



بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

# Robust Regression

نام درس : مدل خطی

- ارائه دهندگان :

فاطمه افیونیان

مهرنوش شوشتری

# نیاز برای برآورد نیرومند

- در مدل رگرسیون خطی  $Y = X\beta + \varepsilon$  مشاهدات و خطاها دارای توزیع نرمال هستند در نتیجه می توان از روش حداقل مربعات برای برآورد  $\beta$  استفاده کرد. در حالیکه اگر مشاهدات، توزیع غیر نرمال داشته باشند به خصوص که شاخه انتهایی منحنی طویل تر یا ضخیم تر از نرمال باشد ممکن است روش حداقل مربعات مناسب نباشد.
- شاخه های ضخیم توزیع معمولاً نقاط دورافتاده ایجاد می کند که این نقاط دورافتاده بر برآورد حداقل مربعات تاثیر می گذارد. در تاثیر گذاری، دورافتاده ها برازش حداقل مربعات را تا حد زیادی به سمت خود می کشند و در نتیجه تعیین و تشخیص این دور افتاده ها مشکل می شود زیرا باقیمانده های مربوط به آن ها بطور ساختگی و مصنوعی کوچک هستند.
- روش رگرسیون نیرومند برای کاهش اثر مشاهداتی به کار می رود که اگر روش حداقل مربعات به کار گرفته شود، تاثیر گذاری بالایی خواهند داشت یعنی روش نیرومند باقیمانده های مرتبط با دورافتاده های بزرگ را کنار می گذارد که باعث می شود تشخیص دورافتاده ها خیلی ساده تر شود.
- یک رگرسیون نیرومند علاوه بر حساس نبودن نسبت به دورافتاده ها، زمانیکه مشاهدات توزیع نرمال دارد، دارای کارایی ۹۵-۹۰ درصد نسبت به روش برآورد حداقل مربعات می باشد.

# نقاط دور افتاده (Out Lier)

نقاطی هستند که  $|\hat{\varepsilon}_i|$  آن ها بزرگ باشد.

مشکلات به وجود آمده از داشتن نقاط دور افتاده:

باعث افزایش SSE می شود:

$$\hat{\sigma}^2 = \frac{SSE}{n-p} \rightarrow MSE = \frac{SSE}{n-p} \rightarrow \sum_{i=1}^n \hat{\varepsilon}_i^2$$

باعث افزایش MSE می شود:

$$T_j = \frac{\hat{\beta}_j - 0}{\sqrt{MSE C_{jj}}}$$

افزایش MSE باعث کم شدن مقدار  $T_j$  می شود:

$$H_0 \rightarrow t_{\alpha/2, n-p} < |T_j| \text{ If رد}$$

در نهایت شانس رد شدن فرض  $H_0$  کمتر می شود.

افزایش MSE باعث می شود که بعضی از متغیرهای توضیحی مهم به اشتباه حذف شوند، همچنین باعث افزایش طول فاصله اطمینان می شود در نتیجه دقت کم می شود، طول فاصله پیش بینی و فاصله اطمینان تک تک ضرایب رگرسیونی نیز افزایش می یابد.

## راه های تشخیص داده های دور افتاده:

1- باقیمانده های استاندارد نشده ( $\hat{\varepsilon}_i$ ):

$$\hat{\varepsilon}_i = y_i - \hat{y}_i$$

2- باقیمانده های استاندارد شده ( $d_i$ ):

$$d_i = \frac{\hat{\varepsilon}_i}{\sqrt{MSE}}$$

3- باقیمانده های استیودنت شده ( $r_i$ ):

$$r_i = \frac{\hat{\varepsilon}_i}{\sqrt{MSE \left(1 - \frac{1}{n} - \frac{(x_i - \bar{x})^2}{sxx}\right)}}$$

راه حل برخورد با اینگونه مشاهدات (داده های پرت):

1- حذف مرحله به مرحله داده های دور افتاده

2- استفاده از توزیع های دم کلفت تر برای برآورد ML پارامترها

3- استفاده از روش کمترین قدرمطلق خطا برای برآورد پارامترها



## نقاط اهرم گون (High Leverage)

اگر بردار مشاهدات متغیرهای توضیحی داده  $i$  ام دور از سایر مشاهدات باشد، داده  $i$  ام اهرم گون است.

روش تشخیص داده های اهرم گون:

$$h_{ii} > \frac{2p}{n} \quad ; \quad p=k+1, \quad n = \text{تعداد مشاهدات}$$

در این صورت مشاهده  $i$  ام اهرم گون است.

نقاط اهرمی خوب: نقاطی که از طرح کلی داده ها پیروی می کنند و درجهت  $X$  پرت هستند.  $|r_i| < 2$

نقاط اهرمی بد: نقاطی که از طرح کلی داده ها پیروی نمی کنند و درجهت  $Y$  پرت هستند.  $|r_i| > 2$

## داده های تاثیرگذار (Influential):

داده هایی هستند که خصوصیات کلیدی مدل ( $\hat{\beta}$  و  $\hat{\beta}_j$  و  $\sigma^2$  و  $\hat{y}_i$ ) را تحت تاثیر خود قرار می دهند.

راه های تشخیص داده های تاثیرگذار:

اگر خصوصیت کلیدی مدنظر ما  $\hat{\beta}$  باشد از آماره کوک برای تشخیص داده های تاثیرگذار استفاده می کنیم:

$$D_i = \frac{(\hat{\beta}_i - \hat{\beta})' x' x (\hat{\beta}_i - \hat{\beta})}{p \text{ MSE}}$$

مشاهده  $i$  ام تاثیرگذار : If  $D_i > 1$



برآوردگرهای  
رگرسیون نیرومند

برآوردگر M

برآوردگر LTS

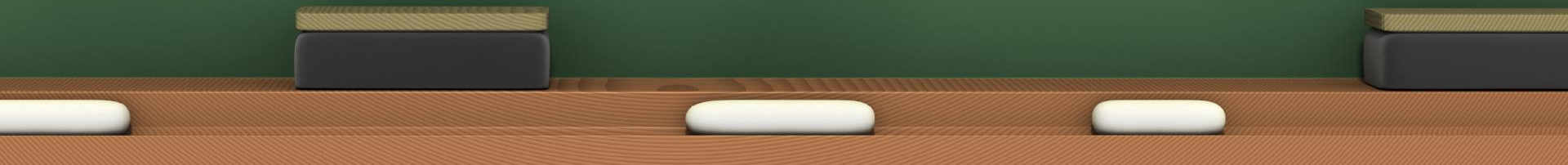
رگرسیون نرم-  
L1

برآوردگر LMS

برآوردگر GM

برآوردگر R

برآوردگر MM



## رگرسیون نرم-L1

مینیمم کردن مجموع قدرمطلق های خطاها، معمولاً مسئله رگرسیون نرم-L1 نامیده می شود. حداقل مربعات، مسئله رگرسیون نرم-L2 می باشد.

مسئله رگرسیون نرم-L1 حالت خاصی از رگرسیون نرم- $L_p$  می باشد که در آن پارامترهای مدل طوری انتخاب می شوند که  $\sum_{i=1}^n |\varepsilon_i|^p$  ( $1 \leq p \leq 2$ ) مینیمم شود.

$p=1.5$  یک انتخاب قراردادی است که منجر به برآورد هایی می شود که وقتی خطاها نرمال نیستند از برآورد حداقل مربعات بهتر است. وقتی هم که خطاها نرمال باشند به کارگیری  $p=1.5$  برآوردگرهایی را با کارایی به اندازه 90% حداقل مربعات می دهد. اما نمی توان این روش را به عنوان جایگزین برای حداقل مربعات در نظر گرفت زیرا به شدت تحت تاثیر مشاهده پرت درجهت  $X$  است.

## برآوردگر M

متداول ترین روش برآورد رگرسیون نیرومند، روش برآورد M است  
در حالت کلی می توان دسته ای از برآورد گرهای نیرومند را تعریف کرد که یک تابع  $\rho$  از  
باقی مانده ها را مینیمم می کند:

$$\min \sum_{i=1}^n P(e_i) = \min \sum_{i=1}^n P(y_i - x_i' \beta)$$

که در آن  $x_i$  نشان دهنده  $i$ -امین سطر  $X$  می باشد. اینگونه برآوردگر یک برآوردگر-M  
نامیده می شود که در آن  $M$  به عنوان ماکزیمم درستنمایی است یعنی تابع  $\rho$  به تابع  
درستنمایی (برای یک انتخاب مناسب توزیع خطا) بستگی دارد.

برای مثال اگر روش حداقل مربعات به کار برده شود نتیجه گرفته می شود که توزیع  
خطا نرمال است. در این صورت:

$$P(z) = \frac{1}{2} z^2 \quad -\infty < z < \infty$$

تابع  $p$  باید ویژگی های زیر را داشته باشد:

$$p(0) = 0$$

مقارن است:  $p(e) = p(-e)$

همواره نامنفی است:  $p(e) \geq 0$

یکنوا است:  $\text{if } e_1 \geq e_2 \quad \text{Then } p(e_1) \geq p(e_2)$



برآوردگر M-لزو ما پایا مقیاس نیست (یعنی اگر باقیمانده های  $y_i - x_i' \beta$  در یک ثابت ضرب شوند جواب جدید ممکن است مانند جواب قدیم نباشد). برای بدست آوردن شکل جدید پایا-مقیاس از این برآوردگر معمولا معادله زیر را حل می کنیم:

$$\min \sum_{i=1}^n P\left(\frac{e_i}{s}\right) = \min \sum_{i=1}^n P\left(\frac{y_i - x_i' \beta}{s}\right)$$

که در آن  $s$  یک برآوردگر نیرومند مقیاس است. یک جواب عمومی برای  $s$  عبارتست از:

$$s = \frac{\text{median}|e_i - \text{median}(e_i)|}{0.675}$$

اگر  $n$  بزرگ باشد و توزیع خطا نرمال باشد ثابت  $0.6745$ ،  $S$  را تقریباً به یک برآوردگر نارایب تبدیل می کند.

برای مینیمم کردن مشتقات جزئی مرتبه اول  $P$  نسبت به  $\beta_j$  ( $j=0,1,\dots,k$ ) را مساوی صفر قرار می دهیم بدین ترتیب یک شرط لازم برای بدست آوردن مینیمم حاصل می شود . با این کار دستگاهی از  $P = k+1$  معادله تشکیل می گردد که عبارتند از :

$$\sum_{i=1}^n x_{ij} \psi \left( \frac{e_i}{s} \right) = 0$$

که در آن  $\psi = \rho'$  و  $x_{ij}$  عبارتست از  $i$  امین مشاهده از  $i$  امین متغیر رگرسیونی و  $x_{i0} = 1$  می باشد. در حالت کلی تابع  $\psi$  غیر خطی و بایستی با روش تکرار حل شود. اگر تکنیک های متعدد غیر خطی بتواند بکار گرفته شود، حداقل مربعات موزون بطور وسیعی مورد استفاده قرار می گیرد. این روش معمولاً به بیتون و توکی [۱۹۷۴] نسبت داده می شود. برای اینکه بطور مکرر حداقل مربعات دوباره موزون شده را بکار ببریم فرض کنیم که یک برآورد اولیه  $\widehat{\beta}_0$  در دسترس باشد و  $S$  یک بردار مقیاس است. در این صورت  $P = k+1$  معادله مذکور بصورت زیر نوشته میشود:

$$\sum_{i=1}^n x_{ij} \psi\left(\frac{e_i}{s}\right) =$$

$$\sum_{i=1}^n \frac{x_{ij} \psi\{[(y_i - x_i' \beta)/s]/(y_i - x_i' \beta)\} \{(y_i - x_i' \beta)\}}{s} = 0$$

به طوری که:

$$\sum_{i=1}^n x_{ij} w_{i0} (y_i - x_i' \beta) = 0$$

$$w_{i0} = \begin{cases} \frac{\psi[(y_i - x_i' \widehat{\beta}_0)/s]}{(y_i - x_i' \widehat{\beta}_0)/s} & y_i \neq x_i' \widehat{\beta}_0 \\ 1 & y_i = x_i' \widehat{\beta}_0 \end{cases} \quad \text{که در آن:}$$

با نماد ماتریسی بصورت زیر نوشته می شود:

$$x'w_0x\beta = x'w_0y$$

که در آن  $w_0$  یک ماتریس قطری  $n \times n$  از "وزن ها" با اعضای قطر  $w_{n0}, \dots, w_{10}, w_{20}$  که در صفحه قبل مشخص گردیده اند. معادله بالا را به عنوان معادلات حداقل مربعات نرمال معمولی می شناسیم. در نتیجه برآورد قدم اول عبارتست از:

$$\widehat{\beta}_1 = (x'w_0x)^{-1}x'w_0y$$

در قدم بعدی وزن ها را مجددا محاسبه می کنیم اما  $\widehat{\beta}_1$  را به جای  $\widehat{\beta}_0$  به کار می بریم. در حالت کلی داریم:

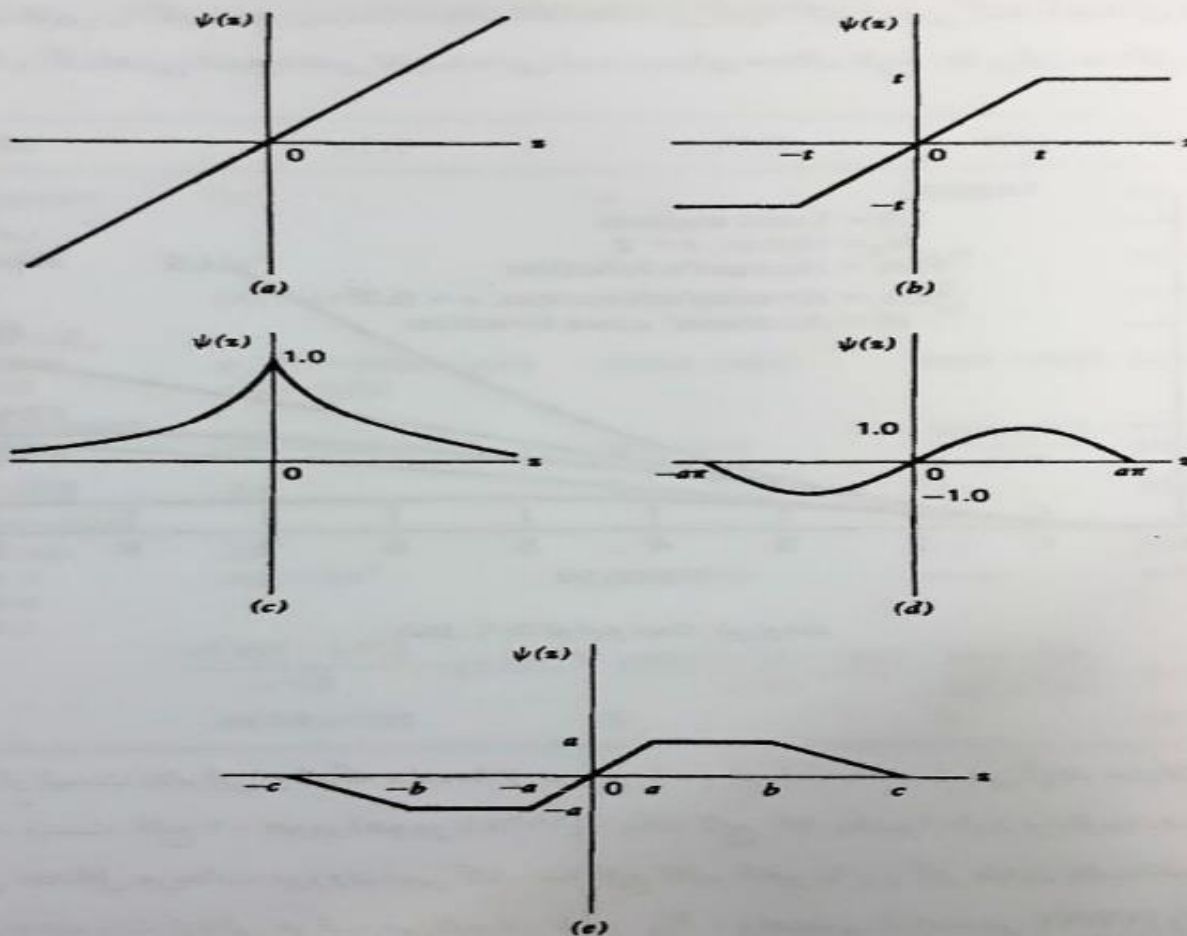
$$\widehat{\beta}_{j+1} = (x'w_jx)^{-1}x'w_jy$$



# در این شکل تعدادی توابع محک نیرومند متداول نشان داده شده است:

Criterion	$\rho(z)$	$\psi(z)$	$w(z)$	Range
Least squares	$\frac{1}{2}z^2$	$z$	1.0	$ z  < \infty$
Huber's $t$ function $t=2$	$\frac{1}{2}z^2$ $ z t - \frac{1}{2}t^2$	$z$ $t \text{ sign}(z)$	1.0 $\frac{t}{ z }$	$ z  \leq t$ $ z  > t$
Ramsay's $E_a$ function $a=0.3$	$a^{-2}[1 - \exp(-a z ) \cdot (1 + a z )]$	$z \exp(-a z )$	$\exp(-a z )$	$ z  < \infty$
Andrew's wave function $a=1.339$	$a[1 - \cos(z/a)]$ $2a$	$\sin(z/a)$ 0	$\frac{\sin(z/a)}{z/a}$ 0	$ z  \leq a\pi$ $ z  > a\pi$
Hampel's 17A function $a=1.7$ $b=3.4$ $c=8.5$	$\frac{1}{2}z^2$ $a z  - \frac{1}{2}a^2$ $\frac{a(c z  - \frac{1}{2}z^2)}{c-b} - (7/6)a^2$ $a(b+c-a)$	$z$ $a \text{ sign}(z)$ $a \text{ sign}(z)(c- z )$ 0	1.0 $a/ z $ $\frac{a(c- z )}{ z (c-b)}$ 0	$ z  \leq a$ $a <  z  \leq b$ $b <  z  \leq c$ $ z  > c$

در این شکل نیز نمودار توابع محک نیرومند ارائه شده است:



شکل ۶-۹ توابع تاثیر نیرومند: (a) حداقل مربعات، (b) تابع  $t$ -هوبرت، (c) تابع  $E_a$ -رامسی، (d) تابع موج اندرو و (e) تابع 17A هامیل

## برآوردگر R

روش برآورد R بر اساس رتبه ها است. در این روش جانشین کردن یک عامل رتبه در تابع حداقل مربعات  $S(\beta) = \sum_{i=1}^n (y_i - x_i' \beta)^2$  به وسیله رتبه هایش را مورد توجه قرار می دهیم. اگر  $R_i$  رتبه  $y_i - x_i' \beta$  باشد می خواهیم

$\sum_{i=1}^n (y_i - x_i' \beta) R_i$  مینیمم شود. رتبه ها را که اعداد صحیح هستند با تابع زیان  $a(i)$  ;  $(i=1,2,\dots,n)$  جانشین می کنیم بنابراین تابع هدف برابر می شود با:

$$\min_{\beta} \sum_{i=1}^n (y_i - x_i' \beta) a(R_i)$$

برآوردگر R از نظر محاسباتی از برآوردگر M مشکل تر است اما در شرایط واقعی به طور مجانبی با برآوردگر M معادل است.

## برآوردگر LTS

برآوردگر کمترین مربعات پیراسته:

$$\hat{\beta} = \sum_{i=1}^n e_{(i)}^2$$

که در آن:

$$e_{(1)}^2 \leq e_{(2)}^2 \leq \dots \leq e_{(n)}^2$$

در این رابطه بزرگترین عدد صحیح کوچک تر یا مساوی  $\frac{n}{2} + 1$  است. اگر خطاها دارای توزیع نرمال باشند این برآوردگر از برآوردگر حداقل مربعات ناکارتر است.

## برآوردگر LMS

در این روش میانه مربعات باقی مانده ها را بدست می آوریم سپس با مینیمم کردن آن برآورد ضرایب را بدست می آوریم:

$$\hat{\beta} = \min \text{med} (y_i - x_i \beta)^2$$



## برآوردگر GM

برآورد GM در سال ۱۹۷۵ توسط مالوس معرفی شد. اساس این روش مبتنی بر به کارگیری توابع وزنی است و برآورد از حل معادله زیر بدست می آید:

$$\sum_{i=1}^n w_i \psi\left(\frac{r_i(\hat{\beta})}{\hat{\sigma}}\right) x_i = 0$$

اغلب برای سادگی  $\sigma$  را معلوم فرض کرده و یا برآورد آن را در معادله قرار می دهیم.

## برآوردگر MM

این برآوردگر از حل معادله زیر به دست می آید:

$$L(\beta) = \sum_{i=1}^n \rho_i \left( \frac{r_i(\hat{\beta})}{\hat{\sigma}_n} \right)$$

که در آن:

$$\rho < \rho_0 \rightarrow L(\hat{\beta}) \leq L(\beta_0)$$

# یک مثال واقعی

مقایسه نتایج مدل رگرسیون معمولی و رگرسیون نیرومند در مدل بندی عوامل مرتبط با بیماری پره دیابت:

1. گروه بهداشت عمومی، دانشگاه علوم پزشکی نیشابور
2. گروه آموزشی آمار، دانشکده علوم ریاضی، دانشگاه فردوسی مشهد
3. گروه آمار زیستی، مرکز تحقیقات عوامل اجتماعی موثر بر سلامت، دانشگاه علوم پزشکی مشهد
4. گروه اپیدمیولوژی، مرکز تحقیقات عوامل اجتماعی موثر بر سلامت، دانشگاه علوم پزشکی مشهد
5. گروه جمعیت شناسی، کمیته تحقیقات دانشجویی، دانشگاه آزاد اسلامی واحد تهران





**زمینه و هدف:** با توجه به اینکه خطر ابتلا به دیابت در افراد پره‌دیابتیک بسیار بالا است، تعیین عوامل موثر بر پره‌دیابت دارای اهمیت می‌باشد. این مطالعه با هدف مقایسه نتایج مدل رگرسیون معمولی و رگرسیون نیرومند در مدل‌بندی عوامل مرتبط با بیماری پره‌دیابت انجام شد.

**روش بررسی:** این مطالعه که از نوع مقطعی-تحلیلی است روی ۶۴۶۰ نفر از افراد بالای ۳۰ سال، شرکت‌کننده در طرح غربالگری دیابت دانشگاه علوم پزشکی مشهد، از مهر تا آذر ۱۳۸۹ انجام شد. با توجه به میزان قندخون ناشتای افراد، ۵۴۱۴ نفر سالم و ۱۰۴۶ نفر به‌عنوان پره‌دیابتیک شناسایی شدند. سن، جنس، نمایه توده بدن، فشارخون سیستولیک، فشارخون دیاستولیک و نسبت کمر به باسن در مورد هر فرد اندازه‌گیری شد. مدل رگرسیون معمولی روی داده‌ها برازش شد. سپس داده‌های پرت مشخص و سه مدل نیرومند Mallow، WBY و BY برازش شد. آنگاه مدل‌ها با هم مقایسه گردیدند.

**یافته‌ها:** متغیرهای سن، نمایه توده بدن و فشارخون سیستولیک در همه مدل‌ها از لحاظ آماری معنادار شدند ( $P < 0.01$ ) و متغیر نسبت کمر به باسن معنادار نشد ( $P > 0.1$ ) تعداد ۵۵۲ داده‌ی پرت با خطای بد رده‌بندی در مدل معمولی وجود داشت. مقادیر کای دو پیرسون و سطح زیرمنحنی را که در مدل Mallow به‌طور تقریبی فرقی با مدل معمولی نداشت. اما در مدل‌های WBY و BY به نسبت بیشتر بود.

**نتیجه‌گیری:** با توجه به نتایج این پژوهش، سن بالا، نمایه توده بدنی و فشارخون بالا در ابتلا به بیماری پره‌دیابت موثر می‌باشند. همچنین مدل‌های رگرسیون نیرومند WBY و BY برازش بهتر و توان پیشگویی بالاتری نسبت به رگرسیون معمولی در مدل‌بندی عوامل گفته‌شده در ارتباط با پره‌دیابت دارند.

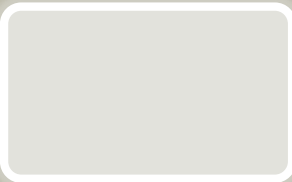
# منابع



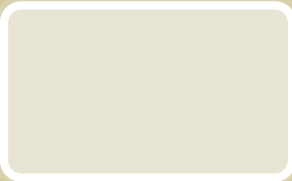
مقدمه ای بر تحلیل رگرسیون خطی  
(داگلاس مونتهگمری و الیزابت پک)



Linear Regression Analysis  
(Seber, 2003)



رگرسیون استوار (نیرومند)  
(سید مهدی امیرجهانشاهی و حسینعلی نیرومند)



Introduction to Robust Estimation Techniques  
(Andreas Ruckstuhl, 2016)

بہ تشکر از  
مہرمانے و توجہ تان

