

مقایسه روش های برآوردیابی و انتخاب متغیر در مدل های رگرسیونی با استفاده از داده های شبیه سازی

علی احمدی^۱، هوشنگ طالبی^۲

چکیده:

در این مقاله ضمن معرفی روش های جدیدی که در زمینه برآوردیابی و انتخاب متغیر در مدل های رگرسیونی مطرح شده اند، براساس یک بررسی مبتنی بر شبیه سازی از مدل های مختلف، به بررسی عملکرد این روش ها و همچنین مقایسه آنها با روش های معمول مثل روش انتخاب پیشرو و روش رگرسیون مرزی خواهیم پرداخت. **واژه های کلیدی:** شبیه سازی، برآوردیابی، انتخاب متغیر، گارت نامنفی، لاسو، لارس، الاستیک نت، اسکار.

۱ مقدمه

استین [۴] نشان دادند، برای تابع زبان مربع خطا، روش های انقباضی کارایی برآوردگر را افزایش می دهند. در این روش ها ضرایب رگرسیونی را با اعمال محدودیت روی دامنه تغییرات آنها برآورد می کنند. وجود چنین محدودیتی واریانس برآوردگر را کاهش داده ولی همراه با ایجاد اریبی برای برآوردگر خواهد بود (دادوستد واریانس-اریبی^۵)، بطوریکه می توان امیدوار بود در نهایت میانگین مربعات خطا کاهش یافته باشد. مدل خطی زیر را در نظر بگیرید:

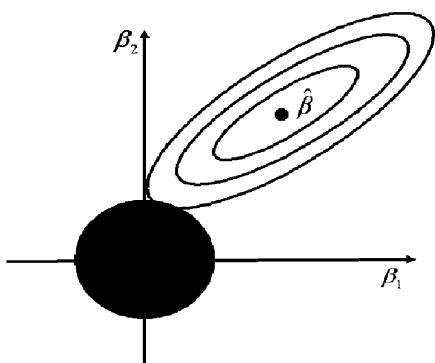
$$y = X\beta + \epsilon \quad (1)$$

که در آن $y_{(n*1)}$ بردار مشاهدات مرکزی شده متغیر وابسته، $X_{(n*P)}$ ماتریس مقادیر استاندارد شده متغیرهای پیشین، $\beta_{(P*1)}$ بردار ضرایب رگرسیونی و $\epsilon_{(n*1)}$ بردار

انتخاب متغیر و برآورد کردن ضرایب در مدل رگرسیونی اساسی ترین بخش در مدل سازی است. روش برآوردیابی حداقل مربعات، روش انتخاب متغیر به صورت پیشرو و یا حتی روش رگرسیون مرزی^۳ در مواجهه با داده هایی که از ویژگی های متفاوتی برخوردار باشند، عملکرد قابل اطمینانی از خود نشان نمی دهند. عدم پایداری، دقت پیشبینی کم و انتخاب نادرست متغیرها نمونه هایی از آسیب های مدل در هنگام استفاده از این روش ها هستند. بعلاوه این آسیب ها زمانی که همبستگی بین پیشبین ها زیاد باشد تشدید نیز می شوند. روش های انقباضی^۴ بعنوان راهکاری در جهت کاهش این آسیب ها بخصوص زمانی که همبستگی بین پیشبین ها زیاد باشد مورد توجه قرار گرفته است. جیمز و

^۱ کارشناس ارشد آمار
^۲ عضو هیات علمی دانشگاه اصفهان
^۳ Ridge Regression
^۴ Shrinkage Methods
^۵ Bias-Variance Tradeoff

و نواحی بیضوی، مجموع توان دوم مانده ها با مرکزیت برآورد



شکل ۱. عملکرد روش رگرسیون مرزی در حالت $P = 2$

حداقل مربعات معمولی (OLS) در رابطه

$$\|y - X\beta\|^2 = n\hat{\sigma}^2 + (\beta - \hat{\beta}^{OLS})' X' X (\beta - \hat{\beta}^{OLS})$$

را نشان می دهد. این شکل به وضوح ناتوانی روش مرزی در صفر برآورد کردن ضرایب را نشان می دهد، زیرا برخورد دو ناحیه نمی تواند در نقطه ای باشد که در آن یکی از ضرایب صفر است. روش های جدید انقباضی میل دادن ضرایب به سمت صفر را به گونه ای انجام می دهند که بعضی از ضرایب، مربوط به متغیرهای بی اثر، دقیقاً صفر برآورد شده و به این ترتیب متغیر مربوط به آن ضرایب از مدل خارج خواهد شد. بنابراین در این روش ها برآوردیابی و انتخاب متغیر تواما صورت می پذیرد. باید توجه داشت که در تمام روش های انقباضی ابتدا به ازاء مقادیر مختلف پارامتر کنترل برآورد ضرایب محاسبه شده و سپس با استفاده از معیارهای ارزیابی مدل، مثل $AIC, BIC, C_P, Cross - Validation, \dots$ ، برآورد بهینه از میان مجموعه برآوردهای بدست آمده

متغیرهای تصادفی خطاست. در این مدل ستون های ماتریس X را x_j و عناصر آن را با x_{ij} نمایش می دهیم. در دو دهه اخیر روش های انقباضی مختلفی برای برآورد ضرایب رگرسیونی ارائه شده است. در بخش دوم برخی از این روش ها و ویژگی های آنها را مرور خواهیم کرد. برای مقایسه این روش هاد شرایط مختلف، در بخش سوم مدل های متنوعی را شبیه سازی کرده و ضرایب رگرسیونی را با استفاده از روش های ذکر شده در بخش دوم برآورد می کنیم. در بخش چهارم نتایج بدست آمده را تحلیل و درباره آسیب ها و مزایای این روش ها در شرایط مختلف بحث می کنیم.

۲ روش های انقباضی

روش های انقباضی از جمله روش های جدید در امر برآوردیابی و انتخاب متغیر در مدل های رگرسیونی هستند. روش رگرسیون مرزی یکی از اولین روش های انقباضی است. این روش ضرایب رگرسیونی را بصورت زیر برآورد می کند:

$$\hat{\beta}_{Ridge} = \underset{\beta}{\operatorname{argmin}} \|y - X\beta\|^2 \quad \text{s.t.} \quad \sum_{j=1}^P \beta_j^2 \leq t \quad (2)$$

که در آن $t \geq 0$ پارامتر کنترل^۶ بوده و میزان انقباض تحمیل شده به ضرایب را کنترل می کند. ایراد این روش استفاده از یک تابع تاوان^۷ توان دوم از ضرایب رگرسیونی است که مانع از صفر برآورد شدن ضرایب و در نتیجه عدم حذف متغیر از مدل می شود ([۵]).

در شکل (۱) ناحیه دایره ای، ناحیه تاوان $(\beta_1^2 + \beta_2^2 \leq t)$

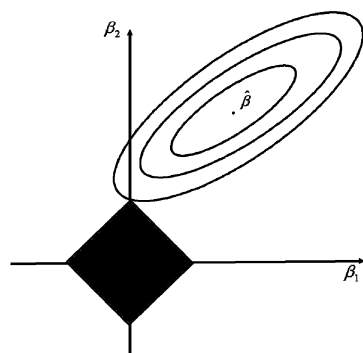
^۶ Tuning Parameter
^۷ Penalty Function

تابع تاوان مجموع قدرمطلق ضرایب استفاده می کند:

$$\hat{\beta}_{Lasso} = \underset{\beta}{\operatorname{argmin}} \|y - X\beta\|^2 \quad (4)$$

$$s.t. \sum_{j=1}^P |\beta_j| \leq t$$

این تغییر در نوع تابع تاوان امکان صفر برآورد شدن برخی از ضرایب را پدید می آورد.



شکل ۲. عملکرد روش لاسو در حالت $P = 2$

در شکل ۲ ناحیه لوزی شکل ناحیه تاوان ($|\beta_1| + |\beta_2| \leq t$) و نواحی بیضوی مجموع توان دوم مانده ها با مرکزیت برآورد حداقل مربعات را نشان می دهد. این شکل به وضوح توانایی روش لاسو در صفر برآورد کردن ضرایب را نشان می دهد، زیرا برخورد دو ناحیه می تواند در نقطه ای صورت گیرد که در آن یکی از ضرایب مقدار صفر را دارد. روش لاسو به لحاظ پایداری و دقت پیشبینی برآوردگرها عملکرد قابل قبولی را از خود نشان داده است. این روش همچنین بعنوان روشی پایه ای در ساخت روش های پیچیده تر نیز مورد استفاده قرار گرفته است.

انتخاب می شود. در ادامه این بخش روش های مختلف را به صورت مختصر مرور می کنیم.

۱.۲ روش گارت نامنفی

برآورد ضرایب رگرسیونی در روش گارت نامنفی^۸ [۲] به صورت زیر محاسبه می شوند:

$$\hat{c} = (\hat{c}_1, \dots, \hat{c}_P)' = \underset{c}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \sum_{j=1}^P c_j \hat{\beta}_j x_{ij})^2 \quad (3)$$

$$s.t. \sum_{j=1}^P c_j \leq s \quad \text{and} \quad c_j \geq 0$$

که در آن $s \geq 0$ و $\hat{\beta}_j$ برآورد حداقل مربعات β_j بوده که به عنوان برآورد اولیه استفاده شده است. در این روش برآورد نهایی به فرم $\hat{\beta}_j^G = \hat{c}_j \hat{\beta}_j$ خواهد بود. ضرایب \hat{c}_j عامل های انقباض هستند که مقادیر کوچک یا صفر آنها باعث کوچک یا صفر شدن $\hat{\beta}_j$ ها در مدل خواهد شد. این روش از نظر دقت پیشبینی و پایداری به خوبی روش مرزی بوده و در ضمن با صفر برآورد کردن بعضی از ضرایب عمل انتخاب متغیر را نیز انجام می دهد. ایراد اصلی این روش استفاده از برآورد حداقل مربعات بعنوان برآورد اولیه است. زیرا معایب برآوردگر حداقل مربعات ممکن است برآورد نهایی را نیز تحت تاثیر قرار دهد این عیب را معمولاً می توان با استفاده از دیگر برآوردگرها به عنوان برآورد اولیه برطرف کرد.

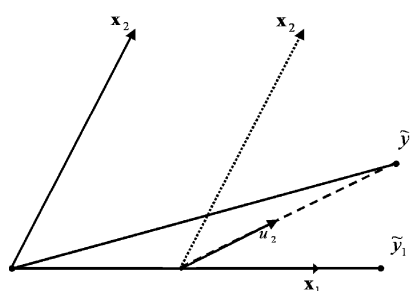
۲.۲ روش لاسو

روش لاسو^۹ [۵] اساساً مشابه روش رگرسیون مرزی است، با این تفاوت که بجای استفاده از تابع تاوان درجه دواز

^۸ Nonnegative Garrote
^۹ (Least Absolute Shrinkage and Selection Operator)Lasso

۳.۲ روش لارس

حالت $P = 3$ با کمک شکل ۳ تشریح می کنیم:
در شکل ۳ $X = (x_1, x_2)$ و \tilde{y}_2 برآورد حداقل مربعات
بردار y است. بنابراین $X'(y - \hat{\mu}) = X'(\tilde{y}_2 - \hat{\mu})$.



شکل ۳. عملکرد روش لارس در حالت $p = 2$

الگوریتم لارس با $\hat{\mu}_0 = 0$ آغاز می شود. با توجه به شکل
۳ $\tilde{y}_2 - \hat{\mu}_0$ زاویه کمتری با متغیر x_1 نسبت به x_2 می
سازد، یعنی همبستگی x_1 با y بیشتر است، بنابراین x_1
وارد مدل می شود. سپس $\hat{\mu}_0$ در جهت x_1 تصحیح می
شود یعنی در اولین گام خواهیم داشت: $\hat{\mu}_1 = \hat{\mu}_0 + \hat{\gamma}_1 x_1$
که در آن $\hat{\gamma}_1$ طوری انتخاب می شود که به ازاء آن:
 $\hat{\mu}_1 - \tilde{y}_2$ نیمساز زاویه بین $x_1 - \hat{\mu}_1$ و x_2 باشد. اگر u_2
برداریکه ای باشد که در طول نیمساز قرار گیرد آنگاه
برآورد بعدی لارس در گام دوم به فرم $\hat{\mu}_2 = \hat{\mu}_1 + \hat{\gamma}_2 u_2$
خواهد بود. در این رابطه $\hat{\gamma}_2$ طوری انتخاب می شود که
همبستگی x_1 و x_2 با مانده ها صفر شود، یعنی $\tilde{y}_2 = \hat{\mu}_2$.
برآوردهای بدست آمده در روش لارس تا حدود زیادی
به روش لاسو نزدیک است و برخلاف روش لاسو، در
روش لارس برآورد ضرایب نیازی به استفاده از الگوریتم
های پیچیده ریاضی ندارد.

روش لارس^{۱۰} [۳] با رویکردی شبیه روش انتخاب پیشرو
نسبت به ورود متغیرها به مدل عمل می کند. بدین
ترتیب که در هر گام یک متغیر پیشین به صورت زیر وارد
مدل می شود:

۱- ابتدا همه ضرایب صفر در نظر گرفته می شوند،
سپس متغیری انتخاب می شود که بیشترین همبستگی را
با متغیر وابسته داشته باشد، آنرا x_1 می نامیم.

۲- سپس بزرگترین گام ممکن در جهت این متغیر (با
تغییر ضریب مربوط به آن)، تا جایی که همبستگی
متغیری دیگر مثل x_2 با مانده ها برابر با همبستگی
متغیر x_1 با مانده ها شود برداشته می شود (این برخلاف
روش پیشرو است که در آن تا جایی پیش می رویم که
همبستگی متغیر x_1 با مانده ها صفر شود).

۳- در مرحله بعد در مسیر متساوی الزاویه بین x_1 و x_2
پیش می رویم، تا جایی که همبستگی متغیری دیگر مثل
 x_3 با مانده برابر با همبستگی دو متغیر قبلی با مانده ها
شده و بتواند وارد مدل شود (توجه کنید که با حرکت در
مسیر متساوی الزاویه، همبستگی دو متغیر x_1 و x_2 با
مانده ها به یک میزان کاهش می یابد).

به همین ترتیب در هر مرحله در مسیر متساوی الزاویه بین
متغیرهای داخل مدل حرکت کرده تا اینکه همبستگی
متغیری دیگر با مانده برابر با همبستگی متغیرهای داخل
مدل با مانده شده و وارد مدل شود. این روند تا ورود
همه متغیرها به مدل و صفر شدن ماکزیمم قدرمطلق
همبستگی ها ادامه می یابد. مراحل الگوریتم لارس را در

^{۱۰} LARS (Least Angle Regression)

۴.۲ روش الاستیک نت

برآورد اولیه در روش الاستیک نت^{۱۱} [۷] به صورت زیر محاسبه می شود:

$$\hat{\beta}_{NEN} = \underset{\beta}{\operatorname{argmin}} (\|y - X\beta\|^2 + \lambda_2 \|\beta\|^2 + \lambda_1 \|\beta\|_1) \quad (5)$$

که در آن λ_1 و λ_2 مقادیری نامنفی و با در نظر گرفتن $\|\beta\|_1 = \sum_{j=1}^P |\beta_j|$ می توان برآورد الاستیک نت اولیه را به فرم توانی زیر نوشت:

$$\hat{\beta}_{NEN} = \underset{\beta}{\operatorname{argmin}} \|y - X\beta\|^2 \quad (6)$$

$$s.t. \alpha \|\beta\|^2 + (1 - \alpha) \|\beta\|_1 \leq t \text{ and } t \geq 0$$

ناحیه توان در این روش ترکیبی اکیدا محدب از نواحی توان در روش های لاسو و مرزی است. این فرم از تابع توان علاوه بر قابلیت صفر برآورد کردن برخی از ضرایب، توانایی برابر برآورد کردن ضرایب متغیرهایی که اثر یکسان روی متغیر پاسخ دارند و یا به شدت همبسته هستند را نیز دارد. به این ترتیب از این روش می توان برای گروه بندی پیشین ها به خصوص زمانی که تعداد آنها زیاد باشد استفاده کرد. برآورد اصلاح شده الاستیک نت که دارای دقت پیشینی بالاتری نسبت به برآورد اولیه است بصورت $\hat{\beta}_{EN} = (1 + \lambda_2)\hat{\beta}_{NEN}$ تعریف می شود که در آن $\hat{\beta}_{EN}$ و $\hat{\beta}_{NEN}$ به ترتیب برآوردگر الاستیک نت و برآوردگر الاستیک نت اولیه (خام)^{۱۲} هستند. روش الاستیک نت در اکثر مواقع از روش های گارت نامنفی، لاسو، لارس و مرزی دقت پیشینی بالاتری دارد.

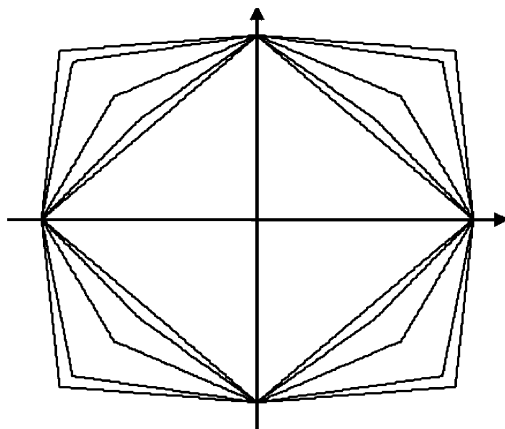
۵.۲ روش اسکار

روش اسکار^{۱۳} [۱] برآورد ضرایب رگرسیونی را بصورت زیر انجام می دهد:

$$\hat{\beta}_{OSCAR} = \underset{\beta}{\operatorname{argmin}} \|y - X\beta\|^2 \quad (7)$$

$$s.t. \sum_{j=1}^P |\beta_j| + c \sum_{j < k} \max(|\beta_j|, |\beta_k|) \leq t$$

تابع توان به کار رفته در این روش به لحاظ هندسی هشت گوشه است و این روش همانند روش الاستیک نت، توانایی برابر برآورد کردن پیشین های به شدت همبسته و یا آنهایی که اثری یکسان روی متغیر پاسخ دارند را دارد. یعنی علاوه بر برآوردیابی و انتخاب متغیر قابلیت گروه بندی آنها نیز در این روش وجود دارد.



شکل ۴. ناحیه توان در روش اسکار در حالت $p = 2$

شکل ۴ نواحی توان را به ازاء مقادیر مختلف c نشان می دهد. همانطور که در اشکال ۱ و ۲ هم نشان داده شد، در اینجا در صورتی که همبستگی بین دو پیشین زیاد باشد، برخورد نواحی بیضوی با ناحیه توان می تواند در گوشه ای از ناحیه توان رخ دهد که در آن، ضرایب با مقداری

^{۱۱} Elastic Net

^{۱۲} Naive Elastic Net

^{۱۳} (Octagonal Shrinkage and Clustering Algorithm for Regression) OSCAR

می کنیم که در آن $n = 100$ ، $\epsilon \sim N_n(0, \Omega^2 I)$ ،
 $cov(x_i, x_j) = 0.5$ ، $i \neq j$ و $Var(x_i) = 1$ و

در این مدل $\beta = (0, \dots, 0, 2, \dots, 2, 0, \dots, 0, 2, \dots, 2)'$
 گروه هایی از متغیرهای موثر و غیرموثر وجود دارد که
 همبستگی بین آنها متوسط و یکسان است.

شبه سازی چهارم: 100 مجموعه داده
 از مدل ۱ تولید می کنیم که در آن
 $n = 40$ ، $\epsilon \sim t_n(0, I, 5)$ ، $cov(x_i, x_j) = 0.5^{|i-j|}$ و
 $\beta = (3, 1.5, 0, 0, 0, 2, 0, 0)'$ این مدل، داده هایی را
 نشان می دهد که در آنها ضمن وجود مشاهدات پرت،
 متغیرهای موثر اول و دوم با همدیگر همبستگی متوسط و
 با متغیر موثر ششم همبستگی کمی دارند.

شبه سازی پنجم: 100 مجموعه داده از مدل ۱
 تولید می کنیم که در آن $\epsilon_i = \rho \epsilon_{i-1} + u_i$ ، $u_i \sim N_n(0, 0.5)$
 $\beta = (3, 1.5, 0, 0, 0, 2, 0, 0)'$ ، $\rho = 0.4$ و $n = 40$ ،
 $cov(x_i, x_j) = 0.5^{|i-j|}$ این مدل داده هایی را نشان می دهد که در آنها ϵ_i ها دارای
 خودهمبستگی هستند.

شبه سازی ششم: 100 مجموعه داده از مدل ۱
 تولید می کنیم که در آن $n = 30$ ، $\epsilon \sim N_n(0, 3^2 I)$ و
 $\beta = (3, 1.5, 2, 0, 0, 0, 2, 0)'$ در این مدل ساختار
 ماتریس کواریانس X به فرم زیر است:

جدول ۱. ساختار ماتریس کواریانس X

1	0.9	0.85	0.7	0.4	0.4	0.4	0.4
0.9	1	0.9	0.8	0.4	0.4	0.4	0.4
0.85	0.9	1	0.8	0.4	0.4	0.4	0.4
0.7	0.8	0.8	1	0.4	0.4	0.4	0.4
0.4	0.4	0.4	0.4	1	0.4	0.4	0.4
0.4	0.4	0.4	0.4	0.4	1	0.4	0.4
0.4	0.4	0.4	0.4	0.4	0.4	1	0.8
0.4	0.4	0.4	0.4	0.4	0.4	0.8	1

برابر برآورد شوند. روش اسکار دارای دقت پیشبینی بالا
 و خاصیت گروهی قوی می باشد.

۳ شبیه سازی

در این بخش با استفاده از شبیه سازی به ارزیابی عملکرد
 و مقایسه روش های پیش گفته می پردازیم. از آنجایی
 که روش های انقباضی برای غلبه بر معایب روش های
 قدیمی، بخصوص زمانی که همبستگی بین پیشبین ها
 زیاد باشد، طراحی شده اند مدل های شبیه سازی شده را
 بیشتر براساس ساختارهای مختلف ماتریس کواریانس X
 مورد مطالعه قرار می دهیم. ابتدا ویژگی های داده های
 شبیه سازی شده را بیان و سپس نتایج کاربرد روش ها را
 ارائه می دهیم.

شبه سازی اول: 100 مجموعه داده از مدل ۱
 تولید می کنیم که در آن $n = 30$ ، $P = 8$ ،
 $\epsilon \sim N_n(0, 3^2 I)$ ، $cov(x_i, x_j) = 0.7^{|i-j|}$ و
 $\beta = (3, 2, 1.5, 0, 0, 0, 0, 0)'$ در این مدل، متغیرهای
 موثر با همدیگر همبستگی نسبتا بالایی دارند. متغیرهای
 بی اثر نیز با همدیگر تا حدودی همبسته بوده ضمن اینکه
 متغیر موثر سوم و متغیر بی اثر چهارم نیز با همدیگر
 همبستگی نسبتا بالایی دارند.

شبه سازی دوم: شبه سازی مدل دوم
 مشابه مدل اول است با این تفاوت که در آن
 $\beta = (3, 0, 0, 1.5, 0, 0, 0, 2)'$ در نتیجه در این مدل
 ، متغیرهای موثر همبستگی ناچیزی بایکدیگر داشته ولی
 با متغیرهای بی اثر همبستگی نسبتا بالایی دارند.

شبه سازی سوم: 100 مجموعه داده از مدل ۱ تولید

مرزی، لارس و الاستیک نت را تکرار می کنیم. بعلاوه از معیار $ME = (\hat{\beta} - \beta)' V(\hat{\beta} - \beta)$ ، که در آن V ماتریس کواریانس X است، برای انتخاب برآورد بهینه در بین مجموعه برآوردهایی که در هر تکرار برای هر روش محاسبه می شوند استفاده شده است (برای جزئیات بیشتر در مورد معیار ME به [۵] مراجعه کنید). نتایج حاصل از شبیه سازی ها در جداول ۱، ۲ و ۳ ارائه شده اند (در این جداول عبارت «د.م.ا.د.» نشان دهنده درصد انتخاب مدل درست و عبارت «م.ت.ض.غ.ص.» نشان دهنده میانه تعداد ضرایب غیر صفر برآورد شده توسط هر روش در تکرارهای مختلف هر مدل است).

داده های شبیه سازی شده در این مدل مشابه داده هایی هستند که در آنها گروه هایی از متغیرهای به شدت همبسته وجود دارند. در این مدل متغیرهای پیشین اول تا سوم موثر بوده و با متغیر بی اثر چهارم تشکیل یک گروه به شدت همبسته را می دهند. متغیر هفتم نیز موثر بوده و با متغیر بی اثر هشتم تشکیل یک گروه به شدت همبسته را می دهند.

برای داده های تولید شده در این شبیه سازی ها به ترتیب از روش های برآوردیابی حداقل مربعات معمولی، انتخاب پیشرو، روش مرزی، گارت نامنفی، روش لاسو، لارس، الاستیک نت و اسکار ضرایب رگرسیونی را برآورد و همچنین روش گارت نامنفی با برآوردهای اولیه

جدول ۱. نتایج مربوط به شبیه سازی های اول و دوم

مدل ۲			مدل ۱			روش برآوردیابی
م.ت.ض.غ.ص.	د.م.ا.د.	میانه ME (انحراف معیار)	م.ت.ض.غ.ص.	د.م.ا.د.	میانه ME (انحراف معیار)	
۸	۰	۳.۲۹۶۸ (۰.۲۵۱۴)	۸	۰	۳.۳۴۰۹ (۰.۳۴۱۸)	<i>OLS Estimation</i>
۳	۴۷	۱.۶۶۶۹ (۰.۱۹۸۳)	۳	۱۵	۲.۰۶۱۸ (۰.۱۴۸۰)	<i>Forward Method</i>
۸	۰	۱.۸۰۴۷ (۰.۰۹۹۳)	۸	۰	۱.۵۷۰۱ (۰.۱۹۷۲)	<i>Ridge Regression</i>
۴	۱۱	۱.۵۴۰۷ (۰.۱۸۵۵)	۴	۸	۱.۸۹۲۳ (۰.۱۵۴۳)	<i>N - Garrote(OSL)</i>
۵	۳	۱.۴۶۷۰ (۰.۱۲۵۷)	۵	۱۳	۱.۳۵۳۸ (۰.۱۴۸۲)	<i>Lasso</i>
۵	۵	۱.۶۰۹۲ (۰.۱۳۴۶)	۴	۲۳	۱.۴۲۴۷ (۰.۱۴۶۲)	<i>Lars</i>
۵	۹	۱.۳۲۶۴ (۰.۱۲۳۸)	۳	۴۴	۰.۹۷۰۴ (۰.۱۱۸۷)	<i>Elastic Net</i>
۴	۷	۱.۵۵۶۹ (۰.۱۰۰۶)	۶	۱۲	۱.۰۴۴۷ (۰.۱۳۷۹)	<i>Oscar</i>
۴	۱۴	۱.۳۸۵۱ (۰.۱۲۷۹)	۴	۱۸	۱.۲۴۲۰ (۰.۱۲۷۲)	<i>N - Garrote(Ridge)</i>
۴	۱۹	۱.۳۲۵۴ (۰.۱۶۰۵)	۳	۲۹	۱.۳۸۹۵ (۰.۱۴۲۶)	<i>N - Garrote(Lars)</i>
۴	۲۲	۱.۲۹۴۵ (۰.۱۲۰۲)	۳	۵۱	۱.۱۳۷۹ (۰.۰۸۹۹)	<i>N - Garrote(EN)</i>

جدول ۲. نتایج مربوط به شبیه‌سازی‌های سوم و چهارم

مدل ۴			مدل ۳			روش برآوردیابی
م.ت.ض. غ.ص.	د.م.ا.د.	میانۀ ME (انحراف معیار)	م.ت.ض. غ.ص.	د.م.ا.د.	میانۀ ME (انحراف معیار)	
۸	۰	۰.۵۳۹۵ (۰.۰۲۶۱)	۴۰	۰	۱۵۲.۲۰ (۵.۷۱۸)	OLS Estimation
۳	۶۰	۰.۳۰۱۲ (۰.۰۲۷۳)	۹	۰	۸۳.۶۸ (۲.۳۱۹)	Forward Method
۸	۰	۰.۴۲۳۷ (۰.۰۲۶۶)	۴۰	۰	۲۰.۸۶ (۰.۹۹۵)	Ridge Regression
۴	۲۵	۰.۲۴۸۸ (۰.۰۲۹۹)	۱۵	۰	۷۰.۶۲ (۱.۹۰۵)	N - Garrote(OSL)
۶	۹	۰.۳۲۲۵ (۰.۰۲۲۷)	۲۰	۰	۴۲.۶۷ (۱.۳۷۵)	Lasso
۵	۱۳	۰.۳۴۱۳ (۰.۰۲۵۶)	۲۱	۰	۴۴.۵۹ (۱.۴۸۶)	Lars
۴	۳۷	۰.۱۸۰۱ (۰.۰۲۳۳)	۲۰	۰	۴۱.۹۱ (۱.۴۷۲)	Elastic Net
۶	۱۲	۰.۲۹۴۸ (۰.۰۳۴۴)	۳۱	۰	۲۱.۷۷ (۱.۷۵۷)	Oscar
۴	۲۴	۰.۲۴۵۴ (۰.۰۳۰۱۰)	۱۹	۰	۵۱.۲۹ (۱.۷۷۳)	N - Garrote(Ridge)
۴	۳۵	۰.۲۳۶۳ (۰.۰۳۳۹)	۱۷	۰	۴۵.۰۲ (۱.۵۶۹)	N - Garrote(Lars)
۳	۵۴	۰.۲۲۳۹ (۰.۰۳۰۸)	۱۷	۰	۴۲.۰۷ (۱.۴۹)	N - Garrote(EN)

جدول ۳. نتایج مربوط به شبیه‌سازی‌های پنجم و ششم

مدل ۶			مدل ۵			روش برآوردیابی
م.ت.ض. غ.ص.	د.م.ا.د.	میانۀ ME (انحراف معیار)	م.ت.ض. غ.ص.	د.م.ا.د.	میانۀ ME (انحراف معیار)	
۸	۰	۳.۷۵۱ (۰.۲۴۶)	۸	۰	۰.۳۰۱۷ (۰.۰۲۹۲)	OLS Estimation
۴	۱	۲.۴۰۳ (۰.۲۱۴)	۳	۵۷	۰.۲۳۳۱ (۰.۰۲۹۹)	Forward Method
۸	۰	۱.۶۸۶ (۰.۱۲۴)	۸	۰	۰.۲۳۳۸ (۰.۰۲۳۱)	Ridge Regression
۵	۱	۲.۲۷۹ (۰.۲۰۱)	۴	۲۷	۰.۲۰۷۵ (۰.۰۳۰۵)	N - Garrote(OSL)
۵	۵	۱.۸۸۹ (۰.۱۷۰)	۶	۲۰	۰.۱۹۷۷ (۰.۰۲۵۲)	Lasso
۵	۶	۱.۹۱۷ (۰.۱۷۸)	۶	۲۲	۰.۲۲۱۴ (۰.۰۲۴۰)	Lars
۵	۱۷	۱.۳۲۶ (۰.۱۳۳)	۴	۴۴	۰.۱۱۱۷ (۰.۰۱۸۱)	Elastic Net
۷	۱	۱.۳۵۱ (۰.۱۱۴)	۷	۹	۰.۱۹۵۲ (۰.۰۳۱۸)	Oscar
۵	۴	۲.۰۵۳ (۰.۱۸۱)	۴	۳۲	۰.۲۰۴۲ (۰.۰۳۰۱)	N - Garrote(Ridge)
۵	۸	۱.۸۱۵ (۰.۱۸۱)	۴	۳۶	۰.۲۰۴۶ (۰.۰۳۱۸)	N - Garrote(Lars)
۴	۱۷	۱.۷۹۸ (۰.۱۸۲)	۳	۵۱	۰.۲۰۳۸ (۰.۰۳۰۰)	N - Garrote(EN)

۴ تحلیل نتایج

خواهیم پرداخت. لازم به ذکر است که در این بررسی از

در این بخش بصورت مختصر به تحلیل نتایج بدست

آمده و ذکر معایب و محاسن هر یک از این روش‌ها

روش ها است. در این جداول مشاهده می شود که عیب اصلی این روش یعنی غیرصفر برآورد کردن ضرایب و در نتیجه عدم انتخاب متغیر در همه جا به چشم می خورد. می توان از این روش برای انتخاب برآورد اولیه در روش گارت نامنفی استفاده کرد.

روش گارت نامنفی با برآورد اولیه حداقل مربعات از نظر دقت پیشبینی از روش پیشرو بهتر است اما به دلیل استفاده از برآورد حداقل مربعات بعنوان برآورد اولیه از معایب این برآوردگر مستثنی نیست. البته همانطور که قبلا هم اشاره شد می توان از دیگر برآوردگرها بعنوان برآورد اولیه استفاده کرد که این تکنیک در اینجا مورد استفاده قرار گرفته است. نتایج نشان می دهد که استفاده از برآوردگرهایی که نسبت به برآورد حداقل مربعات بهتر هستند (مثل برآوردگر رگرسیون مرزی) بعنوان برآورد اولیه، باعث بهتر شدن نتایج در این روش می شود. اما باید توجه داشت که استفاده از روش گارت نامنفی بعد از استفاده از روش انقباضی دیگری، تضمینی برای بهتر کردن دقت پیشبینی نمی دهد. این مسیله ناشی از انقباض مضاعف^{۱۵} است (یعنی اعمال آریبی بیش از حد به برآوردگر در حالی که باعث کاهش واریانس به اندازه کافی نشود).

بررسی نتایج بدست آمده نشان می دهد روش لاسو در مقایسه با روش گارت نامنفی (با برآورد اولیه حداقل مربعات) و روش رگرسیون مرزی در اکثر مواقع دقت پیشبینی بالاتری داشته و در مقایسه با روش پیشرو کمتر پرخور است. مشاهده برآوردهای بدست آمده برای

میان مقادیر ME بدست آمده که حاصل از محاسبه این معیار برای ۱۰۰ مجموعه داده از هر مدل و در هر روش است، برای مقایسه خطای پیشبینی روش ها استفاده شده است. از درصد تعداد دفعاتی که هر روش مدل درست را انتخاب کرده و همچنین میانه تعداد ضرایب غیر صفر برآورد شده توسط هر روش در هر مدل برای مقایسه توانایی روش ها در انتخاب مدل صحیح استفاده شده است.

انتخاب متغیر در مدل معمولا باعث بالا رفتن دقت پیشبینی می شود. مدل هایی که در آنها برآورد OLS برای همه ضرایب متغیرها محاسبه شده و عملا انتخاب متغیری در آنها صورت نگرفته است دارای کمترین میزان دقت پیشبینی هستند. در مورد روش انتخاب پیشرو معروف است که این روش پرخور^{۱۴} است، میانه تعداد ضرایب غیرصفر برآورد شده در جداول ۱ تا ۳ نیز گویای همین واقعیت است. این مسئله بخصوص در مدل ۳ مشهودتر است. مشاهده ضرایب بدست آمده نشان می دهد که این روش در انتخاب و تمایز بین متغیرهای موثر و غیرموثر همبسته توانایی کافی را ندارد. همچنین با زیاد شدن تعداد پیشبین ها و یا همبستگی بین آنها، از دقت پیشبینی و توانایی انتخاب مدل درست در این روش کاسته می شود.

مقایسه جداول ۱ تا ۳ برای روش رگرسیون مرزی نشان می دهد که این روش، معمولا نسبت به روش انتخاب پیشرو دقت پیشبینی بالاتری دارد. این روش بخصوص در مدل ۳، دارای بهترین دقت پیشبینی در بین سایر

متغیر ارائه می دهیم. باید توجه داشت که مشخص کردن ساختار ماتریس کواریانس متغیرهای پیشین برای انتخاب درست یک روش برآوردیابی و انتخاب متغیر می تواند مفید باشد.

برای مدل هایی با تعداد متغیرهای متوسط یا کم (مثلا $10 \leq P$) روش الاستیک نت هم در حالتی که همبستگی بین پیشبین ها زیاد باشد و هم در حالتی که این همبستگی چندان قابل توجه نباشد عملکرد قابل قبولی از خود نشان داده است. البته در حالتی که همبستگی بین پیشبین ها زیاد نباشد می توان از روش های دیگری مثل لارس و لاسو نیز استفاده کرد. بخصوص از روش لارس که هم از نظر دقت پیشبینی تقریبا مشابه روش لاسو بوده و هم از نظر شیوه محاسبات به سادگی روش پیشرو است. در این حالت اگر تعداد نمونه ها زیاد باشد می توان از روش گارت نامنفی نیز استفاده کرد (به دلیل خواص جانبی آن [۶]). البته پیشنهاد می شود از برآورد مرزی بعنوان برآورد اولیه در روش گارت نامنفی استفاده شود. در این صورت هم از خواص خوب روش مرزی در دقت پیشبینی بهره برده ایم و هم خاصیت تنک بودن روش گارت و در ضمن تا حد امکان از انقباض مضاعف دوری جسته ایم. اما برای حالتی که همبستگی بین پیشبین ها زیاد باشد تنها استفاده از روش های الاستیک نت و اسکار توصیه می شود. در این حالت اگر تعداد متغیرهای پیشبین خیلی زیاد بوده و گروه بندی آنها مورد علاقه باشد می توان از روش اسکار استفاده کرد، زیرا این روش خاصیت گروه بندی قویتری نسبت به روش الاستیک نت دارد. در حالتی که تعداد متغیرهای پیشبین زیاد

روش لاسو نشان می دهد که یکی از ایرادات این روش عدم ثبات کافی در انتخاب متغیرهای موثر، زمانی که داده ها شامل گروه هایی از متغیرهای پیشین به شدت همبسته باشند، است.

نتایج بدست آمده در مورد روش لارس، همانطور که قبلا هم اشاره شد، بسیار به روش لاسو نزدیک است. یک چنین ارتباطی از نظر تئوری هم بین این دو روش وجود دارد (رجوع شود به [۱]).

روشی که از هر نظر نسبت به روش های دیگر برتری داشته روش الاستیک نت است. این روش بجز در مدل ۳، که مربوط به متغیرهای گروهی است، در سایر مدل ها از دیگر روش ها بهتر عمل کرده است. روش الاستیک نت هم از نظر دقت پیشبینی و هم از نظر انتخاب مدل صحیح از دیگر روش ها بهتر عمل کرده است. این روش از نظر تمایز بین متغیرهای موثر و غیر موثر همبسته هم بهتر از سایر روش ها عمل کرده است. روش اسکار نیز در اکثر مواقع از نظر دقت پیشبینی به خوبی روش الاستیک نت بوده و حتی در مدل ۳ از این روش نیز بهتر عمل کرده است. مشاهده برآوردهای بدست آمده در مدل های مختلف نشان می دهد که روش اسکار دارای اثر گروهی قویتری نسبت به روش الاستیک نت بوده است. البته تمایل این روش به برابر برآورد کردن متغیرهای به شدت همبسته باعث ضعف این روش در انتخاب مدل درست، وقتی همبستگی بین متغیر موثر و غیر موثر بالا باشد، شده است.

جمع بندی و ارائه چند پیشنهاد:

بر اساس نتایج شبیه سازی های انجام شده پیشنهادهایی را در مورد استفاده از روش های برآوردیابی و انتخاب

باشد (مثلا $p \geq 10$)، در دو حالت همبستگی زیاد یا کم استفاده شود. پیشنهاد می شود تنها از روش الاستیک نت و یا اسکار

مراجع

- [1] Bondell, H. Reich, B.(2008), Simultaneous Regression Shrinkage, Variable Selection, and Supervised Clustering of Predictors with OSCAR, *Biometrics*, 64, 115-123
- [2] Breiman, L. (1995), Better subset regression using the nonnegative garrote, *Technometrics*, 37, 373-384.
- [3] Efron, B., Hastie, T., Johnstone, I. and Tibshirani, R (2004). Least Angle Regression, *Ann.Statist*, 32, 407-499.
- [4] James, W., C.Stein (1961), Estimation with quadratic loss, *Proceeding of the forth berkely symposium*, 1, 361-379.
- [5] Tibshirani, R.(1996), Regression shrinkage and selection via the lasso, *J. Roy. Statist.Soc.Ser*, B.58, 267-288.
- [6] Yuan, M., Lin, Y.(2007), On the nonnegative garrote estimator, *J. Roy. Statist.Soc.Ser*, B.69, 143-161.
- [7] Zou, H., Hastie, T.(2005), Regularization and variable selection via the elastic net, *J.R Statist.Sco*, B.67, 301-320.