

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

ع. ۱۱۸

۲۲۸۱ / ۴ / ۲۶



دانشگاه صنعتی اصفهان
دانشکده علوم ریاضی

وزارتخانه صنعت، معدن و تجارت
جمهوری اسلامی ایران

برآورد در رگرسیون خطی ساده زمانی که
ناهمگونی خطای مدل به فرم نامعلومی باشد

پایان نامه کارشناسی ارشد آمار
صدیقه میرزائی صالح آبادی

۴. ۸۱۸

استاد راهنما
دکتر احمد پارسیان
دکتر علی زینل همدانی

۱۳۸۰

۴. ۸۱۸



دانشگاه صنعتی اصفهان

دانشکده علوم ریاضی

پایان نامه کارشناسی ارشد رشته آمار خانم صدیقه میرزائی صالح آبادی

تحت عنوان

برآورد در رگرسیون خطی ساده زمانی که

ناهمگونی خطای مدل به فرم نامعلومی باشد

در تاریخ ۸۰/۷/۱۱ توسط کمیته تخصصی زیر مورد بررسی و تصویب نهائی قرار گرفت.

دکتر احمد پارسیان و دکتر علی زینل همدانی

۱- استاذ راهنمای پایان نامه

دکتر علی رجالی

۲- استاد مشاور پایان نامه

دکتر عبدالرسول برهانی حقیقی

۳- استاد داور ۱

دکتر سروش علیمرادی

۴- استاد داور ۲

دکتر امیر نادری

سرپرست تحصیلات تکمیلی دانشکده

الحمد لله رب العالمين

...رحمت واسعه پروردگار فرصتی مغتنم داد تا به اقتضای توان خود از محضر اساتیدی گرانقدر بهره جسته و ره توشه‌ای برگیرم، بی‌پیرایه بهترین ایام زندگی‌م را در محضر این عزیزان که نه فقط رهنمون به دانشم بودند، بلکه درس اخلاق و تواضع را آموختند سپری کرده‌ام. تمام مباحث من در طول تحصیل نه دست‌یازیدن به درجه‌ای از دانش، بلکه فراسوی آن شاگردی نزد معلمانی بوده‌است که خود دریایی از معرفت هستند و سهم من پرتویی از تشعشع معرفت ایشان بر اندیشه بوده‌است، انسانهای بزرگواری که حق‌قدردانی بر عهده اینجانب دارند:

معلمان ارجمندم، دکتر احمد پارسیان و دکتر علی زینل همدانی

مشاور پرتلاشم، دکتر علی رجالی

ناظرین محترم پایان‌نامه‌ام، دکتر عبدالرسول برهانی حقیقی و دکتر سروش علیمرادی

اساتید همراهم، دکتر ایوب ساعی، دکتر محمد صالحی

سرپرست محترم تحصیلات تکمیلی دانشکده، دکتر امیر نادری

اساتید گرانقدر، کارشناسان، مسؤولین و کارکنان محترم دانشکده علوم ریاضی

دوستان و همکلاسیهای خوبم در دوران کارشناسی و کارشناسی ارشد به ویژه آقای مسعود ماهوش که

زحمت تایپ این پایان‌نامه را به عهده داشتند و خانمها مهتاب بقائی، فرخنده سجادی و مریم قائمی.

و خانواده بزرگووارم

پدر ارجمند و مادر مهربانم - دوگوهر پر تلالؤ زندگی‌ام - که وجودم برایشان

همه رنج بود و وجودشان برایم همه مهر، مویشان سپیدی گرفت تا رویم

سپید بماند، آنان که فروغ نگاهشان و گرمی کلامشان سرمایه‌های جاودانی

زندگی‌م هستند. آنان که راستی قامت در شکستگی قامتشان تجلی یافت.

در برابر وجودگرامیشان زانوی ادب بر زمین می‌نهم و با دلی مملو از عشق و

محبت بر دستانشان بوسه می‌زنم...

و خواهران عزیز و برادران خوبم.

امید است که خداوند توفیق و رسیدن به کمال را جبران محبت‌هایشان قرار دهد.

صدیقه میرزائی صالح آبادی

تابستان ۱۳۸۰

کلیه حقوق مادی مترتب بر نتایج
مطالعات، ابتکارات و نوآوریهای ناشی
از تحقیق موضوع پایان نامه متعلق به
دانشگاه صنعتی اصفهان است.

معلمان بزرگوارم:

نهال دانشی را که دیروز در وجودم نهادید به لطف حق و صبر و تلاش
مستمر شما عزیزان امروز به درختی جوان بدل گشته و اولین شاخه اش
به گل نشست است، اینک این اولین شکوفه اش را هرچند ناچیز به شما
تقدیم می کنم.

فهرست مطالب

صفحه	عنوان
هشت	فهرست مطالب
۱	چکیده
فصل اول: مقدمه	
۴	(۱-۱) تاریخچه
۵	(۲-۱) روشهای تجزیه و تحلیل رگرسیون
۵	(۱-۲-۱) رگرسیون خطی ساه
۸	(۲-۲-۱) رگرسیون وزنی
۱۳	(۳-۲-۱) رگرسیون توانمند
۱۵	الف- برآوردگر M
فصل دوم: هموارسازی، مونت کارلو، بوت استرپ	
۱۹	(۱-۲) هموارسازی
۲۰	(۱-۱-۲) هموارسازی نمودار پراکنش
۲۱	(۲-۱-۲) فنون هموارسازی
۲۱	الف- هموارسازی هسته
۲۶	ب- هموارسازی نزدیکترین K همسایه
۲۹	ج- هموارساز خط متحرک
۳۱	(۲-۲) مطالعه مونت کارلو
۳۵	(۳-۲) بوت استرپ
۳۶	(۱-۳-۲) روش بوت استرپ
۳۷	(۲-۳-۲) تابع توزیع تجربی

۳۷ قاعده‌های با جایگذاری (۳-۳-۲)
۳۷ برآورد بوت استرپ خطای استاندارد (۴-۳-۲)
۳۹ فاصله‌های اطمینان بوت استرپ (۵-۳-۲)
۴۰ فاصله اطمینان بوت استرپ-t (۶-۳-۲)
۴۱ فاصله اطمینان بر اساس صدک‌های بوت استرپ (۷-۳-۲)

فصل سوم: روشهای برآورد واریانس خطا در مدل‌های رگرسیونی

۴۵ هموارساز فاصله متحرک (۱-۳)
۴۸ برآوردگر کراندار اثرگذار M با وزن شوئیپ (۲-۳)
۵۱ روش رابینسون (۳-۳)
۵۲ یک مثال (۴-۳)

فصل چهارم: شبیه‌سازی برای ارزیابی روشها

۵۶ شبیه‌سازی (۱-۴)
۵۷ مفهوم شبیه‌سازی (۱-۱-۴)
۶۰ شرایط استفاده از شبیه‌سازی (۲-۱-۴)
۶۰ استفاده از شبیه‌سازی برای بررسی روشهای برآورد (۲-۴)
۷۷ فاصله اطمینان بوت استرپ (۳-۴)

فصل پنجم: نتیجه‌گیری

۸۰ بحث و بررسی نتایج (۱-۵)
۸۲ پیشنهادات (۲-۵)

پیوست

۸۵ پیوست ۱
۸۹ پیوست ۲
۱۱۴ پیوست ۳
۱۶۹ منابع
۱۷۳ چکیده انگلیسی

چکیده:

عدم کارایی روش حداقل مربعات برای برآورد پارامترها، وقتی که خطای مدل دارای توزیعی با دم سنگین باشد و یا خطای مدل دارای واریانس ثابت نباشد، مسأله‌ای آشنا در رگرسیون خطی است.

در چنین حالتی فاصله اطمینان استاندارد پارامترهای مدل، احتمال پوششی کاملاً متفاوت از آنچه راکه باید باشد خواهد داشت. در این پایان‌نامه ابتداء معرفتی سه روش پیشنهادی برای حل این مشکل، با استفاده از تکنیک شبیه‌سازی می‌پردازیم و کارایی آنها را با هم مقایسه می‌نمائیم. نشان داده می‌شود که روش برآورد کراندار اثر گذار M با وزن شوئیپ، دارای بالاترین کارایی است، اگرچه در بعضی موارد ممکن است دو روش دیگر هم نتایج مناسبی را ارائه دهند.

در پایان، برای شیب خط رگرسیونی، یک فاصله اطمینان، با استفاده از تکنیک بوت استرپ و هرکدام از روشهای معرفی شده، به دست می‌آوریم.

فصل اول

مقدمه:

مطالعه روابط بین متغیرها در زمینه‌های زیادی از فعالیتهای علمی متداول شده است. تجزیه و تحلیل رگرسیونی یکی از گسترده‌ترین روشهای آماری است که در بررسی روابط بین متغیرها مورد استفاده قرار می‌گیرد. این روش به علت ظرافت مباحث نظری آن و به دلیل فراهم آوردن یک روش ساده مفهومی برای بررسی رابطه بین متغیرها در علوم مختلف قابل توجه است.

تحلیل رگرسیونی و استفاده از مدل‌های رگرسیونی ممکن است برای اهداف مختلفی از جمله توصیف داده‌ها، برآورد پارامترها، کنترل متغیرها و پیش‌بینی به کار رود. شاید مهمترین کاربرد رگرسیون در این است که بر اساس مشاهدات، متغیر پاسخ را در محدوده داده‌ها پیش‌بینی نماید که این امر در انجام برنامه‌ریزی‌ها می‌تواند بسیار مفید باشد. از طرف دیگر، از لحاظ نظری بخش عمده‌ای از تحلیل‌های رگرسیونی در ارتباط با مدل‌های ریاضی است، مدل‌هایی که تعیین‌کننده شکل و نوع ارتباط بین متغیرها هستند.

برای تجزیه و تحلیل رگرسیونی لازم است فرضیهایی در نظر گرفته می‌شوند. اعتبار نتایج حاصل از رگرسیون، مشروط به صحت این فرضیات است. فرض اساسی که در تجزیه و تحلیل رگرسیون خطی مبنای کار قرار می‌گیرد این است که ارتباط بین متغیر پاسخ Y و متغیرهای مستقل X_1, X_2, \dots, X_p به صورت زیر است:

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon \quad (1-1)$$

که در آن ϵ عامل خطای تصادفی مدل است. در اینجا معمولاً پارامترهای مدل به روش حداقل مربعات برآورد می‌شوند.

اگر علاوه بر برآورد پارامترها و برازش مدل نیاز به استنباطهای آماری درباره پارامترهای مدل باشد، برای سهولت در کار به جز فرض بالا فرضهای زیر نیز باید برقرار باشند:

(۱) خطاهای تصادفی ϵ مستقل از هم باشند.

(۲) خطاهای تصادفی ϵ دارای توزیع نرمال با میانگین صفر و واریانس ثابت σ^2 باشند.

اما در بسیاری از مطالعات عملی مشاهده می‌شود که برخی از فرضیات فوق قابل حصول نیستند. در آن صورت برای رسیدن به فرضهای مذکور از روشهای متفاوتی از جمله انتخاب تبدیلات مناسب (دوجمله‌ای، پواسون، دوجمله‌ای منفی و...) یا انتخاب وزن مناسب برای مشاهدات و در نظر گرفتن رگرسیون وزنی و... استفاده می‌شود (۱). در بسیاری موارد از جمله در بررسی‌های اقتصادی یا مطالعات روانشناسی با عدم تثبیت واریانس خطاهای تصادفی ϵ مواجه می‌شویم که در این صورت می‌توان با در نظر گرفتن وزن برای مشاهدات و استفاده از رگرسیون وزنی واریانس خطا را ثابت نمود.

این تحقیق تلاشی برای شناخت و مقایسه چند روش ابداع شده برای نحوه انتخاب وزن مناسب برای مشاهدات در رگرسیون است به طوری که برآورد پارامترها و استنباطهای آماری درباره آنها از اعتبار لازم برخوردار باشند.

۱.۱) تاریخچه

استفاده از روشهای رگرسیونی از همان موقعی که لژاندر^۱ (۱۸۰۵) و گوس^۲ (۱۸۰۹) داده‌های «مسیر حرکت اشیاء نجومی» را بررسی می‌کردند مورد علاقه محققان بوده است. در دهه‌های بعد، با گسترش روشهای رگرسیونی در علوم اجتماعی متوجه شدند که بعضی از فرضیات مدل (حتی به طور تقریبی) برقرار نیستند. یول^۳ (۱۸۹۹) یکی از اولین کسانی بود که در مقاله خود تحت عنوان «یک بررسی در علل تغییرات در گدائی در انگلستان، مخصوصاً در دود دهه اخیر» با این مشکل مواجه گردید و برای رفع این اشکال در داده‌های موجود تغییراتی ایجاد نمود که سبب تعدیل در فرضیات مدل شد. پس از آن فعالیت‌های زیادی در این زمینه انجام شده و هنوز هم ادامه دارد. اما سابقه تثبیت واریانس خطای تصادفی به سال ۱۹۶۴ برمی‌گردد. در این سال هوبر^۴ با تعریف تابع Ψ - هوبر و در نظر گرفتن وزن مناسبی برای مشاهدات، رگرسیون وزنی را مورد مطالعه قرار داد.^{۳۱}

از سایر مطالعات انجام شده در زمینه ثابت کردن واریانس خطای تصادفی در مدل رگرسیونی می‌توان به مطالعات هوبر (۱۹۷۳)^۴، باکس^۵ و هیل^۶ (۱۹۷۴)^۵، شوئیپ^۷ (۱۹۷۵)^۶، موستلر^۸ و توکی^۹ (۱۹۷۷)^۷، هیل و هلند^{۱۰} (۱۹۷۷)^۸، جابسون^{۱۱} و فولر^{۱۲} (۱۹۸۰)^{۱۰} که با معرفی وزنه‌های متفاوت، برآوردی برای واریانس خطا در نظر گرفته‌اند و نیز مطالعات اندروز^{۱۳} (۱۹۷۴)^{۱۱}، مالوز^{۱۴} (۱۹۷۵)^{۱۲}، بیکل^{۱۵} و داکسام^{۱۶} (۱۹۸۱)^{۱۳}، هوبر (۱۹۸۱)^{۱۴}، کارول^{۱۷} و راپرت^{۱۸} (۱۹۸۲)^{۱۵}، همپل^{۱۹} (۱۹۸۶)^{۱۶}، مارونا^{۲۰} و مورجنتالر^{۲۱} (۱۹۸۶)^{۱۷}، رایبنسون^{۲۲} (۱۹۸۷)^{۱۸}، ویتز^{۲۳} (۱۹۹۲)^{۱۹}، کوهن^{۲۴}، دالال^{۲۵} و توکی (۱۹۹۳)^{۲۰} و ویلکاکس^{۲۶} (۱۹۹۵)^{۲۱} که با معرفی وزنه‌های مختلف، برآورد پایداری برای واریانس خطا ارائه داده‌اند، اشاره نمود.

هدف اصلی این تحقیق مطالعه اهمیت فرض ثابت بودن واریانس خطا در مباحث رگرسیونی و بررسی چند روش مهم معرفی وزن و در نظر گرفتن رگرسیون وزنی در برقراری این فرض، در فصل سوم و نیز شبیه سازی و مقایسه این روشها در فصل چهارم است.

- | | | |
|--------------|---------------|------------------|
| 1- Legender | 2- Gauss | 3- Yule |
| 4- Huber | 5- Box | 6- Hill |
| 7- Schweppe | 8- Mostteller | 9- Tukey |
| 10- Holland | 11- Jobson | 12- Fuller |
| 13- Andrews | 14- Mallows | 15- Bickel |
| 16- Doksom | 17- Carroll | 18- Ruppert |
| 19- Hample | 20- Marrona | 21- Morgenthaler |
| 22- Robinson | 23- Wiens | 24- Cohen |
| 25- Dalal | 26- Wilcox | |

۲.۱) روشهای تجزیه و تحلیل رگرسیون

در این پایان نامه از واژه‌هایی مانند رگرسیون خطی ساده^۱، روش حداقل مربعات^۲، رگرسیون وزنی^۳، برآوردگر رگرسیونی^۴ و ... استفاده شده است. برای سهولت در مطالعه این پایان نامه، در این بخش به اختصار از آنها یاد می‌کنیم.

۱.۲.۱) رگرسیون خطی ساده

ساده‌ترین مدل رگرسیونی، رگرسیون خطی ساده شامل یک متغیر وابسته (پاسخ) و یک متغیر مستقل است. هدف پیدا کردن یک رابطه خطی بین X و Y است. اگر برای یافتن این رابطه، ضریب X با استفاده از حداقل نمودن مربعات خطا، برآورد شود، در این صورت این روش، روش حداقل مربعات در رگرسیون خطی ساده نامیده می‌شود.

فرضهای آماری در رگرسیون خطی ساده عبارتند از:

الف) شرط وجودی: به این معنی که $\mu_{Y|X}$ و $\sigma_{Y|X}^2$ مقادیر متناهی داشته باشند.

ب) شرط استقلال: به این معنی که Y_1, Y_2, \dots, Y_n مستقل از یکدیگر باشند.

ج) شرط خطی بودن: به این معنی که رابطه بین X و Y خطی باشد.

د) شرط همگونی واریانس: به این معنی که واریانس خطای مدل ثابت باشد.

ه) شرط نرمال بودن: به این معنی که Y_i ها دارای توزیع نرمال باشند.

روش کار به صورت زیر است:

اگر زوج مرتب (x_i, y_i) $i = 1, 2, \dots, n$ از رابطه خطی

$$y = X\beta + \varepsilon \tag{۲-۱}$$

$$\underline{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad \underline{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix} \quad \text{و} \quad X = \begin{pmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_n \end{pmatrix} \quad \text{و} \quad \underline{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$$

که در آن

آنگاه با می‌نیمم کردن تابع

$$Q(\beta) = \varepsilon^T \varepsilon = (\underline{y} - X\underline{\beta})^T (\underline{y} - X\underline{\beta}) \tag{۳-۱}$$

نسبت به β ، برآورد حداقل مربعات برای β به صورت زیر بدست می‌آید:

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

این برآورد دارای توزیع نرمال با میانگین β و ماتریس واریانس-کواریانس $(X^T X)^{-1} \sigma^2$ است.

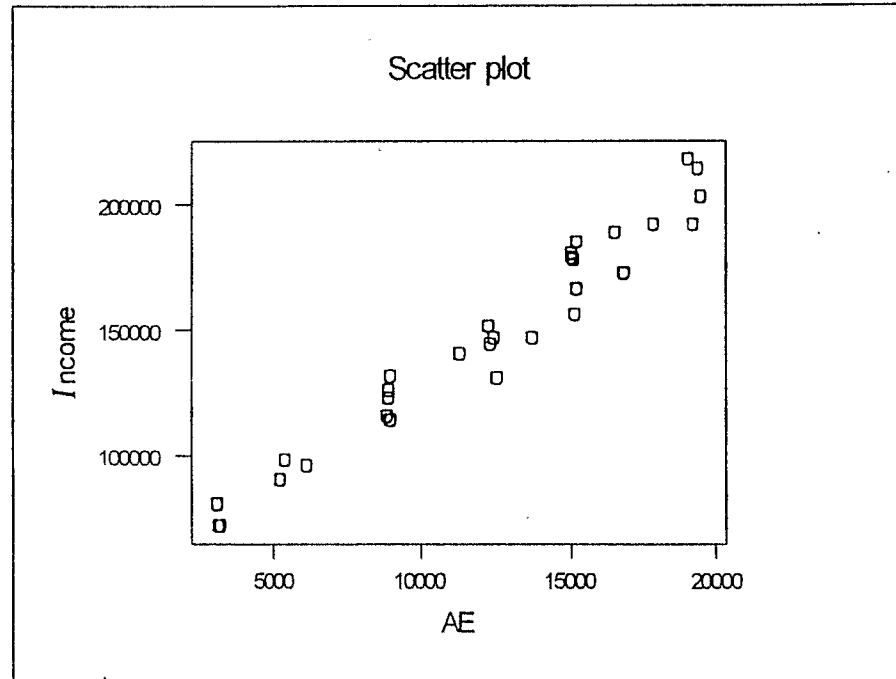
مثال ۱-۱: مدیر رستورانی معتقد است که درآمد ماهانه از فروش غذا و هزینه تبلیغات برای معرفی آنها، رابطه‌ای مستقیم دارند. برای بررسی این مسأله اطلاعات مربوط به درآمد ماهانه حاصل از فروش غذای ۳۰ رستوران و نیز هزینه آگهی‌های سالانه را جمع‌آوری کرده است (۲۰). داده‌های جمع‌آوری شده در جدول ۱-۱ آورده شده‌اند:

جدول (۱-۱) داده‌های مربوط به درآمد ماهانه و هزینه آگهی‌های سالانه مثال (۱-۱) (۲۰)

مشاهده	درآمد	هزینه‌های تبلیغاتی
۱	۸۱۴۶۴	۳۰۰۰
۲	۷۲۶۶۱	۳۱۵۰
۳	۷۲۳۳۴	۳۰۸۵
۴	۹۰۷۴۳	۵۲۲۵
۵	۹۸۵۸۸	۵۳۵۰
۶	۹۶۵۰۷	۶۰۹۰
۷	۱۲۶۵۷۴	۸۹۲۵
۸	۱۱۴۱۳۳	۹۰۱۵
۹	۱۱۵۸۱۴	۸۸۸۵
۱۰	۱۲۳۱۸۱	۸۹۵۰
۱۱	۱۳۱۴۳۴	۹۰۰۰
۱۲	۱۴۰۵۶۴	۱۱۳۴۵
۱۳	۱۵۱۳۵۲	۱۲۲۷۵
۱۴	۱۴۶۹۲۶	۱۲۴۰۰
۱۵	۱۳۰۹۶۳	۱۲۵۲۵

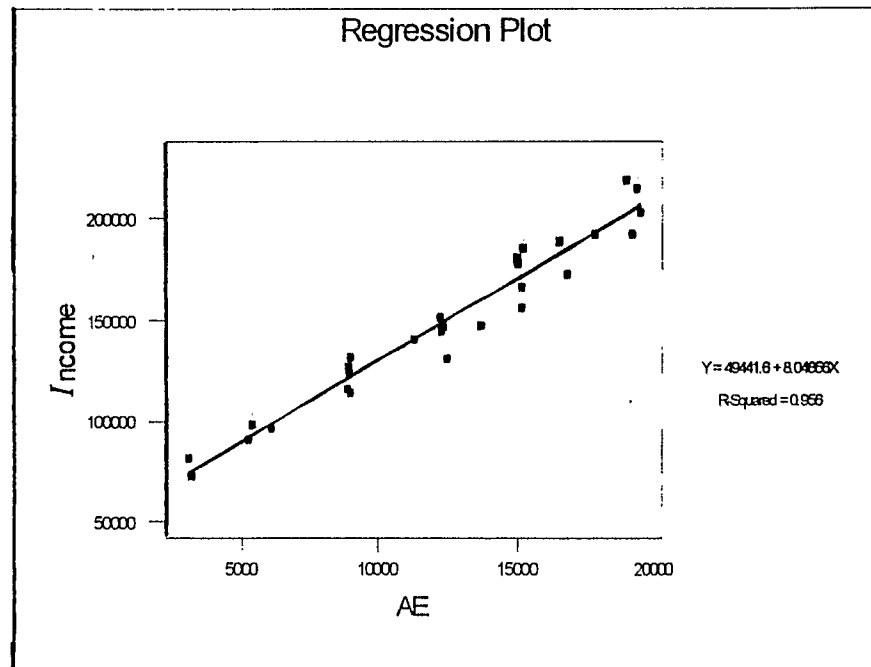
مشاهده	درآمد	هزینه‌های تبلیغاتی
۱۶	۱۴۴۶۳۰	۱۲۳۱۰
۱۷	۱۴۷۰۴۱	۱۳۷۰۰
۱۸	۱۷۹۰۲۱	۱۵۰۰۰
۱۹	۱۶۶۲۰۰	۱۵۱۷۵
۲۰	۱۸۰۷۳۲	۱۴۹۹۵
۲۱	۱۷۸۱۸۷	۱۵۰۵۰
۲۲	۱۸۵۳۰۴	۱۵۲۰۰
۲۳	۱۵۵۹۳۱	۱۵۱۵۰
۲۴	۱۷۲۵۷۹	۱۶۸۰۰
۲۵	۱۸۸۸۵۱	۱۶۵۰۰
۲۶	۱۹۲۴۲۴	۱۷۸۳۰
۲۷	۲۰۳۱۱۲	۱۹۵۰۰
۲۸	۱۹۲۴۸۲	۱۹۲۰۰
۲۹	۲۱۸۷۱۵	۱۹۰۰۰
۳۰	۲۱۴۳۱۷	۱۹۳۵۰

در این مثال درآمد ماهانه، متغیر پاسخ و هزینه آگهی ها، متغیر مستقل است. قبل از هر چیز نمودار پراکنش این دو متغیر را رسم نموده و مشاهده می‌کنیم که رابطه‌ای مستقیم بین این دو وجود دارد. (شکل ۱-۱)



شکل ۱-۱) نمودار پراکنش درآمد در مقابل هزینه آگهی، مثال (۱-۱)

با توجه به داده‌ها، مدل برازش شده عبارت است از: $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X \equiv Income = 49441.6 + 8.04668X$ که خط برازش شده به داده‌ها در شکل (۲-۱) نمایش داده شده است.



شکل ۲-۱) نمودار خط برازش شده به داده‌ها، مثال (۱-۱)