

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



سمینار درس مدل های فنی

stepwise method

ارائه دهندگان :

مینا عزیزی

سمیرا شاه نوشی

دی ماه ۹۷

* مسئله انتخاب متغیر

* اهمیت انتخاب متغیرهای مناسب

* روش های انتخاب متغیر

* (انتخاب پیشرو) Forward selection

* (حذف پسرو) Backward elimination

* (رگرسیون قدم به قدم) Stepwise method

* معایب روش قدم به قدم

* منابع ومراجع

* مسئله انتخاب متغیر (variable selection)

* پیدا کردن زیرمجموعه ای از متغیرهای رگرسیونی مناسب برای مدل، مسئله **انتخاب متغیر** نامیده می شود.

A **variable selection** method is a way of selecting a particular set of independent variable for use in a regression model.

* اهمیت انتخاب متغیرهای مناسب

انتخاب متغیرهایی با بیشترین تاثیر بر روی متغیر پاسخ

* مدل واقعی معمولا بسیار پیچیده، مجهول و ناشناخته است.

* همواره مدل های با تعداد متغیر کم دلخواه تحلیل گران آماری بوده است. زیرا:

* ۱- روابط بین متغیرهای علمی را به صورت ساده بیان کرده و قابلیت تفسیر ساده تری دارند.

* ۲- هزینه های محاسباتی و صرف زمان، اهمیت انتخاب متغیرهای مناسب را بیشتر می کند.

*۳- چنانچه متغیرهای توضیحی نامناسب وارد مدل شوند باعث افزایش MSE می شود و در نتیجه:

* ← کاهش دقت برآورد ضرایب رگرسیونی ودقت پیش بینی

* ← افزایش طول فواصل اطمینان وفواصل پیش بینی

* ← کاهش دقت آزمون ها (ممکن است متغیر توضیحی مناسبی به اشتباه حذف شود.)

$Y_{(n \times 1)}$
متغیر پاسخ

X_1, \dots, X_k
متغیرهای توضیحی

$$Y = X\beta + \varepsilon$$



فرض می کنیم k متغیر توضیحی و متغیر پاسخ را در اختیار داریم و می خواهیم مدل مناسبی برازش دهیم.

یکی از راه های انتخاب مدل این است که تمام 2^p مدل ممکن را برازش دهیم (همه رگرسیون های ممکن) سپس با استفاده از معیارهای مناسبت مدل بهترین آن ها را انتخاب کنیم. اما این کار در حالتی که p **بزرگ** باشد بسیار مشکل است. پس به دنبال روشی برای یافتن **مدل مناسب** هستیم.

از جمله افرادی که در این زمینه تلاش کرده اند می توان کاکس و سنل (۱۹۷۴)، هاکینگ (۱۹۷۲)، میرز (۱۹۹۰)، هاکینگ و لاموت (۱۹۷۳) و تومپسون (۱۹۷۸) را نام برد.

روشهای انتخاب متغیر*

Tests-based

Penalty-based

Screening-based

*

*Tests-based

*

1) Stepwise(forward & backward)

2)Autometrics

*Forward selection

* این روش با فرض اینکه هیچ متغیر رگرسیونی غیر از عرض از مبدا در مدل نیست، آغاز می شود. سپس با هر یک از k متغیر توضیحی و عرض از مبدا به صورت جداگانه مدل رگرسیونی برازش می دهیم و برای هر کدام آزمون صفر بودن ضریب رگرسیونی مربوط به آن را انجام می دهیم. در واقع داریم:

$$* Y_i = \beta_0 + \beta_1 X_{i1} + \epsilon_i$$

$$* Y_i = \beta_0 + \beta_2 X_{i2} + \epsilon_i$$

*

.

*

.

*

.

$$* Y_i = \beta_0 + \beta_k X_{ik} + \epsilon_i$$

گام اول

*

* $Y_i = \beta_0 + \beta_1 x_{i1} + \epsilon_i$

$H_0 : \beta_1 = 0$

$H_1 : \beta_1 \neq 0$

$$F = \frac{(RSS_r - RSS_f)/1}{RSS_f/(n-2)}$$

* $Y_i = \beta_0 + \beta_2 x_{i2} + \epsilon_i$

$H_0 : \beta_2 = 0$

$H_1 : \beta_2 \neq 0$

$$F = \frac{(RSS_r - RSS_f)/1}{RSS_f/(n-2)}$$

*

*

*

*

* $Y_i = \beta_0 + \beta_k x_{ik} + \epsilon_i$

$H_0 : \beta_k = 0$

$H_1 : \beta_k \neq 0$

$$F = \frac{(RSS_r - RSS_f)/1}{RSS_f/(n-2)}$$

$$F_{IN} = F_{\alpha}(1, n-2) < F_{\text{بزرگترین آماره}}$$

yes

no

$$Y_i = \beta_0 + \beta_j x_{ij} + \epsilon_i$$

stop

* با فرض اینکه بزرگترین آماره F مربوط به آزمون صفر بودن ضریب زامین متغیر توضیحی بوده است.

* در روش انتخاب پیشرو $\alpha = 0.25$ توصیه می شود.



*
$$H_0 : \beta_1 = 0$$

*
$$Y_i = \beta_0 + \beta_j x_{ij} + \beta_1 x_{i1} + \epsilon_i$$

*
$$H_1 : \beta_1 \neq 0$$

$$F = \frac{(RSS_r - RSS_f)/1}{RSS_f/(n-3)}$$

*
$$H_0 : \beta_2 = 0$$

*
$$Y_i = \beta_0 + \beta_j x_{ij} + \beta_2 x_{i2} + \epsilon_i$$

*
$$H_1 : \beta_2 \neq 0$$

$$F = \frac{(RSS_r - RSS_f)/1}{RSS_f/(n-3)}$$

*

*

*

*
$$H_0 : \beta_{k-1} = 0$$

*
$$Y_i = \beta_0 + \beta_j x_{ij} + \beta_{k-1} x_{i(k-1)} + \epsilon_i$$

*
$$H_1 : \beta_{k-1} \neq 0$$

$$F = \frac{(RSS_r - RSS_f)/1}{RSS_f/(n-3)}$$

$$F_{IN} = F_{\alpha}(1, n-3) <$$

بزرگترین آماره F

yes

no

$$Y_i = \beta_0 + \beta_j x_{ij} + \beta_l x_{il} + \epsilon_i$$

stop

* با فرض اینکه بزرگترین آماره F مربوط به آزمون صفر بودن ضریب امین متغیر توضیحی بوده است.

در این روش در هر گام متغیری که آماره آزمون مربوط به صفر بودن ضریب آن از بقیه **بزرگتر** است (دارای بیشترین همبستگی ساده با متغیر پاسخ) برای **ورود** به مدل کاندید می شود و این روند تا زمانی ادامه می یابد که بزرگترین آماره F به دست آمده از F_{IN} تجاوز نکند و یا اینکه همه متغیرهای توضیحی به مدل وارد شده باشند و مدل به دست آمده در گام قبل به عنوان مدل نهایی گزارش می شود.

*

۱- این روش برای مسائلی که در ابتدایازبه جمع آوری اطلاعات دارند مناسب است.

۲- انجام این روش ساده و قابل فهم است.

مزایای انتخاب
پیشرو

متغیری که در این روش وارد مدل شده، قابل حذف نیست.

معایب انتخاب
پیشرو

* Backward elimination

* در این روش ابتدا مدلی شامل کلیه k متغیر رگرسیونی نامزد برآزش داده می شود، سپس آزمون صفر بودن همه ضرایب رگرسیونی را به صورت مجزا انجام می دهیم.

$$* Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \epsilon_i$$

$$* Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \epsilon_i$$

*

$$H_0 : \beta_1 = 0$$

*

$$F = \frac{(RSS_r - RSS_f) / 1}{RSS_f / (n - p)}$$

*

$$H_1 : \beta_1 \neq 0$$

*

.

.

*

.

.

*

.

.

*

$$H_0 : \beta_k = 0$$

*

$$F = \frac{(RSS_r - RSS_f) / 1}{RSS_f / (n - p)}$$

*

$$H_1 : \beta_k \neq 0$$

$$F_{\text{out}} = F_{\alpha}(1, n-p) >$$

کوچکترین آماره F

yes

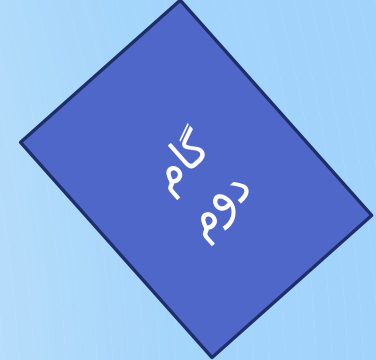
no

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{j-1} x_{i(j-1)} + \beta_{j+1} x_{i(j+1)} + \dots + \beta_k x_{ik} + \epsilon_i$$

stop

* با فرض اینکه کوچکترین آماره F مربوط به آزمون صفر بودن ضریب زامین متغیر توضیحی بوده است.

* در روش حذف پسرو $\alpha = 0.1$ توصیه می شود.



$$* Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{j-1} x_{i(j-1)} + \beta_{j+1} x_{i(j+1)} + \dots + \beta_k x_{ik} + \epsilon_i$$

$$* \left\{ \begin{array}{l} H_0 : \beta_1 = 0 \end{array} \right.$$

*

$$F = \frac{(RSS_r - RSS_f) / 1}{RSS_f / (n - (p - 1))}$$

$$* \left\{ \begin{array}{l} H_1 : \beta_1 \neq 0 \end{array} \right.$$

*

.

*

.

*

.

$$* \left\{ \begin{array}{l} H_0 : \beta_k = 0 \end{array} \right.$$

*

$$F = \frac{(RSS_r - RSS_f) / 1}{RSS_f / (n - (p - 1))}$$

$$* \left\{ \begin{array}{l} H_1 : \beta_k \neq 0 \end{array} \right.$$

*

$$F_{\text{out}} = F_{\alpha}(1, n - (p - 1)) > F_{\text{آماره کوچکترین}}$$

yes

no

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{j-1} x_{i(j-1)} + \beta_{j+1} x_{i(j+1)} + \dots + \beta_{l-1} x_{i(l-1)} + \beta_{l+1} x_{i(l+1)} + \dots + \beta_k x_{ik} + \epsilon_i$$

stop

* با فرض اینکه کوچکترین آماره F مربوط به آزمون صفر بودن ضریب امین متغیر توضیحی بوده است.

* در این روش در هر گام متغیری که آماره آزمون مربوط به صفر بودن ضریب β_n از بقیه کوچکتر است برای حذف از مدل کاندید می شود و این روند تا زمانی ادامه می یابد که کوچکترین آماره F به دست آمده از F_{out} کوچک تر نباشد و یا اینکه همه متغیرهای توضیحی از مدل حذف شده باشند و مدل به دست آمده در گام قبل به عنوان مدل نهایی گزارش می شود.

۱- این روش برای مسائلی که در آغاز فرضی بنا شده و سپس در جهت اثبات آن باشیم مناسب است.

۲- انجام این روش ساده و قابل فهم است.

مزایای حذف
پسرو

متغیری که در این روش از مدل حذف شده، قابل برگشت نیست.

معایب حذف
پسرو

*Stepwise method

* به دلیل وجود معایب ذکر شده در روش انتخاب پیشرو و حذف پسرو ، روش جدیدی به نام **stepwise** ارائه می شود که در واقع ترکیبی از دو روش مذکور می باشد.

* این روش توسط **Efroymsen** در سال ۱۹۶۰ ارائه شد.

* روش قدم به قدم با انتخاب پیشرو شروع شده و با حذف پسرو ادامه می یابد و در واقع رگرسیون قدم به قدم تعدیل انتخاب پیشرو می باشد.

* در رگرسیون قدم به قدم در هر مرحله بعد از وارد شدن متغیر جدید به مدل بررسی می کنیم که آیا **در حضور بقیه متغیرها** به این متغیر احتیاجی هست یا خیر. چرا که ممکن است متغیر رگرسیونی اضافه شده در مرحله قبل به لحاظ ارتباط با متغیرهای رگرسیونی که اکنون در معادله اند، زائد باشد.

* در واقع در این روش با مدل شامل عرض از مبدا شروع کرده و روشی همانند الگوریتم ارائه شده در صفحه بعد را دنبال می کنیم تا به مدل نهایی دست یابیم.

گام اول

* $Y_i = \beta_0 + \beta_1 x_{i1} + \epsilon_i$

$H_0 : \beta_1 = 0$

$H_1 : \beta_1 \neq 0$

* $Y_i = \beta_0 + \beta_2 x_{i2} + \epsilon_i$

$H_0 : \beta_2 = 0$

$H_1 : \beta_2 \neq 0$

*
.

*
.

*
.

*
.

* $Y_i = \beta_0 + \beta_k x_{ik} + \epsilon_i$

$H_0 : \beta_k = 0$

$H_1 : \beta_k \neq 0$

*

$$F_{IN} = F_{\alpha}(1, n-2) <$$

بزرگترین آماره F

yes

no

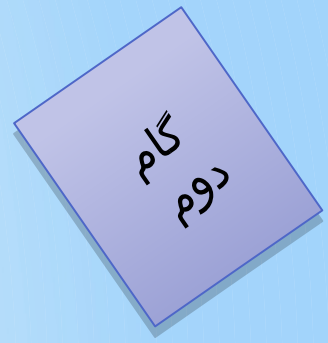
$$Y_i = \beta_0 + \beta_j x_{ij} + \epsilon_i$$

stop

* با فرض اینکه بزرگترین آماره F مربوط به آزمون صفر بودن ضریب زامین متغیر توضیحی بوده است.

$$\begin{array}{l}
 * \\
 * \\
 *
 \end{array}
 \left\{ \begin{array}{l}
 Y_i = \beta_0 + \beta_j x_{ij} + \beta_1 x_{i1} + \epsilon_i \\
 H_0 : \beta_1 = 0 \\
 H_1 : \beta_1 \neq 0
 \end{array} \right.$$

$$F = \frac{(RSS_r - RSS_f)/1}{RSS_f/(n-3)}$$



$$\begin{array}{l}
 * \\
 * \\
 * \\
 *
 \end{array}
 \left\{ \begin{array}{l}
 Y_i = \beta_0 + \beta_j x_{ij} + \beta_2 x_{i2} + \epsilon_i \\
 H_0 : \beta_2 = 0 \\
 H_1 : \beta_2 \neq 0
 \end{array} \right.$$

$$F = \frac{(RSS_r - RSS_f)/1}{RSS_f/(n-3)}$$

$$\begin{array}{l}
 * \\
 * \\
 * \\
 *
 \end{array}
 \left\{ \begin{array}{l}
 \cdot \\
 \cdot \\
 \cdot \\
 \cdot
 \end{array} \right.$$

$$\begin{array}{l}
 * \\
 * \\
 *
 \end{array}
 \left\{ \begin{array}{l}
 Y_i = \beta_0 + \beta_j x_{ij} + \beta_{k-1} x_{i(k-1)} + \epsilon_i \\
 H_0 : \beta_{k-1} = 0 \\
 H_1 : \beta_{k-1} \neq 0
 \end{array} \right.$$

$$F = \frac{(RSS_r - RSS_f)/1}{RSS_f/(n-3)}$$

$$F_{\text{max}} > F_{\text{IN}} = F_{\alpha}(1, n-3)$$

$$Y_i = \beta_0 + \beta_j x_{ij} + \beta_l x_{il} + \epsilon_i$$

stop

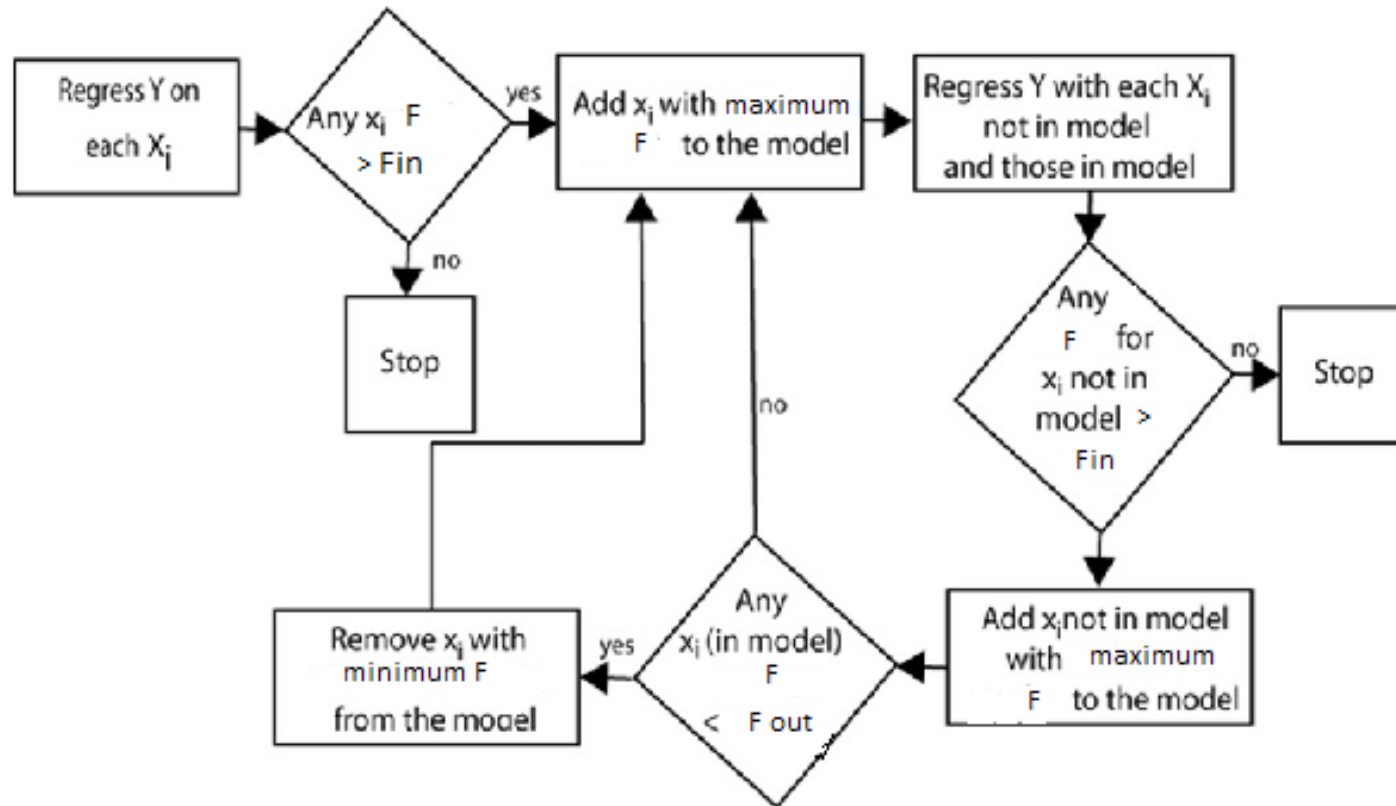
$$F_{\text{min}} < F_{\text{out}} = F_{\alpha}(1, n-3)$$

$$Y_i = \beta_0 + \beta_j x_{ij} + \epsilon_i$$

گام بعد

*

* الگوریتم روش قدم به قدم :



* انتخاب مقادیر برای F_{IN} و F_{out} در حد زیادی یک موضوع ترجیحی تحلیلگراست . می توان مقادیر F_{IN} و F_{out} را طوری انتخاب کرد که همه متغیرهای رگرسیونی به وسیله انتخاب پیشرو وارد ویا به وسیله حذف پسرو کنار گذاشته شوند.

* بعضی تحلیلگران ترجیح می دهند در رگرسیون قدم به قدم

* $F_{IN} = F_{out} = 4$ انتخاب کنند

* به طور معمول $F_{IN} < F_{out}$ انتخاب می کنیم ،چرا که به طور نسبی اضافه کردن یک متغیر رگرسیونی نسبت به حذف آن مدل را مشکل تر می کند.

* مثال:

* داده های مربوط به یک نمونه ۵۰ تایی از متغیر پاسخ y و ۶ متغیر توضیحی X_1, \dots, X_6 را وارد نرم افزار SPSS کرده و به صورت زیر به روش step wise مدل مناسبی می یابیم.

$$n=50$$

$$k=6$$

$$p=7$$

The screenshot shows the IBM SPSS Statistics software interface. The menu bar includes View, Data, Transform, Analyze, Direct Marketing, Graphs, Utilities, Add-ons, Window, and Help. The Analyze menu is open, showing options like Reports, Descriptive Statistics, Tables, Compare Means, General Linear Model, Generalized Linear Models, Mixed Models, Correlate, Regression, Loglinear, Neural Networks, Classify, and Dimension Reduction. The Regression option is highlighted, and a sub-menu is open showing options like Automatic Linear Modeling..., Linear..., Curve Estimation..., Partial Least Squares..., and Binary Logistic... The data table in the background has columns y, x1, x4, x5, x6, and var. The value 55.00 is visible in the top right of the data area.

y	x1	x4	x5	x6	var
20.00	1.0	45.00	45.00	60.00	
70.00	2.0	88.00	3.00	3.00	
23.00	3.0	34.00	56.00	3.00	
23.00	45.0	56.00	6.00	3.00	
12.00	3.0				
12.00	34.0				
12.00	345.0				
23.00	54.0				
23.00	56.0				

Linear Regression

- x1
- x2
- x3
- x4
- x5
- x6

Dependent:

y

Block 1 of 1

Previous

Next

Independent(s):

- x1
- x2
- x3

Method:

- Enter
- Stepwise
- Remove
- Backward
- Forward

Selection Variable:

Case Labels:

WLS Weight:

Statistics...

Plots...

Save...

Options...

Bootstrap...

OK

Paste

Reset

Cancel

Help

Variables Entered/Removed ^a			
Model	Variables Entered	Variables Removed	Method
1	X1		. Stepwise (Criteria: Probability-of-F- to-enter <= .050, Probability-of-F- to-remove >= .100).

a. Dependent Variable: Y

* خروجی نرم افزار :

$$y_i = \beta_0 + \beta_1 x_{i1} + \varepsilon_i^*$$

Coefficients ^a						
Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.
		B	Std. Error	Beta		
1	(Constant)	479.145	40.527		11.823	.000
	X1	.388	.048	.757	8.014	.000

a. Dependent Variable: Y

$$\tilde{y}_i = 479.145 + .388 x_{i1}^*$$

* معایب رگرسیون قدم به قدم :

* اگرچه این روش بسیار **سریع، ساده، آسان و قابل فهم** است اما لزوماً بهترین مدل را ارائه نمی‌دهند و معایبی نیز دارد از جمله اینکه:

* در هر گام بعد از برازش مدل باید بررسی شود:

* ۱- برقراری فرضیات مدل

* ۲- عدم وجود داده پرت

* ۳- عدم وجود مشکل هم خطی چندگانه

* و تکرار عدم بررسی این موارد در هر گام باعث **کاهش اعتبار مدل** می‌شود.

* هم‌چنین در هر گام آزمونی با میزانی از خطا انجام می‌گردد و در هر مرحله این خطاها تکرار می‌شود و این امر باعث **زیاد شدن خطای مدل نهایی** می‌گردد.

* در مواردی که تعداد متغیرهای توضیحی (p) **زیاد** باشد گام‌ها طولانی شده و این روش کارایی چندانی ندارد.

* (۲۰۰۱) Harrell به بررسی معایب این روش پرداخته و موارد زیر را عنوان کرده است:

- * ۱- آریبی بر آورد گره‌های پارامترها بسیار بالاست.
- * ۲- مساله همبستگی تشدید پیدا می کند.
- * ۳- تخمین پارامترها خیلی دور از صفر خواهد بود.
- * ۴- تخمین واریانس بر آورد پارامترها دقیق نمی باشد.
- * ۵- آزمون فرضها و فواصل اطمینان نیز اشتباه بدست می آیند.

*۶- مقادیر R^2 بسیار اریب هستند.

*۷- خطاهای استاندارد برآورد پارامترها بسیار کوچک هستند.

*۸- فواصل اطمینان حول برآوردهای پارامترها باریک هست.

منابع ومراجع*

- *1- George A.F. Seber and Alan J. Lee _Linear Regression Analysis_second Edition
- *2- Loann Desboulets _ A Review on Variable Selection in Regression Analysis
- *3- L. Flom, National Development and Research Institutes, New York, NY David L. Cassell, Design Pathways, Corvallis, OR 4_ Stopping stepwise: Why stepwise and similar selection methods are bad, and what you should use_Peter Introduction to Linear Regression Analysis _Montgomery,Douglas C; Peck, Elizabeth
- *5- Harrell, F. E. (2001), Regression modeling strategies: With applications to linear models, logistic regression, and survival analysis, Springer-Verlag, New York.



با تشکر از توجه شما