

استاد راهنما :
آقای دکتر خسرو ملکی

مقدمه ای بر اقتصادسنجی

AN INTRODUCTION TO ECONOMETRICS

A. A. Walters

آ. آ. والترز

ترجمه :
خسرو ملکی

۱۰۶۵۴

مقدمه

انتخاب کتاب "مقدمه ای بر اقتصاد سنجی" با توجه به سوابق مترجم در امر ترجمه و بر اساس راهنمایی های ارزشمند جناب آقای دکتر ملاح صورت گرفت. ترجمه مزبور که از صفحات ۸۰ تا ۱۴۴۰ کتاب است جهت ارائه بعنوان پایان نامه دوره لیسانس سیاسی دانشکده اقتصاد و علوم سیاسی دانشگاه ملی می باشد. امید است این ترجمه که در زمینه جدیدی از اقتصاد است مورد استفاده پژوهشگران قرار بگیرد.

در اینجا لازم میدانم مجدداً از کلیه استادان که همواره ارشاد کننده اینجانب بوده اند بویژه آقای دکتر ملاح برای کمکهاییکه در ترجمه این پایان نامه نموده اند تشکر نماید.

خسرو ملک

۱۹۶۵

فهرست مندرجات

صفحه

مقدمه رگرسیون و همبستگی

فصل چهارم - روابط بین متغیرها

۱		بخش ۱ رگرسیون و روابط نظری
۹		بخش ۲ کمترین مجذورات
۱۵		بخش ۳ واریانس مانده ها و ضریب همبستگی
۱۷		بخش ۴ همبستگی های غیرخطی
۲۰		بخش ۵ برآورد مقدار آلتا
۲۰		بخش ۶ محاسبه رگرسیون
۲۲		بخش ۷ رگرسیون و روابط علت و معلولی
۲۶		بخش ۸ خطا در متغیرها
۳۲		سوالات برای بحث

فصل پنجم - رگرسیون متعدد (مرکب)

۳۹		بخش ۱ ضرائب رگرسیون
۴۴		بخش ۲ همبستگی چند تائی
۴۶	R^2	بخش ۳ ضریب همبستگی متعدد
۴۸		بخش ۴ ضریب همبستگی نسبی
۵۵		بخش ۵ نتیجه گیری
۶۰		سوالات برای بحث

رگرسیون^۱ و همبستگی^۲

در قسمت دوم ابزار اصلی رگرسیون معرفی و مورد بحث قرار داده شده است. قسمت عمده اقتصاد مربوط میشود به پیش‌بینی تأثیر متغیرها بر روی یکدیگر. رگرسیون و همبستگی روشهایی هستند که برای تجزیه و تحلیل مسائل وابستگی و فتنیکه اثرات خطی میانند بکار گرفته میشوند. فصل عمده ایسمن قسمت فصل چهارم است که بعنوان زیربنائی است که مابقی این کتاب بر روی آن بنا شده است. تبحر یا حداقل آشنائی نزدیک به مدلهای ابتدائی رگرسیون که دارای دو متغیر میباشند قدم اول و لازمی است برای آشنائی به مسائل رگرسیون. ولی هدف اصلی فصل چهارم تفسیر متضاد و بررگرسیون و همبستگی است. هدف طرح عمل رگرسیون — نمایان نمودن افلا — چند نمونه از — استثنائات میباشد. فصل پنجم نمونه دو متغیری را کاملتر کرده و به نمونه های سه یا چند متغیری میرد از د. در فصل پنجم چندان افکار جدیدی وارد نشده (با مقایسه با فصل چهارم) ولی متأسفانه لزوماً "از جبهه میزان متناهی — استفاده گردیده است. با وجود این، از آنجائیکه کار عملی در رگرسیون شامل حداقل رگرسیون نسبی و متعدد میگردد درک کامل جهانی اصلی لازم است. فصل ششم در مورد نمونه گیری خصوصیات مدلهای رگرسیون و همبستگی بحث می نماید و لزومی ندارد که برای درک مابقی کتاب این فصل خوانده شود.

Regression — ۱

Correlation — ۲

رگرسیون، روابط نظری؛ اغلب دانش‌های علمی را میتوان بصورت برتراری رابطه یا نسبتی بین دو یا چند مقدار تنظیم و توجیه کرد. تحقیق بمنظور دستیابی به برقراری نظم کلی بین متغیرها یکی از گرایش‌های عمده علوم فیزیکی و اجتماعی را تشکیل میدهد و روش‌هایی که از طریق آنها میتوان چنین روابطی را کشف و ترسیم نمود ابزار اصلی "علم اقتصاد سنجی" را تشکیل میدهند. در این فصل کاربرد این ابزار بطور نسبتاً مبسوطی مورد بحث قرار میگیرد.

بعنوان نمونه مسئله قدیمی اقتصاد کلاسیک جدید (نئوکلاسیک) را در نظر بگیرید یعنی شکل و شیب منحنی تقاضا. نظریه تقاضای مصرف‌کننده باین پیش‌بینی واکنشیت مصرف‌کننده مستقل (یا یک خانواده) در قبال تغییر قیمت بازار مربوط میگردد. چنین فرض می‌کنیم که "سلیقه‌ها" ثابت است، درآمد واقعی مصرف‌کننده تغییر نمی‌کند و اینکه کالای مورد نظر چنان جزء کوچکی از کل بودجه مصرف‌کننده را تشکیل میدهد که تغییرات قیمت آن مصرف‌کننده را از لحاظ مادی بطور چشم‌گیری در وضع بهتریابد تری قرار نمیدهد. وانگهی فرض بر این است که مصرف‌کننده بطور یکنواختی (Consistent) عمل میکند که البته تعریف "یکنواختی" بر حسب اینکه نظریه ما تا چه حد ممکن است پیچیده باشد، تغییر میکند. ولی صرف نظر از ریزه کاریهای آن، "قانون تقاضا" بر اساس این نظریه بطور کم و بیش روشنی چنین مطرح میگردد که: اگر قیمت کالائی کاهش یابد مقدار بیشتری از آن کالا خریداری میشود. این قانون که بیشتر برای یک فرد مصرف‌کننده در نظر گرفته شده است میتوان بسهولة آن را در شرایط بازار نیز تعمیم داد. برای رسیدن به منحنی تقاضای تنها کاریکه باید انجام داد این است که مقادیر کالائی که در قیمت معینی توسط افراد خریداری شده است با هم جمع کنیم. چون طبق پیش‌بینی قانون فوق الذکر با کاهش قیمت، افراد مقدار خرید خود را افزایش میدهند پس میتوان چنین نتیجه‌گیری کرد که بطور کلی حجم خریدها افزایش خواهد یافت. پس میتوان گفت که قانون فوق الذکر هم در مورد افراد صدق میکند و هم در مورد بازار.

این چنین فرض می‌کنیم که کالای مورد نظر توسط یک انحصارگر عرضه میشود. این شخص قیمت کالای را خود تعیین میکند و سپس بازار را در مقابل این قیمت مورد مشاهده قرار میدهد. بدیهی است که انحصارگر مایل است که راجع به ماهیت تقاضا اطلاعاتی را کسب کند چون علاقمند با انتخاب

قیمتی است که بیشترین سود را برای شما حاصل کند. حال فرض می‌کنیم که این انحصارگر "آزمایش‌های" خود را به مرحله اجرا درآورد. باین معنی که هرماه قیمت را در سطوح مختلفی تعیین کند و سپس مقدار کالایی که در آن قیمت بفروش رفته است مورد مشاهده قرار دهد. اگر این مشاهده تغییر قیمت و فروش کالای را روی منحنی تقاضا رسم کنیم منحنی بدست آمده دارای شیب نزولی است که قانون یا نظریه فوق الذکر آن را پیش‌بینی کرده است.

ولی حتی در این شرایط آزمایشی مشاهدات ما دقیقاً با نظریه فوق الذکر تطبیق نمی‌کند. "سلیقه‌ها" ممکن است از یک زمان نسبت بزمان دیگر تغییر کند و اینکه قیمت کالاهای جانشین و مکمل نیز تغییر کنند و یا درآمد‌ها افزایش یافته باشند تمام فرضیات مربوط به نظریه فوق الذکر ممکن است هنگام آزمایش نقض شوند. بنابراین درآمد مورد مشاهده لااقل دارای دو جز اصلی است: پیش‌بینی نظری و "خطاهای" تجربی که شامل خطاهای اندازه‌گیری نیز می‌شود. ولی کسی نمی‌تواند ادعا کند (یا لااقل در اقتصاد چنین ادعا کند) که اگر خطاهای تجربی وجود نداشت اطلاعات داده شده نمی‌توانست همواره با نظریه فوق الذکر دقیقاً مطابقت پیدا کند.

حتی در شرایط آزمایشی کامل نظریه فوق ممکن است تنها یک پیش‌بینی ناقص از اطلاعات فرض شده را بعمل آورد. از این روی همیشه یک جز سوم نیز در نتایج حاصله از آزمایش وجود دارد و آن عبارت است از: خطاهای ناشی از اینکه نظریه خود به تنهایی نمایشگر یک توصیف تقریبی از واقعیت‌ها است.

خطاهای موجود در آزمایش چنان حائز اهمیت اند که همواره بایستی آنها را مورد توجه دقیق قرار داد. سودمندی این نظریه به مقدار زیادی این امید را در ما منمیزند که اهمیت خطاهای نسبت به قدرت توجیه کننده سیستماتیک نظریه ناچیز است. در واقع روش تجربی نه تنها بایستی روشنگر نظریه مورد بحث باشد — یعنی رابطه سیستماتیک بین متغیرها — بلکه بایستی کیفیت خطاهای آن نیز توجیه کند.

مثال فرضی آن انحصارگری را در نظر بگیرید که می‌خواست "منحنی تقاضای" خود را کشف کند. فرض کنیم که نظریه مورد بحث پیش‌بینی می‌کند که رابطه بین قیمت و مقدار فروخته شده رابطه خطی است که دارای شیب منفی است. ما می‌توانیم "نظریه" یا چند نظریه رگرسیون را بطور زیر نشان دهیم:

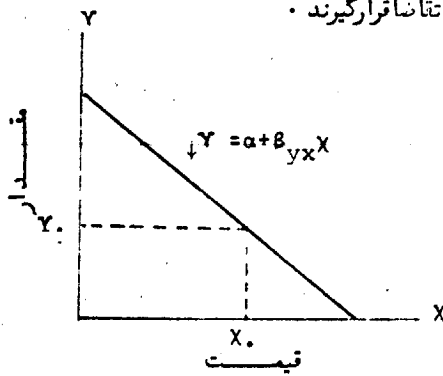
+ توجه داشته باشید که ما نظریه تقاضا را خیلی دقیق تراز حالت معمولی بکار می‌بریم. در اقتصاد سنجی معمولاً ضروری ترویا مطلوب تر است که از فرضیات محدودتری استفاده کنیم تا آنجا که بطور ضمنی توسط صرف نظریه بیان می‌شود.

(۱-۱) $\gamma = \alpha + \beta_{yx} X$ برای مقدار معینی از X که در آن $\alpha > 0$ ، $\beta_{yx} < 0$ است . در معادله فوق γ مقدار کالای خریداری شده؛ X قیمت α ، β_{yx} نیز مقدار ثابتی هستند . ترتیب ریزنویس β_{yx} yx نشانه آن است که γ متغیر تابع و X متغیر مستقل است یعنی با معلوم بودن مقدار X ، میتوان مقدار γ را پیش بینی کرد . حال فرض میکنیم که نظریه ماتریح میکند که بین رابطه خطی وجود دارد ولی صرف نظر از پیش بینی اینکه خط دارای شیب منفی است، نظریه مقدار عددی این شیب و یا محل دقیق آن را با نشان نمیدهد . فرضیه ماصرفاً مبتنی بر این است که یک رابطه خطی منفی وجود دارد ، ولی این کار آزمایی است که باید کشف نماید که آیا ۱- اطلاعاتی وجود دارند که با این فرضیه تطبیق نمایند ؟ و یا دقیق تر اینکه ۲- محل و شیب خط تقاضا را برآورد - نمایند . ترسیم هندسی این نظریه در شکل (۱-۱) نشان داده شده است . رابطه تقاضا به شرح زیر تفسیر میشود :

" با معلوم بودن قیمت X_0 سپس مقدار فروش رفته γ_0 خواهد بود "

این رابطه عبارت است از پیش بینی مقدار فروخته شده روی محور عمودی و بر اساس قیمت داده شده روی محور افقی . مقدار کالای فروخته شده متغیر تابع است که بر اساس اطلاع از قیمت که متغیر مستقل است ، پیش بینی میشود . (خوانندگان توجه خواهند داشت که ما جای محور ها را از روش معمولی نشان دادن منحنی های تقاضا که در کتب کشورهای آنکلو ساکسون بکار برده میشود ، عوض کرده ایم . هیچگونه تغییر عمده ای در اینجا داده نشده است و ترسیم چنین نموداری با موازین پذیرفته شده در علوم و آمار بهتر مطابقت پیدا میکند .)

تنها در یک آزمایش کامل وین نظریه کامل است که تمام نتایج مشاهده شده روی خط قرار میگیرند . درین دنیای واقعی بهترین چیزی که میتوان انتظار داشت این است که قیمت ها و مقدارها مشاهده شده در نزدیکی خط تقاضا قرار گیرند .

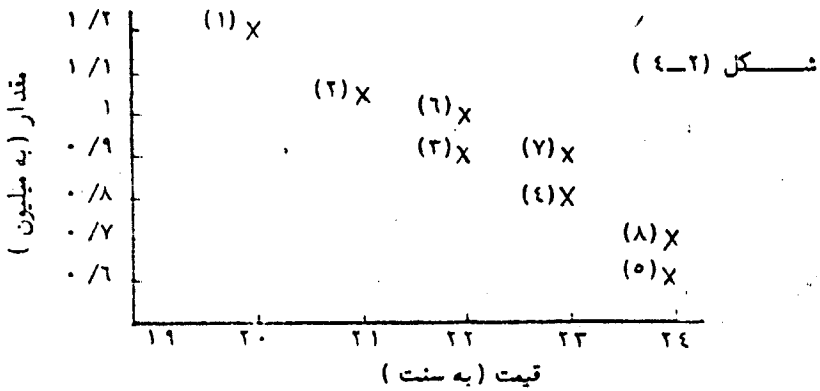


شکل (۱-۱)

ما میتوانیم مقادیر خریداری شده در بازار را همانطوریکه با قیمت های مختلفی که توسط یک انحصارگر تعیین میشود رسم کنیم. فرض میکنیم که مقادیر زیر مورد مشاهده قرار گرفته است.

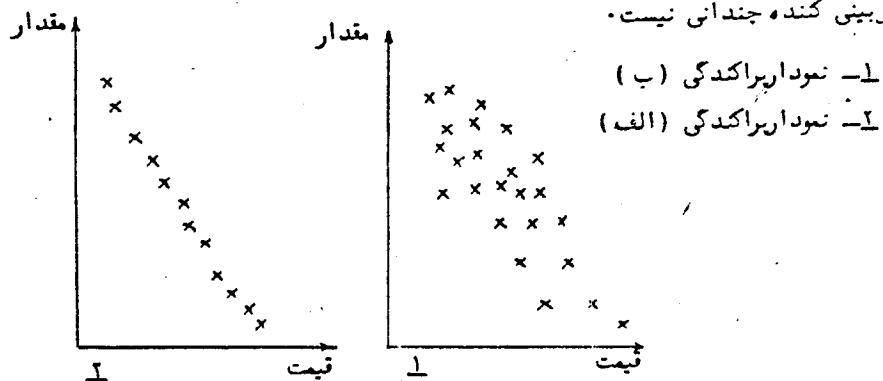
ماه	قیمت (به سنت)	مقدار
۱	۲۰	۱/۲۰۰/۰۰۰
۲	۲۱	۱/۰۰۰/۰۰۰
۳	۲۲	۹۰۰/۰۰۰
۴	۲۳	۸۰۰/۰۰۰
۵	۲۴	۶۰۰/۰۰۰
۶	۲۲	۱/۰۰۰/۰۰۰
۷	۲۳	۹۰۰/۰۰۰
۸	۲۴	۷۰۰/۰۰۰

اگر اطلاعات فوق را روی منحنی ایکه دارای دو متغیر است رسم کنیم، (نموداریکه اصطلاحاً نمودار پراکندگی نامیده میشود) نمودار زیر بدست میآید.

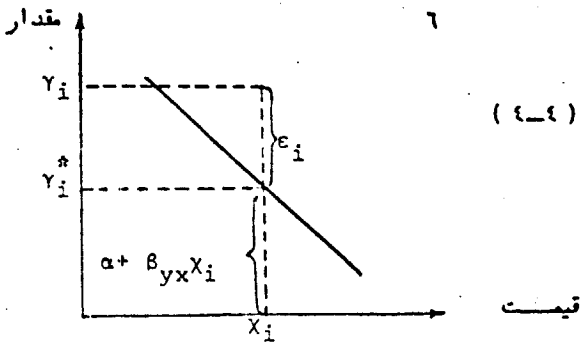


باید توجه داشته باشید که این مشاهدات دقیقاً روی خط قرار نمیگیرند. خط‌هایی را که خود را بر میان نتایج بدست آمده پیدا کرده اند و بایستی مسئله را با توجه به خط‌های موجود مورد بررسی قرار داد. در اینصورت میتوان رابطه خطی فرضیه و در صورت امکان محل خط را تعیین و کشف کرد.

با د نظر گرفتن مثال داده شده بدیهی است که نوعی رابطه منفی بین متغیرها وجود دارد . ولی همچنین مسلم است که خطاهای مربوط به فرضیه و آزمایش زیاد هم ناچیز و قابل چشم پوشی نیست . مشاهدات در امتداد یک خط مستقیم نشان داده نشده بلکه دارای یک پراکندگی سیگاری مانند است که جهت آن از شمال غربی به جنوب شرقی امتداد دارد . هرچه این مسیر سیگارمانند نازک تر باشد ، قدرت پیش بینی نظریه بیشتر و احتمال وجود خطاها کمتر خواهد شد . ولی اگر ایستگاه سیگارمانند و صدای پراکندگی بیشتر و خیلی ضخیم باشد ، خطاهای موجود نسبت به ارزش پیش بینی کننده نظریه ، قابل توجه خواهد بود . چنین وضعی در نمودارهای زیر نشان داده شده است : نمودار (الف) ارقامی را نشان میدهد که دقیقاً با شرایط نظریه (ثوری) مطابقت دارد ، در حالیکه نمودار (ب) موردی را نشان میدهد که خطاها بسیار قابل ملاحظه هستند و نظریه دارای قدرت پیش بینی کننده چندانی نیست .



برای بررسی این اطلاعات بایستی نظریه خطی ساده را دوباره طرح ریزی نمود بطوریکه این دفعه دقیقاً "خطاها" را شامل شود . به عبارت دیگر این بار نظریه بایستی بشکل آماری طرح ریزی شود بطوریکه بتواند اطلاعات داده شده را دقیقاً توصیف نماید . این X_i و Y_i را که به ترتیب نماینده قیمت و مقدار است برای ماه i ام (ith) می نویسیم . بدیهی است که نمیتوانیم به سهولت X_i و Y_i را در معادله (۱-۴) جایگزین نماییم چون Y_i تحت تأثیر خطاها قرار میگیرد . اینک (ϵ_i) را که نماینده خطاها است در Y_i می نویسیم یعنی آنکه "خطاها" در مدل مورد نظر در طی ماه i ام (ith) در متغیر تابع نشان داده میشود . ما در معادله (۲-۴) $Y_i = \alpha + \beta_{yx} X_i + \epsilon_i$. تصریح میکنیم که خطا اضافه شده است . اضافه کردن خطا در معادله در نمودار زیر نشان داده شده است .



شکل (۴-۴)

با معلوم بودن مقدار X_i (قیمت) در ماه i مقدار X_i را مشاهده می‌کنیم. این مقدار ممکن است بدو جز "تجزیه شود". ابتدا آنکه مقدار Y از طریق نظریه پیش‌بینی می‌شود که این صرفاً مقدار Y است که از طریق رابطه خطی پیش‌بینی می‌شود. در این صورت این مقدار را Y_i^* مینامیم. بدین ترتیب خواهیم داشت:

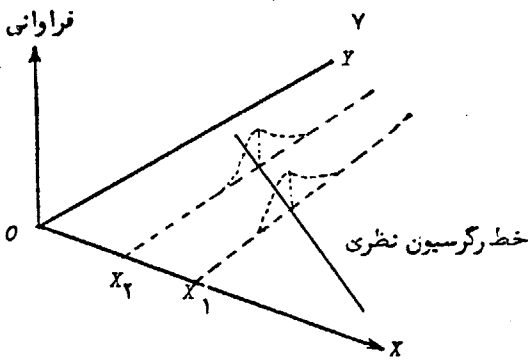
$$Y_i^* = \alpha + \beta_{yx} X_i$$

$$Y_i = Y_i^* + \epsilon_i$$

یا عبارت دیگر، (یا اخلال) خطا + مقدار پیش‌بینی شده = مقدار مشاهده شده

اهمیت قاطع خطا چنان است که گاهی عنوان کم اهمیت تر "اخلال" بآن خطاب می‌شود. (بنابراین از این بی‌مذکمه "خطا" با عبارت "اخلال مشاهده شده" توصیف می‌شود.)

مشکل این است که یک نام مناسب برای "اخلال" پیدا شود. در اقتصاد سنجی از عبارت "اخلال" ها برای نشان دادن اثر عواملی که از نظریه صرف شده است، استفاده می‌شود. در استعمال نظریه تقاضا، می‌توان بعضی از علل اختلاف ها از مدل را تشخیص داد. بعنوان مثال شرایط جوی، جنک و تغییرات سلیقه تعدادی از دلایل مهم برای اختلاف بین پیش‌بینی و واقعیت را تشکیل می‌دهد. موضوع واقعاً مهم است که توزیع اخلال ها از چه قواعدی متابعت می‌کنند؟ روش عادی، تشخیص این مطلب است که اخلال ها مانند متغیرهای اتفاقی "عمل می‌کنند". برای اینکه دقیق باشیم، اخلال ممکن است بعنوان متغیری که بطور نرمال توزیع شده است و دارای میانگین صفر است، توصیف می‌شود. در هر مورد یک انحصارگر قیمت X_i را انتخاب می‌کند ما علاوه برداشتن یک قسمت خطی سیستماتیک، دارای اخلالی (Disturbance) هستیم که از یک توزیع نرمال جمعیت بدست آمده است. بنابراین بطور نموداری می‌توان چنین تصور کرد که توزیع نرمال اخلال در جهت محور Y دور از خط رگرسیون قرار می‌گیرد. اثر مرادونی (Frequency) راروی به محور سوم اندازه گیر کنیم، شکل مدل مسا بصورت شکل (۴-۵) خواهد شد.



تغییر رسمی مدل این است که وقتی انحصارگریت قیمت را تعیین میکند فقط یک بازده وجود دارد که بطور سیستماتیک توسط نظریه پیرامونی میشود که آن مقدار γ^* است. ولی این مقدار واقعی بدست آمده γ برای یک قیمت معین است. ما میتوانیم این موضوع را بطور تجربی تفسیر کنیم.

انحصارگریت γ سنت رامه پس از ماه تعیین میکند و در هر مورد مقدار فروخته شده را ثبت میکند. نظریه مورد بحث بیان میکند که متوسط مقدار فروخته شده طبق قانون زیر بدست میآید:

$$\gamma^* = \alpha + \beta_{yx} \quad (22)$$

این مقدار واقعی نظری α β_{yx} را به ترتیب ۳۰ و ۰/۵ — ثبت میکنیم.

$$\gamma^* = 20 - 0/5 = 19$$

از این روی باقیمت γ سنت متوسط مقدار فروخته شده γ^* است. ولی مقدار فروخته شده برای هر ماه معین ممکن است بعلاوه اختلال (Error) کم و بیش از این رقم ۱۹ تجاوز کند. فروش واقعی یا اضافه کردن تعدادی (e) مقدار متوسط که از توزیع عادی نرمال استخراج شده است، بدست میآید. در هر ماه که γ سنت منظور میشود، تنها چیزی که تغییر میکند استخراج از توزیع عادی نرمال است. گوئی هر ماه تعدادی به ۱ میلیون اضافه شده که بطور اتفاقی از کلامی بیرون کشیده ایم که در آن تکه — هائی گذاشته ایم که نمایشگر فراوانی های نسبی تعدادی است که از توزیع نرمال بدست آمده است. در اینجا باید متکی کنیم و به تحقیق بپردازیم که چرا "خطها و معلول های حذف شده از — مدل" باید از طریق یک متغیر اتفاقی نشان داده شود. بنظر فوق العاده مصنوعی خواهد آمد اگر "حذفیات و خطها" را طوری بحساب آوریم که گوئی از طریق یک قانون تصادفی بوجود آمده اند. این قوانین برای توصیف نتایج آزمایش های تکرارشدنی ای بوجود آمده اند که هر کسی میتواند در صورت امکان مالی، آن آزمایش ها را انجام دهد. شواهد کافی نتوانسته است که روش "خط یا شیر" وارد کند — ولی این امر با تعیین قیمت آزمایش توسط یک انحصارگر بخصوص، یکسان نیست. تنها پاسخ — که پاسخ نامطلوبی نیز هست — این است که این تسمیه اصطلاح اختلال معمولاً

بهترین تسمیه ای است که در حال حاضر وجود دارد. این در واقع شامل چیزی جز این نیست که نامی به جهالت خود داده باشیم و برای آن یک شکل مقداری وسیع در نظر بگیریم.

تسمیه اصطلاح اخلاص اغلب بصورت متغیری که بطور نرمال و اتفاقی توزیع شده است بیان میشود، ولی ضروری نیست. گویانکه اغلب مناسب است که بطور نرمال توزیع شده باشد. در واقع تنها باین احتیاج است که بر حسب یک قانون احتمالی خاص توزیع شود و مهم تر از هر چیز دیگر قواعد توزیع احتمالی نباید به مقدار X بستگی داشته باشد.

صرف نظر از اینکه مقدار X تا چه حد تعیین شود، توزیع اِپسِلین ها (ϵ_s) نباید تحت تاثیر قرار گیرد. در یک مفهوم احتمالی توزیع اِپسِلین ها (ϵ_s) بایستی مستقل از یکسرها (X^s) باشد. این فرضیه حساس به تفصیل در فصل های بعدی مورد بررسی قرار خواهد گرفت ولی در این مرحله ضروری است که به روش تعقلی را ارائه دهیم. برعکس تصور کنید که توزیع اِپسِلین ها (ϵ_s) در واقع به X بستگی دارد. مثلاً فرض کنید که همانطور که X افزایش میابد، مقدار (ϵ) احتمالاً بیشتر مثبت است تا منفی. پس بدیهی است که (ϵ) شامل یک جزئی است که بطور سیستماتیک بمقدار X انتخاب شده بستگی دارد. این جزء از (ϵ) با اطلاع از مقدار X قابل پیش بینی است بنابراین دیگر جزئی از بی اطلاعی ما را تشکیل نمیدهد و نتیجه ظهور این اصل میشود که توزیع اِپسِلین ها (ϵ_s) بایستی مستقل از مقدار X باشد.

قبل از آنکه به تجزیه و تحلیل این مطلب بپردازیم، بهتر است نوعی آسان سازی را ارائه دهیم. چون بایستی دو مقدار ثابت α و β را اندازه گیریم، بهتر است این کار را یکی یکی انجام دهیم باین معنی که اول (β) را اندازه میگیریم و سپس به سنجیدن (α) میپردازیم. بدین ترتیب در حصل مسئله ابتدا باین مسئله (بطور موقتی) مواجه میشویم که چگونه (α) را از معادله حذف کنیم. مقدار (α) محل خطر را تعیین میکند — با بالا و پایین بردن آن در روی محور Y . بنابراین معقول است که موقعیت خطر را همانطور که از میان "میانگین های حسابی نمونه مشاهدهات Y و X عبور میکند" تعریف کنیم — و بعداً مشاهده خواهیم کرد که این امر در عمل چگونه انجام خواهد شد. سپس اگر محورهای X و Y را جابجا کنیم بنحویکه از میان میانگین های X و Y عبور کنند، دو محور و خند در منطقه میانگین ها بایکدیگر تلاقی میکنند.

اگر دارای مقیاس های زیر باشیم

$$\begin{aligned} \varepsilon_i &= \varepsilon_i - \varepsilon_0 & \text{که در آن} & \quad \varepsilon = (i/n) \sum \varepsilon_i \\ \pi_i &= \chi_i - \chi_0 & \text{که در آن} & \quad \chi = (i/n) \sum \chi_i \\ y_i &= \gamma_i - Y_0 & \text{که در آن} & \quad Y = (i/n) \sum \gamma_i \end{aligned}$$

که در آن حروف $\varepsilon_i, y_i, \pi_i$ نشان دهنده "انحرافات میانگین است"، بخش جدا شده (Intercept) ممکن است صفر تلقی گردد — چون خط از میانگین مبداء مختصات (جدید) —

عبور میکند. بنابراین ابتداری پیدا کردن شیب β_{yx} در معادله تکیه میکنیم یعنی معادله —
 $y_i = \beta_{yx} \chi_i + \varepsilon_i$ و س از حل این معادله به بررسی برآورد مقدار (α) میرد ازیم.

کمترین مجذورات: (Least Squares)

با توصیف چنین مدلی اینک به بررسی کشف روش هایی برای برآورد β_{yx} (و بعداً α) از یک نمونه مشاهدات میرد ازیم. فرض بر این است که متغیرها از روش مدل یا الگوی مورد نظر متابعت میکنند (ولی توجه کنید که بخاطر سهولت کار (ε_i) را بدور از انحرافات از میانگین نمونه بجای (ε_i) مینویسیم) (۳-۱)
 $y_i = \beta_{yx} \chi_i + \varepsilon_i$ در این معادله تنها مشکل ما این است که بتوانیم برآورد مقدار β_{yx} را از نمونه مشاهدات (π_i) (قیمت) و (y_i) (مقدار) بدست آوریم. بعلاوه چنین فرض میکنیم که π_i و y_i بدون هرگونه خطائی دقیقاً مشاهده و اندازه گیری میشوند. بر اساس نمونه π_i و y_i باید چنین نتیجه گیری کنیم که بهترین برآورد β_{yx} چیست.

روش های مختلفی برای مطابقت دادن خط رگرسیون با چنین اطلاعاتی وجود دارد. ساده ترین تکنیک این است که اغلب خطی را رسم کنیم که بنظر بیننده به بهترین وجهی نمایشگر رابطه خطی اطلاعات داده شده است. گوا اینکه از این روش بخوبی یاد نشده است ولی این روش مطابقت دادن خط با اعداد از طریق چشم اغلب برای مقاصد متعددی میتواند مفید واقع شود. تنها مشکل اصلی این روش این است که آن روشی عینی نیست — چون وجود بسیاری از مشاهده کنندگان و بسیاری از خطوط امکان پذیر است. وحتى هنوز کاملاً روشن نشده است که کسی بتواند خواص این روش را به بهترین ضرز ممکنه مورد بهره برداری قرار دهد. روش دیگری این است که تنها بالاترین و پایین ترین نقاط π_i و یا شاید بالاترین و پایین ترین نقاط y_i را بهم وصل کنیم. این میتواند یک روش نسبتاً معقولانه ای باشد ولی از آنجائیکه شالوده این مدل را صریحاً یکار نمی برد، بدلیل غیر تجربی ناپایده گرفتن کلیه

مشاهدات در فاصله بین دو منتهایه افراط کارانه بنظر خواهد آمد . روشی که بیشتر مورد تحسین آمارگران قرار گرفته است "روش کترین مجذورات" است . بمعادله زیر توجه کنید :

$$y_i = \beta_{yx} x_i + \epsilon_i$$

که در این مورد ما دارای تعداد n مشاهدات مستقل هستیم x_1, x_2, \dots, x_n و y_1, y_2, \dots, y_n

که در آن میانگین نمونه \bar{y}, \bar{x} صفر است . حال اگر تمام معادله را در x_i ضرب کنیم خواهیم داشت :

$$(i=1, 2, \dots, n) \quad y_i x_i = \beta_{yx} x_i^2 + \epsilon_i x_i \quad (4-1)$$

این روش تعداد (n) معادله $(i=1, 2, \dots, n)$ از نوع معادله (4-1) را در اختیار ما قرار خواهد داد . اینتیت معادله کلی را از جمع معادله (4-1) برای تعداد (n) مشاهدات درست میکنیم ، که حاصل جمع آن معادله زیر میشود :

$$\sum_{i=1}^n y_i x_i = \beta_{yx} \sum_{i=1}^n x_i^2 + \sum_{i=1}^n \epsilon_i x_i \quad (4-5)$$

مقدار $\sum y_i x_i$ حاصل جمع x و y است . وقتیکه معادله فوق بر تعداد مشاهدات (n) تقسیم شود ، معادله "کواریانس" (همپراش) نامیده میشود یعنی کواریانس $(x, y) = (1/n) \sum y_i x_i$ کواریانس شبیه واریانس است ولی بجای مجذور کردن متغیر آن را در دیگری ضرب میکنیم .

در معادله (4-5) مقدار حاصل جمع x و y را درست جب معادله حساب میکنیم چون تمام مقادیر یکسر (x) و (y) و (x^2) را میدانیم . درست معادله بسهولت میتوانیم جمع مجذورات x_i را حساب کنیم یعنی $\sum x_i^2$ (البته مقدار β را نمیدانیم ولی سعی میکنم برآوردی از آن را بدست آوریم) جمله دوم سمت راست معادله را نمیتوانیم حساب کنیم چون مقدار (ϵ_i) را نمیدانیم ، ولی میتوان اندازه این جمله را بر اساس سایر اجزا $\sum x_i^2$ استنتاج کرد . بدیهی است که $\sum x_i^2$ جمع مجذورات است و طبیعی است که دقیقاً میتواند مثبت باشد . از طرف دیگر حاصل جمع

+ یک مورد غیر جالب و عجیب وجود خواهد داشت وقتیکه $\sum x_i^2 = 0$ است و آن وقتی است که $x_1 = x_2 = \dots = x_n = 0$ اگر انحصار گرفتیم کمتر را تغییر ندهد سپس نمیتوان چیزی را راجع بقانون تقاضا استنتاج کرد .

$\sum \epsilon_i x_i$ لزوماً ممکن است مثبت نباشد. چون فرض کرده ایم که (ϵ_i) مستقل از مقدار x توزیع پیدا میکند پس میتوانیم جمع ایکسها (x_i) را در جمع کل وزن های ثابتی فرض کنیم. ایسپلن ها (ϵ_i) تنها بر حسب قانون احتمالات تغییر مییابند. برای هر مقدار معین x میانگین توزیع احتمالی ایسپلن ها (ϵ_i) صفر است. بنا بر این جمع $\sum \epsilon_i x_i$ بصورت جمع وزنی متغیرهای اتفاقی که دارای میانگین صفر است بحساب میآوریم. درست همانطور که در مورد نظریه معمولی نمونه گیری، متحمل ترین مقدار این جمع صفر است، ولی در عمل میتوان انتظار داشت که اگر آزمایش بدفعات زیادی تکرار شود، مقادیر بطور یکنواخت تری توزیع خواهند شد — این مقادیر گاهی مثبت، گاهی منفی و گاهی هم در اطراف میانگین صفر برانگه خواهند شد.

بنا بر این محتمل است که $\sum \epsilon_i x_i$ نسبت به $\sum x_i^2$ کوچکتر شود و کواریانس (ϵ) و (x) احتمالاً نسبت به واریانس x همانطور که اندازه نمونه (n) افزایش میابد، کوچکتر و کوچکتر شود. بمعادله (۴۵)

$$\frac{\sum y_i x_i}{\sum x_i^2} = \beta_{yx} + \frac{\sum \epsilon_i x_i}{\sum x_i^2} \quad (45)$$

بنظر میآید که آخرین جمله آن نسبتاً کوچک است. پس این امر نشان میدهد که مآخیزین جمله را نادیده بگیریم و آن را بعنوان برآوردی برای β_{yx} بکار ببریم.

$$\frac{\sum y_i x_i}{\sum x_i^2} = \hat{\beta}_{yx} \quad (46)$$

در این فرمول کلاهکی که روی قرار دارد نشان میدهد که این برآوردی است از اطلاعات داده شده ولی تفاوت بین برآورد و مقدار واقعی آن بشرح زیر است:

$$\hat{\beta}_{yx} - \beta_{yx} = \frac{\sum \epsilon_i x_i}{\sum x_i^2} \quad (47)$$

اینکه در نمونه گیری میانگین متوجه میشویم همانطور که اندازه نمونه را افزایش دادیم، احتمال یک انحراف معین بین میانگین نمونه گیری و میانگین واقعی، کوچکتر و کوچکتر میشود و میانگین نمونه نیز دارای این گرایش است که بیشتر در اطراف مقدار واقعی تراکم پیدا کند. بهمین ترتیب همانطور که اندازه نمونه افزایش مییابد، مقدار $\sum \epsilon_i x_i / \sum x_i^2$ دارای این گرایش است که نزدیکتر و نزدیکتر بمقدار صفر متمرکز پیدا کند. بدین ترتیب برآورد های $\hat{\beta}_{yx}$ همانطور که اندازه نمونه افزایش مییابد نزدیکتر و نزدیکتر به β_{yx} تراکم پیدا میکند. بنا بر این میتوان چنین نتیجه گیری کرد، در صورتیکه (ϵ) مستقل از x توزیع پیدا کند، برآورد $\hat{\beta}_{yx}$ ارزیاب (Estimator) ثابتی از β_{yx} خواهد بود. بطور کلی همانطور که اندازه نمونه افزایش مییابد، مقدار برآورد شده $\hat{\beta}$ دارای این گرایش خواهد بود که بمقدار واقعی β نزدیکتر و نزدیکتر شود.