \quad (\text{III.6})$$

$g_i$ 'nin  $Q^{(k)}$ 'ya göre birinci dereceli Taylor serisi kullanılarak,

$$\left[ \frac{\partial g_i}{\partial Q} \Big|_{Q^{(k)}} \right]^T (Q^{(k+1)} - Q^{(k)}) = g_i^{(k+1)} - g_i^{(k)} \text{ ifadesi (III.6)'deki eşitlikten,}$$

$$\sum_{i=1}^n \pi_i^{(k)} (1 - \pi_i^{(k)}) \frac{\partial g_i}{\partial Q_r} \Big|_{Q^{(k)}} \left[ g_i^{(k+1)} - \left( g_i^{(k)} + \frac{y_i - \pi_i^{(k)}}{\pi_i^{(k)} (1 - \pi_i^{(k)})} \right) \right] = 0 \quad (\text{III.7})$$

şeklinde yazılabilir. Burada regresyon modelinin ağılıkları  $w_i^{(k)} = \pi_i^{(k)}(1 - \pi_i^{(k)})$  ve yanıt değişkenleri  $v_i^{(k)} = g_i^{(k)} + (y_i - \pi_i^{(k)})/w_i^{(k)}$  olmak üzere (III.7) nolu denklem ağırlıklandırılmış en küçük kareler denklemdir. Böylece  $Q^{(k+1)}$  ağırlıklandırılmış

modelden hesaplanabilir ve bu her bir iterasyondan yeniden hesaplandırılması yoluyla olur [19].

Gauss-Newton algoritmasının Hartley modifikasyonu, ağırlıklandırılmış iki segmentli regresyon modellerinin çözümünde kullanılabilir. Temel iteratif sürecin  $k$ 'ncü derecesiyle iç içe olan metodun  $l$ 'inci iterasyondaki  $Q$ 'nın bir kestirimi  $Q^{(k,l)}$  olsun. Her bir gözlem için  $Q^{(k,l)}$  çevresinde  $g_i(Q)$ 'nin birinci dereceden Taylor serisi  $g_i(Q) \cong g_i(Q^{(k,l)}) + \left[ \frac{\partial g_i}{\partial Q} \Big|_{Q^{(k,l)}} \right]^T (Q - Q^{(k,l)})$  denklemini verir.  $r_i(k,l) = v_i^{(k)} - g_i(Q^{(k,l)})$  artığı,  $r(k,l)$  artık vektör ve  $G(k,l) = \frac{\partial g_i}{\partial Q} \Big|_{Q^{(k,l)}}$ 'den  $(i,r)$  elemanlı bir matris olmak üzere  $r(k,l) = G(k,l)(Q - Q^{(k,l)}) + \varepsilon$  tüm gözlemlerden oluşan bir matris biçiminde veya aşağıdaki gibi açıklanabilir:

$$r_i(k,l) = \left[ \frac{\partial g_i}{\partial Q} \Big|_{Q^{(k,l)}} \right]^T (Q - Q^{(k,l)}) + \varepsilon \quad i=1, \dots, n. \quad (\text{III.8})$$

$\mathbf{W}(k) = \text{diag} \{ \pi_i^{(k)} (1 - \pi_i^{(k)}) \}$  ağırlıklarının diagonal matrisiyle doğrusallaştırılan regresyon modelinin çözümü -ordinary- ağırlıklandırılmış en küçük kareler yöntemi kullanılarak elde edilir.  $\mathbf{d}(k,l) = \left[ \mathbf{G}^T(k,l) \mathbf{W}(k) \mathbf{G}(k,l) \right]^{-1} \mathbf{G}^T(k,l) \mathbf{W}(k) r(k,l)$  dir. Dolayısıyla çözüm, segmentli polinomial regresyon modelleri için Gauss-Newton algoritmasının Hartley modifikasyonu gereğince, bir sonraki kestirici  $Q^{(k,l+1)}$ , iki segmentli regresyon modelinin ağırlıklandırılmış hata toplamını minimize eden  $Q^{(k,l)} + \eta \mathbf{d}(k,l)$  ifadesinde 0 ve 1 arasındaki bir skaler olan  $\eta(k,l)$ 'nin yer aldığı  $Q^{(k,l)} + \eta(k,l) \mathbf{d}(k,l)$  tarafından belirlenir.

Bu iteratif aşamaların her biri,  $Q^{(k,l+1)}$  ve  $Q^{(k,l)}$  arasındaki bağıl değişim önemsiz olduğu zaman sona erer. Son kestirici,  $Q^{(k+1)}$  sonraki iteratif sürecin başlangıç değeri olarak kullanılır.  $w_i^{(k+1)}$  ağırlığı ve  $v_i^{(k+1)}$  yanıt değişkeni  $Q^{(k+1)}$ 'da yeniden hesaplanmalıdır ve bu iteratif işlemleri kullanarak iki segmentli regresyon modelinin tekrar ağırlıklandırılmasına uygun olmalıdır. Bu süreç yakınsama olana kadar tekrarlanır.

Burada sözü edilen algoritmanın doğru bir yakınsamasını elde etmek için  $Q^{(0)}$  başlangıç kestiriminin seçimi çok önemlidir; çünkü (III.3) log-olabilirlik fonksiyonu yerel maksimuma yol açabilir. Başlangıç noktalarını elde etmek için iyi

bir alternatif, grid-search yaklaşımının kullanılmasıdır. Bu yöntem  $x$  maruz kalan değişkenin aralığında yer alan eşit aralıklı noktaları önceden belirler. Her bir nokta için iki segmentli lojitli bir splayn lojistik regreyon bulunabilir ve log-olabilirliği değerlendirilir. Log-olabilirlikteki nokta,  $Q^{(0)}$  başlangıç değeri olarak splayn lojistik regresyondan diğer parametre kestirimleriyle birlikte kullanılan en yüksek değere ulaşır [19].

Kullanılan splayn lojistik model iki segmentli olmasa da parametrelerin bulunmasına kadar olan işlemler Gauss-Newton algoritmasıyla aynıdır. (III.7) nolu denklemde açık bir şekilde gösterilmiş olmasa da bu yöntemle ağırlıklandırılmış olarak çözümlenmektedir. Sonuç olarak iyi bir başlangıç değeriyle Gauss-Newton algoritmasının kullanımı önemli rol oynamaktadır.

Splayn regresyon doz-yanıt ilişkisini kestirmek için kullanılan geleneksel tüm metotlardan daha avantajlıdır. Kompleks dağılımlara uyarlanabilir ve grafiklerde değişime daha elverişlidir [4 ].

## BÖLÜM IV

### IV.1 UYGULAMA-1

Bu uygulamada, 243 kişilik hemoglobin verisi üzerinden splayn yöntem aracılığıyla doğrusal olmayan lojistik regresyon analizi incelendi. Bu çalışmanın amacı, hemoglobin eksikliğinin etkili olduğu hastaların yaşama şansını tahmin edecek bir doğrusal olmayan lojistik regresyon modeli kurmaktır. Öncelikle bağımsız değişken olarak modelde yer alan hemoglobin hakkında bilgi vermek doğru olacaktır.

Alyuvarlara karbondioksit ve oksijen taşıyabilme yeteneği kazandıran hemoglobin molekülleridir. Hemoglobin aynı zamanda kana kırmızı rengini de veren maddedir.

Hemoglobin yapısında bulunan zincirlerden yalnızca birinin bile aminoasit yapısında ortaya çıkan bir değişiklik, anormal hemoglobin üretimine neden olur. Bu yapısal değişiklik çoğunlukla belirti vermezken, bazı durumlarda hemoglobinin oksijen taşımamasını önemli ölçüde etkileyerek hastanın yaşamını tehdit edebilir.

Anemi (Kansızlık), kan hemoglobin düzeyinde veya kırmızı kan hücreleri sayısında azalma ve sonucunda ortaya çıkan bulgulardır. Bir başka ifadeyle, anemi hemoglobin miktarının yaş ve cinsiyete göre dünya sağlık örgütü tarafından kabul edilen kriterlerin altında kalmasıdır. Bu kriterler erişkin erkeklerde 13 g/dL, kadınlarda 12 g/dL nin altı kabul edilir. 6 ay ile 6 yaş arası çocuklarda 11 g/dL nin, 6-14 yaşlarda 12 g/dL nin altı anemidir [24].

Ortaya çıkan şikayetler ve saptanan bulgular doku ve hücrelere yetersiz oksijen taşınmasına bağlı olarak gelişmektedir. Anemili hastalarda yorgunluk, hafif

çarpıntı ve nefes darlığı gelişebilir. İleri düzeyde bir anemide ise, bütün bu bulgular, istirahat halinde görülmesinin yanı sıra; kulak çınlaması, baş dönmesi, baş ağrısı, uyuma güçlüğü, iştahsızlık, kilo kaybı, adet kanamalarının düzensizliği veya fazlalığı, adet görmeme ve iktidarsızlık gibi bulgular ortaya çıkabilir. Aneminin sık görülen bulgularından çarpıntı, anemi yüzünden dokularda oluşan oksijen açlığını gidermek amacıyla, kalbin atım hızını ve her atımda pompaladığı kan miktarını artırması nedeniyle ortaya çıkar. Buna rağmen dokularda yeterli oksijen sağlanamıyorsa, solunum sayısının artması ve nefes darlığı ortaya çıkar. Uzayan anemilerde ve yaşlı kişilerde veya kalp hastalığı olanlarda kalp yetmezliğine ait bulgular gelişebilir. Anemideki en belirgin bulgulardan birisi de solukluktur. Aneminin şiddetine bağlı olarak, ağız ve göz kapağı içindeki deride ilk olarak fark edilebilen solukluk, aneminin ilerlemesi ile avuç içinde, tırnak yataklarında ve deride de belirginleşir. Anemiye yol açan nedene bağlı olarak çok çeşitli bulgular gelişebilir.

En sık rastlanan anemi türleri demir eksikliğine bağlı anemi, folik asit eksikliğine bağlı anemi, Vitamin B-12 eksikliği anemisi dir.

Anemi gibi kan hücrelerinin sayısındaki deformasyona bağlı olarak gelişen ve tehlike içeren bir durum da lösemidir. Bunun tipik bir örneği de hemoglobin miktarındaki azalmayla görülen Akut Myeloid Lösemi (AML)'dir. Olgunlaşmamış bu hücreler kemik iliğinde çok yüksek sayılara ulaşırlar ve normal kan hücrelerinin üretimini azaltırlar. Sonuçta anemi (kansızlık - kırmızı kan hücresi üretiminde azalma) ve sık enfeksiyona yakalanma (beyaz kan hücresi üretiminde azalma) durumu ortaya çıkabilir. Ergenlik çağında ve 20 li yaşlarda saptanan lösemilerin %50 sini, yetişkinlerdeki lösemilerin de %20 sini AML oluşturur [ 24].

Genel olarak lösemiler tüm kanserlerin %2 sini oluştururlar. Erkeklerde lösemi daha sık gözlenmektedir. Ayrıca beyaz ırkta da daha sıktır. Yetişkinlerde lösemi tanısı konma sıklığı çocuklardan 10 kat daha fazladır ve risk yaşla birlikte artar. Çocuklar arasında ise 4 yaş altında daha sık gözlenir. Bu çalışmada yetişkin kadınları içeren 243 kişilik bir veri kümesinde hemoglobin değerlerine ait bulgular incelendi.

Çalışmada bağımsız değişken hemoglobin değerlerinden oluşmaktadır; bağımlı değişken ise hemoglobin eksikliğinin neden olduğu anemi ve lösemi risklerinin mevcut olduğu  $y=1$  ve risk barındırmayan  $y=0$  şeklinde iki sonuçlu değişkendir.

Burada riske neden olan diğer kan sayımı bulguları dikkate alınmadı, tek değişkenli model üzerinde çalışma gerçekleştirildi.

SAS programında Mulla'nın önerdiği algoritmaya göre hemoglobin verisi çözümlendi. Mulla yaptığı çalışmada bağımsız değişken olan x değişkenini değişim aralıklarına göre gruplandırarak kategorik şekle getirmiştir. Bu çalışmada, hemoglobinin belirgin etki noktaları ise 8.6 g/dL ve 12 g/dL'dir. Bu değişim aralıklarına göre kategorilendirme sonucu her kategori için yeni isimlendirmeler aşağıdaki biçimdedir:

```
if hemoglobin1<8.6 then hgbcats=1;
if hemoglobin1>=8.6 and hemoglobin1<=11.9 then hgbcats=2;
if hemoglobin1>=12.0 then hgbcats=3;

if hemoglobin1>8.6 then hgb2=hemoglobin1-8.6;
                                else hgb2=0;

if hemoglobin1>12.0 then hgb3=hemoglobin1-12.0;
                                else hgb3=0;
```

Burada üç tane bölge vardır. İlk bölge 8.6 g/dL değerinin altında kalan bölgedir ve ikinci bölge  $8.6 \text{ g/dL} \leq \text{hemoglobin} < 12 \text{ g/dL}$  ve son olarak üçüncü bölge ise  $12 \text{ g/dL} \leq \text{hemoglobin}$  şeklindedir. Bu algorithmada ilk bölgedeki bağımsız değerlere 'hemoglobin1' denilirse hgb1 adlandırılmasıyla, ikinci bölgedeki bağımsız değişken kategorilendirme gereği  $\text{hgb2}=\text{hgb1}-8.6$ , üçüncü bölgedeki değişken ise  $\text{hgb3}=\text{hgb1}-12.0$  olarak yazılabilir. Değişken isimlendirmesi bir önceki bölümde bahsedilen Liu ve arkadaşlarının kullandıkları tanımlamaya benzer şekildedir. Doğrusal olmayan bir model yapısı var olduğundan doğrusal olmama durumunu modelde gösterebilmek adına çeşitli algoritma atamaları oluşturularak PROC LOGISTIC, Version 9.1 SAS / STAT paket programı kullanıldı.

Hemoglobin çalışması iki farklı model yapısı üzerinden incelendi. Modellerde değişim noktasının aynı olmasıyla model yapılandırmasının önemini belirtmesi amaçlandı.

Değişim noktalarının ayırdığı bölgeler kategorilendirme yapıları aşağıdaki modelleme yapıldı. Model\_1 çalışması için,

$$\begin{aligned} \text{hgbz1} &= \text{hemoglobin1}; \\ \text{hgbz2} &= (\text{hemoglobin1}^{**2}) - (\text{hgb3}^{**2}); \\ \text{hgbz3} &= (\text{hgb2}^{**2}) - (\text{hgb3}^{**2}); \end{aligned}$$

model\_2 çalışması için de,

$$\begin{aligned} \text{hgbz1} &= \text{hemoglobin1}; \\ \text{hgbz2} &= (\text{hemoglobin1}^{**2}) - (\text{hgb3}^{**2}); \\ \text{hgbz3} &= (\text{hgb2}^{**2}) + (\text{hgb3}); \end{aligned}$$

yapılandırılması oluşturulmuştur. Buna göre 243 gözlemin her iki model için yanıt profili aşağıdaki gibidir:

**Tablo IV.1** HGB Modelleri için Yanıt Profili.

Ordered Value	Outcome	Total Frequency
1	1	46
2	0	197

Her iki model için model uygunluk istatistiği birbirine benzer sonuçlar vermiştir. Hangi modelin daha iyi olacağını belirlemede henüz yeterli bir bilgi olmamakla beraber sonuçları aşağıdaki tabloda yer almaktadır:

**Tablo IV. 2** HGB Model\_1 için Model Uygunluk İstatistiği.

Model Fit Statistics		
Criterion	Intercept Only	Intercept and Covariates
AIC	237.811	87.672
SC	241.304	101.644
-2 Log L	235.811	79.672

**Tablo IV. 3** HGB Model\_2 için Model Uygunluk İstatistiği.

Model Fit Statistics		
Criterion	Intercept Only	Intercept and Covariates
AIC	237.811	88.530
SC	241.304	102.502
-2 Log L	235.811	80.530

Modeller için sıfır hipotez testi aşağıda verilmiştir.

**Tablo IV. 4** HGB Model\_1 için Sıfır Hipotez Testi.

Testing Global Null Hypothesis: BETA=0			
Test	Chi-Square	DF	Pr > ChiSq
Likelihood Ratio	156.1386	3	<.0001
Score	147.3557	3	<.0001
Wald	24.7096	3	<.0001

**Tablo IV. 5** HGB Model\_2 için Sıfır Hipotez Testi.

Testing Global Null Hypothesis: BETA=0			
Test	Chi-Square	DF	Pr > ChiSq
Likelihood Ratio	155.2805	3	<.0001
Score	141.9816	3	<.0001
Wald	28.3060	3	<.0001

Her iki model için dikkat edilmesi gereken belirleyici aşamalardan biri en çok olabilirlik kestirimi analizidir. Model\_1’de parametrelerden birinin kesim noktası ve standart hatasının oldukça yüksek çıkması düşündürücüdür. Modeller arasında şimdilik standart hata kestirimine bakarak model\_2’nin daha iyi sonuç verdiği söylenilebilir. Kestirilmiş katsayının standart hatasına bölünmesi sonucu elde edilen standart normal dağılım olarak tanımlanan Wald istatistiğinin karesinin alınmasıyla Wald ki-kare testi elde edilir ve bu test modele girmesi muhtemel değişkenlerin ne kadar gerekli olduğunu sınar. Bu durumda sıfır hipotezi ile, ilgili değişkenin modelde sıfır olup reddedilmesi gerekip gerekmediğine karar verilir. Burada  $p$  değerleri Wald ki-kare istatistiğinden büyük çıktığı için Sıfır hipotezi reddedilir. Bir başka deyişle ilgili parametrenin sıfır olarak kabul edilmesiyle oluşturulan hipotez koşulu sağlanmaz ve bu da demektir ki parametrenin modelde bulunması, yani sıfır olarak alınmaması gereklidir. Modelde yer alması önemlidir.

**Tablo IV. 6** HGB Model\_1 için En Çok Olabilirlik Kestirimi.

Analysis of Maximum Likelihood Estimates					
Parameter	DF	Estimate	Standard Error	Wald Chi-Square	Pr >ChiSq
Intercept	1	541.9	1023.3	0.2804	0.5964
hgbz1	1	-125.7	238.6	0.2776	0.5983
hgbz2	1	7.3134	13.9093	0.2765	0.5990
hgbz3	1	-8.2594	14.1303	0.3417	0.5589

**Tablo IV. 7** HGB Model\_2 için En Çok Olabilirlik Kestirimi.

Analysis of Maximum Likelihood Estimates					
Parameter	DF	Estimate	Standard Error	Wald Chi-Square	Pr >ChiSq
Intercept	1	90.4324	205.7	0.1933	0.6602
hgbz1	1	-20.1474	48.6131	0.1718	0.6785
hgbz2	1	1.1439	2.8675	0.1592	0.6899
hgbz3	1	-1.9040	3.1569	0.3638	0.5464

Karşılaştıracağımız bir diğer özellik ise Odds oran kestirimidir. Hastalık riskini diğer değişkenlere göre oransal ifade eden Odds kestirimi modelin yorumlanması ve kabul edilebilirliğinde önemli rol oynar. Model\_1'deki nokta kestirimleri güven aralıklarının dışında kalmaktadır. Bu nedenle model\_2'nin tercih edilmesi daha uygun olacaktır.

**Tablo IV. 8** HGB Model\_1 için Odds Ratio Kestirimi.

Odds Ratio Estimates			
Effect	Point Estimate	95% Wald Confidence Limits	
hgbz1	<0.001	<0.001	>999.999
hgbz2	>999.999	<0.001	>999.999
hgbz3	<0.001	<0.001	>999.999

**Tablo IV. 9** HGB Model\_2 için Odds Ratio Kestirimi.

Odds Ratio Estimates			
Effect	Point Estimate	95% Wald Confidence Limits	
hgbz1	<0.001	<0.001	>999.999
hgbz2	3.139	0.011	866.219
hgbz3	0.149	<0.001	72.491

Buradaki odds oran kestiriminde daha iyi model olarak kabul edilen model\_2'de hgbz1 değerinin kestirimi güven aralığı sınırında çıkmaktadır.

HGB model\_2 çalışmasının sonuçları model istatistiği, en çok olabilirlik kestirimi ve odds oran kestirim sonuçlarına dayanarak model\_1'e göre çok daha kabul edilebilir bir çalışmadır. Ayrıca ROC (receiver operating characteristic ) eğrisinin altında kalan alan olarak tanımlanan c ( concordance rate ) ile sembolize edilen uyumluluk oranının 1'e yakın olması modelin kabul edilebilirliğini istatistiksel olarak desteklemektedir. Hemoglobin çalışmasında model\_2'nin c değeri 0.955 olarak sonuç vermiştir. Model\_2 için Hemoglobin SAS Sonucu, her bir bağımsız değişkenle risk durumları, Ek-1'de verilmiştir. Model\_2'nin SAS program sonucu aşağıda gösterilmektedir:

The LOGISTIC Procedure

Model Information

Data Set	WORK. STUDY
Response Variable	outcome outcome
Number of Response Levels	2
Model	binary logit
Optimization Technique	Fisher's scoring

Number of Observations Read	243
Number of Observations Used	243

Response Profile

Ordered Value	outcome	Total Frequency
1	1	46
2	0	197

Probability modeled is outcome=1.

Model Convergence Status

Quasi -complete separation of data points detected.

Model Fit Statistics

Criterion	Intercept Only	Intercept and Covariate
AIC	237.811	88.530
SC	241.304	102.502
-2 Log L	235.811	80.530

Testing Global Null Hypothesis: BETA=0

Test	Chi-Square	DF	Pr > Chi Sq
Likelihood Ratio	155.2805	3	<.0001
Score	141.9816	3	<.0001
Wald	28.3060	3	<.0001

Analysis of Maximum Likelihood Estimates

Parameter	DF	Standard Estimate	Wald Error	Chi-Square	Pr > Chi Sq
Intercept	1	90.4324	205.7	0.1933	0.6602
hgbz1	1	-20.1474	48.6131	0.1718	0.6785
hgbz2	1	1.1439	2.8675	0.1592	0.6899
hgbz3	1	-1.9040	3.1569	0.3638	0.5464

Odds Ratio Estimates

Effect	Point Estimate	95% Wald Confidence Limits
hgbz1	<0.001	<0.001 >999.999
hgbz2	3.139	0.011 866.219
hgbz3	0.149	<0.001 72.491

Association of Predicted Probabilities and Observed Responses

Percent Concordant	93.8	Somers' D	0.910
Percent Discordant	2.8	Gamma	0.943
Percent Tied	3.4	Tau-a	0.281
Pairs	9062	c	0.955

## IV.1 UYGULAMA-2

Uygulama-1’de SAS programında splayn model yapılandırılması incelendi. Uygulama-2’de ise Matlab programında bir başka uygulama olan serum albümin incelenecektir. Kan plazmasının veya serumunun en önemli bölümünü oluşturan albumin, dokuların temel maddelerinden olup önemli bir proteindir. Albümin, suda çözünürlükleri düşük olan yağ asitlerinin kandaki başlıca taşıyıcısıdır. Bunun yanı sıra, oksijen serbest radikallerine bağlanarak bunları kontrol altına alır, ayrıca bilirubin (hem molekülünün yıkımı sırasında ortaya çıkar) gibi suda çözünmeyen bazı zehirli metabolizma ürünlerine bağlanarak onları zararsız kılar. Kandaki albümin seviyesinin düşüklüğü ciddi risk teşkil eder. Dolayısıyla bu çalışmada,  $x$  değişkeni bağımsız değişken olmak üzere 117 kişiden oluşan serum albümin değeridir.  $y$  değişkeni de albümin eksikliğinden kaynaklanan ölüm riskidir. Değişim noktaları da tıbbi açıdan önemli sayılan 2.5g/100mL ve 3.5g/100mL değerleridir. Bu değerlerin altında kalan noktalar riski temsil eder. Norm aralık birinci dereceden 3.5g/100mL’dir. Bu değerden daha düşük seviyeler, karaciğerde sentezlenen protein türevi olan albüminin artık sentezlenmemesi sebebiyle vücutta ödem oluşturur. İleri derecede albümin eksikliğinin ölüme neden olduğu bilinmektedir. Belirleyici değişim noktası 3.5g/100mL’den sonra ciddi tehlike arz eden seviye 2.5g/100mL olarak görülmüştür [17].

Matlab programında Pastor ve Guallar’ın model yapısına benzer bir model serum albümin çalışması için oluşturuldu. Bu modelde değişim noktaları lamda1 ve lamda2 olarak ifade edilmiştir. Buna göre model bu değişim noktalarından öncesi ve sonrası olarak es1 ve es2 olarak belirtilmiştir. Daha açık bir ifadeyle,

```
if x(i)>=2.5, es1(i)=1;
    else es1(i)=0;
end

if x(i)>=3.5, es2(i)=1;
    else es2(i)=0;
end
```

$$es1_i = \begin{cases} 1 & , x_i \geq 2.5 \\ 0 & , \text{diğer durumlarda} \end{cases}$$

$$es2_i = \begin{cases} 1 & , x_i \geq 3.5 \\ 0 & , \text{diğer durumlarda} \end{cases}$$

şeklinde kategorilendirilebilir. Model fonksiyon ise

$$\pi(x_i) = \text{epi}(i) = \frac{\exp(b_0 + b_1(b_2 - x_i)^3 + b_3(x_i - \lambda_1)^2 es1_i + b_4(x_i - \lambda_2) es2_i)}{1 + \exp(b_0 + b_1(b_2 - x_i)^3 + b_3(x_i - \lambda_1)^2 es1_i + b_4(x_i - \lambda_2) es2_i)}$$

biçiminde alınmıştır. Parametreleri kestirmede doğal logaritması alınmış denklem ile çalışmak daha kolay olduğundan Bölüm II'de de belirtildiği gibi olabilirlik denkleminin logaritması alınır. Logaritması alınmış denkleme ML yöntemi uygulanır. Buradan elde edilen denklem sistemi Gauss-Newton yöntemi aracılığıyla çözümlenir; her bir parametreye göre türev alınarak Newton denklemindeki F matrisi ve jakobiyen matrisi oluşturulur. Eğer jakobiyen matrisinin tersi singüler çıkıyorsa burada bir çoklu iç ilişkiden söz edilebilir. Bağımsız değişkenler arasında çoklu iç ilişkinin oluşmasının üç temel nedeni vardır. Deney tasarımındaki planlamadan ya da zayıf gözlemsel verilerden kaynaklanır, ikincisi ise bağımsız değişkenlerin kuvvetleri ya da çarpımları gibi matematiksel işlemler sonucu oluşturulan yeni değişkenlerin bulunduğu model yapısından kaynaklanır. Bir diğeri de bağımsız değişkenler üzerindeki kısıtların neden olduğu çoklu iç ilişkidir. Bu durumda SVD ( Singular Value Decomposition ) ayrışımı uygulanır. Matlab programında yazılmış olan model sistemimiz aşağıda gösterilmektedir:

```

clc; clear
b0=-0.3; b1=0.3; b2=0.6; b3=0.5; b4=4; l amda1=2.5; l amda2=3.5; eps=0.01;
Jacob=zeros(7); i l kteta=zeros(7,1);
i l kteta(1)=b0; i l kteta(2)=b1; i l kteta(3)=b2; i l kteta(4)=b3; i l kteta(5)=b4;
i l kteta(6)=l amda1; i l kteta(7)=l amda2; enbsayi =1;

load veri_zuber.m
[n,k1]=si ze(veri_zuber)
y=veri_zuber(:,1); x=veri_zuber(:,2);

for i ter=1: enbsayi

for i =1:n

    i f x(i)>=2.5, es1(i)=1;
        e l s e es1(i)=0;
    end

    i f x(i)>=3.5, es2(i)=1;
        e l s e es2(i)=0;
    end
end

for i =1:n
    eus(i)=exp(b0+ b1*(b2-x(i))^3+ b3*(x(i)-l amda1)^2*es1(i)+ b4*(x(i)-l amda2)*es2(i));
    eust(i)=1+eus(i);
    epi (i)=eus(i)/eust(i);
    cpi (i)=epi (i)*(1-epi (i));    end

toppi 11=0; toppi 12=0; toppi 13=0; toppi 14=0; toppi 15=0; toppi 16=0; toppi 17=0;

```

```

toppi 21=0; toppi 22=0; toppi 23=0; toppi 24=0; toppi 25=0; toppi 26=0; toppi 27=0;
toppi 31=0; toppi 32=0; toppi 33=0; toppi 34=0; toppi 35=0; toppi 36=0; toppi 37=0;
toppi 41=0; toppi 42=0; toppi 43=0; toppi 44=0; toppi 45=0; toppi 46=0; toppi 47=0;
toppi 51=0; toppi 52=0; toppi 53=0; toppi 54=0; toppi 55=0; toppi 56=0; toppi 57=0;
toppi 61=0; toppi 62=0; toppi 63=0; toppi 64=0; toppi 65=0; toppi 66=0; toppi 67=0;
toppi 71=0; toppi 72=0; toppi 73=0; toppi 74=0; toppi 75=0; toppi 76=0; toppi 77=0;

```

```

for i=1:n

```

```

toppi 11=toppi 11+ (-cpi (i));
toppi 12=toppi 12+ (-(b2-x(i))^3*cpi (i));
toppi 13=toppi 13+ (-3*b1*(b2-x(i))^2*cpi (i));
toppi 14=toppi 14+ (-(x(i)-l amda1)^2*es1(i)*cpi (i));
toppi 15=toppi 15+ (-(x(i)-l amda2)*es2(i)*cpi (i));
toppi 16=toppi 16+ (2*b3*(x(i)-l amda1)*es1(i)*cpi (i));
toppi 17=toppi 17+ (b4*es2(i)*cpi (i));

```

```

toppi 21=toppi 21+ (-(b2-x(i))^3*cpi (i));
toppi 22=toppi 22+ (-(b2-x(i))^6*cpi (i));
toppi 23=toppi 23+ (3*y(i)*(b2-x(i))^2- 3*(b2-x(i))^2*epi (i)- 3*b1*(b2-x(i))^5*cpi (i));
toppi 24=toppi 24+ (-(b2-x(i))^3*(x(i)-l amda1)^2*es1(i)*cpi (i));
toppi 25=toppi 25+ (-(b2-x(i))^3*(x(i)-l amda2)*es2(i)*cpi (i));
toppi 26=toppi 26+ (2*b3*(b2-x(i))^3*(x(i)-l amda1)*es1(i)*cpi (i));
toppi 27=toppi 27+ (b4*(b2-x(i))^3*es2(i)*cpi (i));

```

$\text{toppi } 31 = \text{toppi } 31 + (-3*b1*(b2-x(i))^2*cpi(i));$   
 $\text{toppi } 32 = \text{toppi } 32 + (3*y(i)*(b2-x(i))^2 - 3*(b2-x(i))^2*epi(i) - 3*b1*(b2-x(i))^5*cpi(i));$   
 $\text{toppi } 33 = \text{toppi } 33 + (6*y(i)*b1*(b2-x(i)) - 6*b1*(b2-x(i))*epi(i) - 9*b1^2*(b2-x(i))^4*cpi(i));$   
 $\text{toppi } 34 = \text{toppi } 34 + (-3*b1*(b2-x(i))^2*(x(i)-l\text{ amda1})^2*es1(i)*cpi(i));$   
 $\text{toppi } 35 = \text{toppi } 35 + (-3*b1*(b2-x(i))^2*(x(i)-l\text{ amda2})*es2(i)*cpi(i));$   
 $\text{toppi } 36 = \text{toppi } 36 + (6*b1*b3*(b2-x(i))^2*(x(i)-l\text{ amda1})*es1(i)*cpi(i));$   
 $\text{toppi } 37 = \text{toppi } 37 + (3*b1*b4*(b2-x(i))^2*es2(i)*cpi(i));$

$\text{toppi } 41 = \text{toppi } 41 + (-(x(i)-l\text{ amda1})^2*es1(i)*cpi(i));$   
 $\text{toppi } 42 = \text{toppi } 42 + (-(b2-x(i))^3*(x(i)-l\text{ amda1})^2*es1(i)*cpi(i));$   
 $\text{toppi } 43 = \text{toppi } 43 + (-3*b1*(b2-x(i))^2*(x(i)-l\text{ amda1})^2*es1(i)*cpi(i));$   
 $\text{toppi } 44 = \text{toppi } 44 + (-(x(i)-l\text{ amda1})^4*es1(i)^2*cpi(i));$   
 $\text{toppi } 45 = \text{toppi } 45 + (-(x(i)-l\text{ amda1})^2*(x(i)-l\text{ amda2})*es1(i)*es2(i)*cpi(i));$   
 $\text{toppi } 46 = \text{toppi } 46 + (-2*y(i)*(x(i)-l\text{ amda1})*es1(i) + 2*(x(i)-l\text{ amda1})*es1(i)*epi(i) + 2*b3*(x(i)-l\text{ amda1})^3*es1(i)^2*cpi(i));$   
 $\text{toppi } 47 = \text{toppi } 47 + (b4*(x(i)-l\text{ amda1})^2*es1(i)*es2(i)*cpi(i));$

$\text{toppi } 51 = \text{toppi } 51 + (-(x(i)-l\text{ amda2})*es2(i)*cpi(i));$   
 $\text{toppi } 52 = \text{toppi } 52 + (-(b2-x(i))^3*(x(i)-l\text{ amda2})*es2(i)*cpi(i));$   
 $\text{toppi } 53 = \text{toppi } 53 + (-3*b1*(b2-x(i))^2*(x(i)-l\text{ amda2})*es2(i)*cpi(i));$   
 $\text{toppi } 54 = \text{toppi } 54 + (-(x(i)-l\text{ amda1})^2*(x(i)-l\text{ amda2})*es1(i)*es2(i)*cpi(i));$   
 $\text{toppi } 55 = \text{toppi } 55 + (-(x(i)-l\text{ amda2})^2*es2(i)^2*cpi(i));$   
 $\text{toppi } 56 = \text{toppi } 56 + (2*b3*(x(i)-l\text{ amda1})*(x(i)-l\text{ amda2})*es1(i)*es2(i)*cpi(i));$

```

toppi 57=toppi 57+ (-y(i)*es2(i)+ es2(i)*epi (i)+ b4*(x(i)-l amda2)*es2(i)^2*cpi (i));

toppi 61=toppi 61+ (2*b3*(x(i)-l amda1)*es1(i)*cpi (i));
toppi 62=toppi 62+ (2*b3*(b2-x(i))^3*(x(i)-l amda1)*es1(i)*cpi (i));
toppi 63=toppi 63+ (6*b1*b3*(b2-x(i))^2*(x(i)-l amda1)*es1(i)*cpi (i));
toppi 64=toppi 64+ (-2*y(i)*(x(i)-l amda1)*es1(i)+ 2*(x(i)-l amda1)*es1(i)*epi (i)+ 2*b3*(x(i)-
l amda1)^3*es1(i)^2*cpi (i));
toppi 65=toppi 65+ (-2*b3*(x(i)-l amda1)*(x(i)-l amda2)*es1(i)*es2(i)*cpi (i));
toppi 66=toppi 66+ (2*y(i)*b3*es1(i)- 2*b3*es1(i)*epi (i)- 4*b3^2*(x(i)-
l amda1)^2*es1(i)^2*cpi (i));
toppi 67=toppi 67+ (-2*b3*b4*(x(i)-l amda1)*es1(i)*es2(i)*cpi (i));

toppi 71=toppi 71+ (b4*es2(i)*cpi (i));
toppi 72=toppi 72+ (b4*(b2-x(i))^3*es2(i)*cpi (i));
toppi 73=toppi 73+ (3*b1*b4*(b2-x(i))^2*es2(i)*cpi (i));
toppi 74=toppi 74+ (b4*(x(i)-l amda1)^2*es1(i)*es2(i)*cpi (i));
toppi 75=toppi 75+ (-y(i)*es2(i)+ es2(i)*epi (i)+ b4*(x(i)-l amda2)*es2(i)^2*cpi (i));
toppi 76=toppi 76+ (-2*b3*b4*(x(i)-l amda1)*es1(i)*es2(i)*cpi (i));
toppi 77=toppi 77+ (-b4^2*es2(i)^2*cpi (i));
end

```

```
Jacob=[toppi 11 toppi 12 toppi 13 toppi 14 toppi 15 toppi 16 toppi 17;
        toppi 21 toppi 22 toppi 23 toppi 24 toppi 25 toppi 26 toppi 27;
        toppi 31 toppi 32 toppi 33 toppi 34 toppi 35 toppi 36 toppi 37;
        toppi 41 toppi 42 toppi 43 toppi 44 toppi 45 toppi 46 toppi 47;
        toppi 51 toppi 52 toppi 53 toppi 54 toppi 55 toppi 56 toppi 57;
        toppi 61 toppi 62 toppi 63 toppi 64 toppi 65 toppi 66 toppi 67;
        toppi 71 toppi 72 toppi 73 toppi 74 toppi 75 toppi 76 toppi 77]
```

```
top1=0;
for i=1:n
    top1=top1+(y(i)-epi(i)); end
```

```
top2=0;
for i=1:n
    top2=top2+(y(i)*(b2-x(i))^3- (b2-x(i))^3*epi(i)); end
```

```
top3=0;
for i=1:n
    top3=top3+(3*y(i)*b1*(b2-x(i))^2- 3*b1*(b2-x(i))^2*epi(i)); end
```

```
top4=0;
for i=1:n
    top4=top4+(y(i)*(x(i)-l amda1)^2*es1(i) - (x(i)-l amda1)^2*es1(i)*epi(i)); end
```

```

top5=0;
for i=1:n
    top5=top5+(y(i)*(x(i)-l amda2)*es2(i) - (x(i)-l amda2)*es2(i)*epi(i)); end

top6=0;
for i=1:n
    top6=top6+(-2*y(i)*b3*(x(i)-l amda1)*es1(i)+ 2*b3*(x(i)-l amda1)*es1(i)*epi(i)); end

top7=0;
for i=1:n
    top7=top7+(-y(i)*b4*es2(i)+ b4*es2(i)*epi(i)); end

```

F(1)=top1; F(2)=top2; F(3)=top3; F(4)=top4; F(5)=top5; F(6)=top6; F(7)=top7;

F

TersJacob=inv(Jacob)

teta=ilk teta+TersJacob\*F'

fark=abs(teta-ilk teta);

maxfark=max(fark);

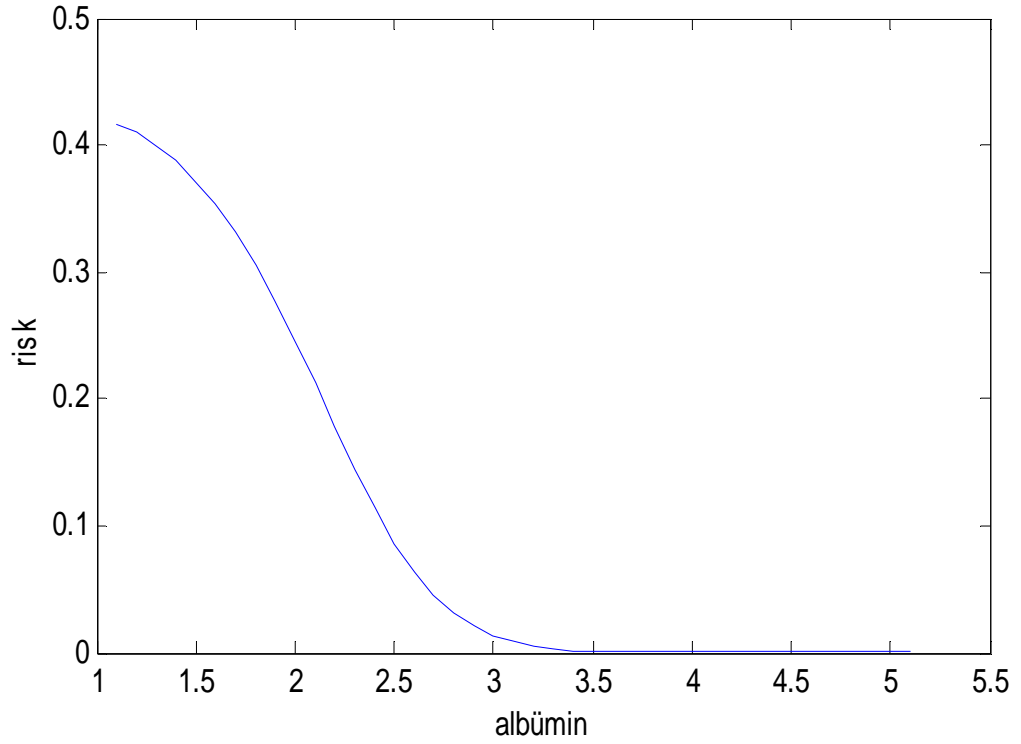
if maxfark<eps, break

end

ilk teta=teta;

b0=teta(1); b1=teta(2); b2=teta(3); b3=teta(4); b4=teta(5); l amda1=teta(6); l amda2=teta(7);

end



**Şekil IV.1** Serum Albüminin Risk Grafiği

Şekil IV.1'den de görüleceği gibi serum albümin değerinin 2.5g/100mL'den düşük olduğu seviyelerde ölüm oranı riski %40'lara kadar ulaşmaktadır. Orta derecede risk olarak kabul edilen 2.5g/100mL ve 3.5g/100mL aralığında ise %14'lere kadar risk görülmekte olup bu aralıktan sonrasında gittikçe azalan risk durumu söz konusudur. Hatta serum albümin için epidemiyolojik olarak ideal aralık 3.5g/100mL'den sonrası olarak kabul görmektedir. Serum albümin için model yapısının Matlab programında çözümlenmesiyle ilgili sonuçlar aşağıda gösterilmektedir:

$$n = 117$$

$$k1 = 2$$

*Jacob =*

*Columns 1 through 3*

-6.1219	26.6493	-13.5509
26.6493	-199.5820	301.2090
-13.5509	301.2090	-93.7363
-0.1253	1.9502	-0.6917
-0.0012	0.0432	-0.0119
0.3030	-3.8839	1.4757
0.0524	-1.4190	0.4240

*Columns 4 through 6*

-0.1253	-0.0012	0.3030
1.9502	0.0432	-3.8839
-0.6917	-0.0119	1.4757
-0.0507	-0.0023	-13.3077
-0.0023	-0.0004	0.0017
-13.3077	-0.0017	9.2965
0.0641	-2.9819	-0.0574

*Column 7*

0.0524  
-1.4190  
0.4240  
0.0641  
-2.9819  
-0.0574  
-0.2095

*F =*

*Columns 1 through 3*

14.0470 -183.3585 63.9666

*Columns 4 through 6*

5.9325 0.3988 -6.6888

*Column 7*

-11.9476

*TersJacob* =

*Columns 1 through 3*

-0.2205	-0.0049	0.0163
-0.0049	0.0012	0.0045
0.0163	0.0045	0.0016
-0.0014	0.0000	-0.0005
0.0007	-0.0000	-0.0017
0.0005	-0.0000	0.0004
-0.0000	0.0000	0.0001

*Columns 4 through 6*

-0.0014	0.0007	0.0005
0.0000	-0.0000	-0.0000
-0.0005	-0.0017	0.0004
-0.0524	0.0002	-0.0749
0.0002	0.0234	-0.0015
-0.0749	-0.0015	0.0003
0.0000	-0.3354	0.0001

*Column 7*

-0.0000  
0.0000  
0.0001  
0.0001  
-0.3354  
0.0001  
0.0000

*teta* =

-1.4760  
0.3035  
0.0903  
0.6323  
7.9319  
2.0905  
3.3680

Dikkat edilirse iterasyon sayısı fazla ilerlememektedir. Doğrusal olmayan sistem çözümlerinde Gauss-Newton yönteminde genelde görülen bir durumdur. Greg Ridgeway makalesinde lojistik regresyonun doğrusal olmaması durumunda parametre kestiriminde iterasyonun fazla gitmeyeceğini belirtmiştir. Bu nedenle doğrusal olmayan yapılarda bu durumun görülmesi muhtemeldir [ 21].

Serum albümin için kurulan modeldeki parametreler  $b_0=-1.4760$  ,  $b_1=0.3035$  ,  $b_2=0.0903$  ,  $b_3=0.6323$  ,  $b_4=7.9319$  ,  $\lambda_1=2.0905$  ,  $\lambda_2=3.3680$  olarak elde edilmiştir. Başlangıç değerlerine oldukça yakın olarak bulunan  $\lambda$  değerleri doğru bir modele yaklaştığımızın da bir kanıtı olarak gösterilebilir. Doğrusal olmayan bu lojistik regresyon modeli, en az iterasyonla en iyi modeli yakalayabilmiştir. Şekil IV.1'den de görüleceği gibi, tıbbi olarak serum albümin için geçerli değerler kurulan modelden de geçerliliğini korumaktadır.

## BÖLÜM V

### SONUÇLAR VE TARTIŞMA

Uygulama-1 çalışmasında model\_2 gerek en çok olabilirlik kestirimi gerekse de odds oran kestirimi bakımından daha iyi sonuç vermiştir. Yalnız burada söz konusu olan bir şey vardır ki o da parametre kestirimlerinin aralıklarıdır. Model\_2'nin parametre kestirimleri ve güven aralıkları hgbz1 hariç daha idealdir.

PROC NLIN gibi SAS paket programı kullanıldığında karşılaşılan problem ele alınan parametreler için her zaman alt ve üst sınırların elde edilmemesi şeklinde olup yakınsama ile ilgilidir. Hataların normal dağılıma sahip olduğu doğrusal olmayan regresyon analizlerinde Gauss-Newton, Marquaart ve Newton yöntemleri çok rahat elde edilmekte ancak doğrusal olmayan lojistik regresyon analizinde hataların binom dağılımına sahip olmasından dolayı sorunla karşılaşılmaktadır. Pek çok durumda istikrarlı çözümler sağlansa bile, doğrusal olmayan lojistik regresyon analizi kullanılarak yapılan modellemeler parametrelerin daha performanslı kestirimlerini sağlar [9].

Hemoglobin çalışmasında kabul edilebilirliği sınavıcı bir başka istatistik ise c (concordance ratio) değeridir. Bu değer 0 ile 1 arasında bir değer alır ve ROC eğrisinin altında kalan alana eşdeğer olduğundan yaptığımız uygulamada 0.955 ile oldukça iyi sonuç vermiştir.

Uygulama-2 Matlab programında yapılan serum albümin çalışması, splayn yöntemin lojistik regresyonda kullanılmasına dair iyi bir örnek olup parametre kestirimleri de anlamlıdır. Mevcut şekilden de desteklenerek epidemiyolojik olarak kabul edilebilir bir model kurulmuştur. Yapılan iterasyonla önemli değişim noktası olarak gösterilen 2.5g/100mL ve 3.5g/100mL değerlerine yaklaşık bir kestirim elde edilmiştir. Bu uygulama göstermiştir ki, serum albümin değerinin 2.5g/100mL'den düşük çıkması ölüm riskini %40 seviyelerine kadar arttırmaktadır. İdeal seviye olan

3.5g/100mL'ye yaklařtıka risk azalmaktadır. Sonu itibariyle doęrusal olmayan lojistik regresyonda model kestirimleri ve yorumların anlamlılıęı daha iyi sonu vermektedir ve programın geliřtirilmesi ile daha ayrıntılı ve geniř aplı bir kullanımın ortaya ıkacaęı dřünölmektedir.

## KAYNAKLAR

- [1] Aydın, D.: “*Semiparametrik Regresyon Modellemede Splayn Düzeltme Yaklaşımı İle Tahmin ve Çıkarsamalar*”, Doktora Tezi, Anadolu Üniversitesi Fen Bilimleri Enstitüsü İstatistik, Eskişehir, Türkiye, (2005).
- [2] Başarır, G.: “*Çok Değişkenli Verilerde Ayırsama Sorunu ve Lojistik Regresyon Analizi*”, Doktora Tezi, Hacettepe Üniversitesi Fen Bilimleri Enstitüsü Uygulamalı İstatistik, Ankara, Türkiye, (1990).
- [3] Bircan, H.: "Lojistik Regresyon Analizi ve Tıp Verileri Üzerine Bir Uygulama", *Kocaeli Üniversitesi Sosyal Bilimler Dergisi*, (2004).
- [4] Boucher, K. M., Slattey, M. L., Berry, T. D., et al.: “*Statistical Methods to Analyze Dose-Response and Trend Analysis in Epidemiologic Studies*”, *Journal of Clinical Epidemiology*, Vol 51, No 12, (1998).
- [5] Dobson, Annette J.: *An Introduction to Generalized Linear Models*, Chapman Hall, London, (1990).
- [6] Erdoğan Eygi, B.: “*Bankaların Mali Performanslarının Lojistik Regresyon ile Analizi ve İleriye Yönelik Tahmin*”, Doktora Tezi, Marmara Üniversitesi Sosyal Bilimler Enstitüsü, İstanbul, Türkiye, (2002).
- [7] Gardside, P. S.; Glueck, C. J.: *The Important role of Modifiable Dietary and Behaviour Characteristic in the Causation and Prevention of Coronary Heart Disease Hospitalization and Mortality*, *Journal of American College of Nutrition*, (1995).
- [8] Hosmer, David W.; Lemeshow, S.: *Applied Logistic Regression*, Second Edition, John Wiley, New York, (2000).
- [9] İyit, N.: “*Lineer Olmayan Lojistik Regresyon Analizinde Model Kurma Stratejileri ve Bir Uygulaması*”, Yüksek Lisans Tezi, Selçuk Üniversitesi Fen Bilimleri Enstitüsü İstatistik, Konya, Türkiye, (2003).

- [10] Kalaycı, Ş.: *SPSS Uygulamalı Çok Değişkenli İstatistik Teknikleri*, Asil Yayın, Ankara, Türkiye, (2006).
- [11] Kleinbaum, David G.: *Logistic Regression, A Self-Learning Text*, Springer-Verlag, New York, (1994).
- [12] Lesaffre, E.; Albert, A.: *A Multiple Group Logistic Regression Diagnostics*, Applied Statistics, (1989).
- [13] Liu, Q.; Lambe, M.; Baik, I.; et al.: “*A Prospective Study of the Transient Decrease in Ovarian Cancer Risk Following Childbirth*”, *Cancer Epidemiology Biomarkers & Prevention* , Vol. 15, 2508-2513, (2006).
- [14] McCullagh, P.; Nelder, J. A.: *Generalized Linear Models*, Chapman Hall, London, (1989).
- [15] Menard, S.: *Applied Logistic Regression Analysis*, Sage Publication, (1995).
- [16] Montgomery, Douglas C.; Peck, Elizabeth A.; Vining, G. Geoffrey: *Introduction to Linear Regression Analysis*, Third Edition, Wiley & Sons, New York.
- [17] Mulla, Z.: “*Spline Regression in Clinical Research*”, *The West Indian Medical Journal*, Vol 56 (1), (2007).
- [18] Özdamar, K.: *Paket Programlar ile İstatistiksel Veri Analizi*, Kaan Kitabevi, Eskişehir, Türkiye, (2004).
- [19] Pastor, R.; Guallar, E.: “*Use of Two-segmented Logistic Regression to Estimate Change-points in Epidemiologic Studies*”, *American Journal of Epidemiology* , Vol. 148, No. 7, 631-642, (1998).
- [20] Rice, John R.: *Numerical Methods, Software, and Analysis*, McGraw Hill Higher Education, (1983).
- [21] Ridgeway, G.: “*Discussion. Additive Logistic Regression: A Statistical View of Boosting*”, *The Annals of Statistics*, Vol 28, No 2, (2000).

- [22] Schimek, Michael G.: *Smoothing and Regression (Approaches, Computation and Application)*, John Wiley & Sons, New York, **(1999)**.
- [23] Tatlıdil, H.: *Uygulamalı Çok Değişkenli İstatistiksel Analiz*, H.Ü. Fen Fakültesi İstatistik Bölümü, Ankara, Türkiye, **(1992)**.
- [24] <http://www.hekimce.com/index.php?klm=hemoglobin> (04/2009)

## EK-1 Hemoglobin SAS Sonucu

x1	ri sk	l cl per	ucl per
5.5	100.000	0.0000	100.000
5.6	100.000	0.0000	100.000
5.7	100.000	0.0000	100.000
5.8	99.999	0.0000	100.000
5.9	99.999	0.0000	100.000
6.0	99.998	0.0000	100.000
6.1	99.996	0.0000	100.000
6.2	99.992	0.0000	100.000
6.3	99.986	0.0000	100.000
6.4	99.976	0.0000	100.000
6.5	99.959	0.0001	100.000
6.6	99.932	0.0005	100.000
6.7	99.888	0.0024	100.000
6.8	99.821	0.0102	100.000
6.9	99.721	0.0395	100.000
7.0	99.574	0.1391	100.000
7.1	99.365	0.4442	100.000
7.2	99.075	1.2812	100.000
7.3	98.684	3.3074	99.999
7.4	98.173	7.5307	99.997
7.5	97.525	14.8545	99.989
7.6	96.729	25.0592	99.962
7.7	95.780	36.2770	99.890
7.8	94.688	46.0559	99.732
7.9	93.474	52.8502	99.457
8.0	92.172	56.3793	99.076
8.1	90.829	57.2257	98.655
8.2	89.500	56.4495	98.247
8.3	88.243	55.2720	97.854
8.4	87.117	54.7698	97.420
8.5	86.174	55.6887	96.866
8.6	85.461	58.3762	96.100
8.7	84.768	61.4092	95.113
8.8	83.843	63.3881	93.959
8.9	82.656	64.3487	92.638
9.0	81.170	64.3220	91.155
9.1	79.340	63.3364	89.514
9.2	77.114	61.4307	87.697
9.3	74.437	58.6701	85.659
9.4	71.252	55.1530	83.320
9.5	67.509	51.0079	80.569
9.6	63.175	46.3830	77.284
9.7	58.248	41.4334	73.342
9.8	52.773	36.3095	68.654
9.9	46.851	31.1487	63.202
10.0	40.649	26.0729	57.081
10.1	34.388	21.1943	50.528
10.2	28.316	16.6284	43.894
10.3	22.674	12.5050	37.563
10.4	17.654	8.9585	31.838
10.5	13.373	6.0903	26.874
10.6	9.869	3.9260	22.683
10.7	7.105	2.4036	19.193

10.8	4.999	1.4016	16.303
10.9	3.443	0.7809	13.907
11.0	2.324	0.4168	11.917
11.1	1.540	0.2136	10.259
11.2	1.003	0.1052	8.872
11.3	0.642	0.0499	7.708
11.4	0.404	0.0228	6.727
11.5	0.250	0.0100	5.898
11.6	0.153	0.0043	5.195
11.7	0.092	0.0017	4.597
11.8	0.054	0.0007	4.086
11.9	0.032	0.0003	3.650
12.0	0.018	0.0001	3.275
12.1	0.008	0.0000	3.405
12.2	0.004	0.0000	3.877
12.3	0.002	0.0000	4.838
12.4	0.001	0.0000	6.598
12.5	0.000	0.0000	9.755
12.6	0.000	0.0000	15.396
12.7	0.000	0.0000	25.187
12.8	0.000	0.0000	40.604
12.9	0.000	0.0000	60.351
13.0	0.000	0.0000	78.781
13.1	0.000	0.0000	90.838
13.2	0.000	0.0000	96.662
13.3	0.000	0.0000	98.930
13.4	0.000	0.0000	99.691
13.5	0.000	0.0000	99.919
13.6	0.000	0.0000	99.980
13.7	0.000	0.0000	99.996
13.8	0.000	0.0000	99.999
13.9	0.000	0.0000	100.000
14.0	0.000	0.0000	100.000
14.1	0.000	0.0000	100.000

## EK-2 Albümin Matlab Sonucu

x	epi
1.1000	0.1433
1.2000	0.1312
1.4000	0.1036
1.6000	0.0744
1.7000	0.0605
1.7000	0.0605
1.8000	0.0478
1.8000	0.0478
1.9000	0.0364
1.9000	0.0364
1.9000	0.0364
2.0000	0.0269
2.0000	0.0269
2.0000	0.0269
2.1000	0.0191
2.1000	0.0191
2.2000	0.0131
2.2000	0.0131
2.3000	0.0086
2.3000	0.0086
2.3000	0.0086
2.3000	0.0086
2.3000	0.0086
2.4000	0.0054
2.4000	0.0054
2.4000	0.0054
2.4000	0.0054
2.5000	0.0036
2.5000	0.0036
2.5000	0.0036
2.5000	0.0036



3.1000 0.0001  
3.1000 0.0001  
3.1000 0.0001  
3.2000 0.0001  
3.2000 0.0001  
3.3000 0.0000  
3.3000 0.0000  
3.3000 0.0000  
3.3000 0.0000  
3.3000 0.0000  
3.3000 0.0000  
3.4000 0.0000  
3.4000 0.0000  
3.4000 0.0000  
3.4000 0.0000  
3.4000 0.0000  
3.5000 0.0000  
3.5000 0.0000  
3.5000 0.0000  
3.5000 0.0000  
3.5000 0.0000  
3.5000 0.0000  
3.5000 0.0000  
3.5000 0.0000  
3.5000 0.0000  
3.5000 0.0000  
3.6000 0.0000  
3.6000 0.0000  
3.6000 0.0000  
3.6000 0.0000  
3.6000 0.0000  
3.7000 0.0000  
3.7000 0.0000  
3.8000 0.0000  
3.8000 0.0000

3.8000	0.0000
3.9000	0.0000
3.9000	0.0000
4.0000	0.0000
4.1000	0.0000
4.1000	0.0000
4.1000	0.0000
4.2000	0.0000
4.2000	0.0000
4.3000	0.0000
4.4000	0.0000
4.4000	0.0000
4.4000	0.0000
4.4000	0.0000
4.7000	0.0000
4.7000	0.0000
4.8000	0.0000
5.1000	0.0000

## ÖZGEÇMİŞ

1984 yılı Üsküdar doğumlu Esra Zeynep ŞENSOY, 1990 yılında İstanbul'da başladığı eğitim hayatını ortaöğretim eğitimi bitimine kadar bu şehirde başarılı bir şekilde sürdürmüş olup, 2002 yılında yabancı dil ağırlıklı lisede okul ikinciliği elde etmiştir. Aynı sene Denizli'de Pamukkale Üniversitesi Fen-Edebiyat Fakültesi Matematik Bölümü'nde okumaya hak kazanarak bu üniversitede görmüş olduğu bir yıllık eğitim hayatından sonra bölüm birincisi olmuş ve Bursa Uludağ Üniversitesi Fen-Edebiyat Fakültesi Matematik Bölümü'ne yatay geçiş yapmıştır. Uludağ Üniversitesi'nde tamamladığı lisans eğitiminden sonra, Marmara Üniversitesi Fen Bilimleri Enstitüsü Uygulamalı Matematik alanında tezli yüksek lisans programına kayıt yaptırmıştır.