

امصائیه تطبیقی برای تحلیل پالیسی

چاپ دوم



محمد شعیب میدری و فهیم احمد یوسفزی



The Asia Foundation

احصائے تطبیقی برای تحلیل پالیسی

محمد شعیب حیدری و فہیم احمد یوسفزی



احصائیه تطبیقی برای تحلیل پالیسی

نویسنده گان

محمد شعیب حیدری و فهیم احمد یوسفزی

دیزاین

سید راشد سادات و روح الله محمدی

چاپ دوم: کابل، 1397

صاحب امتیاز : بنیاد آسیا، دفتر افغانستان

(The Asia Foundation)

تمام حقوق این اثر محفوظ است. تکثیر یا تولید مجدد آن به هر صورت بدون اجازه صاحب امتیاز آن ممنوع میباشد.

فهرست عناوین

بخش اول - معرفی مفاهیم اساسی

10 احصائیه چیست؟
11 جمعیت
11 نمونه
12 تعیین حجم یا اندازه نمونه
14 متغیر و انواع آن
17 فریکونسی
18 تهرین

بخش دوم- معیارهای گرایش مرکزی

20 مود
22 میانه
24 اوسط
28 تهرین

بخش سوم- معیارهای پراگندگی

30 دامنه تغییرات
31 انحراف معیاری
32 انحراف
36 تهرین

بخش چهارم- فرضیه چیست؟

- 38 فرضیه چیست؟
- 39 هدف یک فرضیه
- 39 ویژگی های یک فرضیه خوب
- 40 انواع فرضیه
- 42 تمرین

بخش پنجم- آزمون همبستگی و رگرسیون

- 43 آزمون همبستگی
- 50 رگرسیون
- 51 فرضیه های رگرسیون خطی
- 61 پروبیت رگرسیون
- 62 آزمون میانگین دو جمعیت

سپاسگذاری

در پایان می‌خواهیم از همکاران ما در بنیاد آسیا هر یک داکتر تبسم اکسیر، داکتر زاک وارن، عبدالنواب جلیلی، سید مسعود سادات، محمد شریح شیوان، شمیم سرابی و مهدی فروغ، اشتراک کننده برنامه آموزشی تحلیل دیتا داکتر احمد نوید شمس و محصلین رشته احصائیه پوهنځی اقتصاد پوهنتون کابل حمیدالله قرلق و بی بی رقیه شایق که وقت گرانبهای خویش را برای مرور و ارایه نظریات و پیشنهادات در غنامندی و بهتر شدن این اثر صرف نموده اند ابراز سپاس و امتنان نماییم.

مقدمه

انگیزه اساسی برای تهیه این کتاب در نخست ادامه مواد درسی برای کورس تحلیل دیتا توسط برنامه ستاتا و ضمیمه برای درک بهتر مفاهیمی که در رهنمای برنامه ستاتا تحت عنوان "تحلیل دیتا توسط برنامه ستاتا" است، می باشد.

یکی از اهداف عمده این اثر بلند بردن ظرفیت در نهاد های دولتی و دیگر نهاد های که عملا با موضوعات احصائیه و اقتصاد سنجی (Econometrics) سر و کار دارند میباشد بنابراین، در این اثر کوشش صورت گرفته تا مفاهیم احصائیوی و اقتصاد سنجی (Econometrics) و استفاده از آنها به شکل عملی آن بیان گردد. همچنان مثالهای این اثر به شکلی ترتیب شده است تا برای خواننده قابل درک باشند و با استفاده از مفاهیم احصائیوی به تحلیل و بررسی مشکلات خویش بپردازند. تفاوت دیگر این اثر استفاده از کارتونها است که در هنگام مطالعه این اثر خواننده علاقمندی بیشتر برای خواندن پیدا میکند.

بخشی از اهداف این اثر آشنایی خواننده با نقش احصائیه در طرح و تدوین پالیسی ها می باشد. باید یاد آور شد که نتایج بدست آمده از تحقیقات و سروی ها در شکل دهی و ساختن پالیسی ها با اهمیت می باشد و مورد استفاده قرار می گیرد. هر چند در این اثر در رابطه بین احصائیه و پالیسی سازی صحبت نشده است اما در تالیفات بعدی به این موارد پرداخته خواهد شد.

این کتاب به پنج بخش تقسیم گردیده است. در بخش اول مفاهیم اساسی که در علم احصائیه روزانه مورد استفاده قرار می گیرد به معرفی گرفته است، بخش دوم معیارهای گرایش مرکزی، بخش سوم معیارهای پراکنده گی، بخش چهارم معرفی فرضیه و ساختن فرضیه و در بخش پنجم کتاب آزمون های احصائیوی مانند: همبستگی، رگرسیون خطی و رگرسیون پروبیت و آزمون میانگین دو جامعه (t-test) به تشریح گرفته شده است.

این اثر برای اولین بار در سال 1395 به چاپ رسیده است و چاپ دوم آن در سال 1397 صورت گرفته است. در هردو دوره، کتاب خالی از اشتباهات در نوشتار و تشریح مفاهیم نمیشود و ما را از آن بابت معذور بدارید. برای غامندی هر چه بیشتر کتاب در آینده ما خواهان نظریات و پیشنهادات خواننده گان و مستفیدین محترم میباشیم.

بخش اول

معرفی مفاهیم اساسی

معرفی

در دنیای جدید که عصر تکنالوژی و معلومات بوده، هر روز به اهمیت علم احصائیه که بخش جدا ناشدنی و اساس مهم برای پیشرفت تکنالوژی و ارایه معلومات بوده، افزوده میشود. احصائیه بنیاد و اساس مهم برای پیشرفت در عرصه های مختلف چون (زراعت، اقتصاد، بیولوژی، طب، صنعت، پلانگذاری، تعلیم و تربیه) و سایر علوم میباشد. احصائیه تقریباً در تمام عرصه های زندگی بشری کاربرد وسیع دارد.

احصائیه چیست (Statistics)?

در علم احصائیه موضوعات مورد بحث شامل جمع آوری، ترتیب و تنظیم، خلاصه سازی، تحلیل و ارایه دیتا و همچنان نتیجه گیری و تصمیم گیری بر اساس یافته ها از دیتاهای بدست آمده میباشد. احصائیه بیشتر به جمع آوری معلومات کمی به گونه سیستماتیک و تحلیل و تفسیر معلومات گردآوری شده میپردازد. کلمه احصائیه شامل مفاهیم مانند:

1. حقایق کمی ، مثلاً تعداد افراد که در یک منطقه زندگی میکنند.
2. مطالعه شیوه های جمع آوری دیتا و تحلیل و تفسیر حقایق بدست آمده میباشد.

به صورت عموم احصائیه به دو بخش تقسیم میگردد:

1. احصائیه تشریحی
 2. احصائیه استنباطی
- احصائیه تشریحی** یا عبارت از آن بخشی از احصائیه می باشد که تنظیم، تشریح و نمایش دیتا را در بر می گیرد. مسائل مانند توزیع فریکونسی، میانگین، مود، میدیان (میانه)، معیار های پراگندگی مانند دامنه ی تغییرات، انحراف معیاری و وریانس شامل احصائیه تشریحی اند.
- احصائیه استنباطی** یا تحلیلی عبارت از آن بخشی از احصائیه است که معلومات بدست آمده از نمونه یک جمعیت را تحلیل نموده و پیش بینی و نتیجه گیری لازم را در مورد ارائه می کند. موضوعات مانند تیوری نمونه گیری، آزمون فرضیه ها، رگرسیون (وابستگی) کورلیشن (همبستگی) مربوط احصائیه استنباطی یا تحلیلی می شوند.

جمعیت (Population)

در مباحث احصائیوی تمام واحدهای که در یک مجموعه شامل و حداقل دارای یک صفت مشترک میباشند بنام جمعیت و یا جامعه یاد میشوند. یا به عبارت دیگر تمام مشاهدات شامل در یک مطالعه عبارت از جمعیت میباشند. به گونه مثال تمام کارمندان یک وزارتخانه، تمام شاگردان یک مکتب و یا تمام کتاب های موجود در یک کتابخانه عبارت از مثال های از جمعیت میباشد.

نمونه (Sample)

در مباحث احصائیوی نمونه به مفهوم یک نسبت انتخاب شده کوچکتر از جمعیت میباشد. یا یک مجموعه کوچک از مشاهدات انتخاب شده از جمعیت میباشد. تعداد واحدها یا مشاهدات در یک نمونه عبارت از حجم نمونه میباشد. در تصویر پائین نمونه و جمعیت نمایش داده شده است. جمعیت در تصویر تمام شاگردان صنف می باشد و سه نفر شاگرد که در نزدیک تخته ایستاده اند عبارت از نمونه است که از جمعیت یعنی صنف نماینده گی میکند.



شکل: ۱-۱ نمونه و جمعیت

به ادامه مثال گذشته اگر در یک وزارتخانه به تعداد 2,000 کارمند موجود باشد و از میان آنها 200 نفر برای یک تحقیق انتخاب گردند، 2,000 نفر عبارت از جمعیت و 200 نفر انتخاب شده عبارت از نمونه تمام کارمندان وزارتخانه (جمعیت) میباشد.

هرگاه یک جمعیت بزرگ باشد و وقت کافی، منابع مالی و امکان بررسی تمام افراد شامل در یک جمعیت ممکن نباشد در آن صورت از نمونه گیری استفاده می شود.

چرا نمونه از یک جامعه گرفته میشود؟

- بخاطر کمبود وقت
- مصارف زیاد
- غیر ممکن و غیر معقول بودن مطالعه یک جامعه در کل

تعیین حجم یا اندازه نمونه (Sample Size Selection)

بعد از این که جمعیت مورد نظران را تعیین نمودید در قدم بعدی تصمیم این است که حجم نمونه تان به چه اندازه باشد تا نماینده گی از جمعیت تان کرده بتواند. فرض کنید که شما یک تحقیق در مورد اطفال زیر سن 18 سال که مشغول به کار هستند انجام دهید. با فرض این که مجموع اطفال زیر سن کارگر در افغانستان به 500,000 تن میرسد و بنابر محدودیت وقت و مصارف زیاد نمیتوانید تمام 500,000 کودک کارگر زیر سن 18 سال را سروی کنید که آیا از کار خود راضی هستند و یا خیر، بنابراین مجبور هستید با تعدادی کمتر اطفال سروی را انجام دهید تا از تمام جمعیت نماینده گی کرده بتواند. اما سوالی که در این جا پیش می آید، این است که حجم نمونه باید چه مقدار انتخاب شود؟

اساساً حجم نمونه مربوط به این میشود که شما به چه اندازه میخواهید نتایج بدست آمده از نمونه مورد نظران با جامعه مطابقت داشته باشد و یا به چه اندازه دیتای تان دقیق باشد تا نماینده گی از تمام جمعیت مورد نظر نموده بتواند.

در اینجا اساساً باید به دو نکته ذیل توجه کرد. 1 حاشیه خطا (Margin of error) 2 سطح اطمینان Confidence level

1. حاشیه خطا (Margin of Error)

حاشیه خطا عبارت از خطای نمونه گیری میباشد به طور مثال شما یک تعداد معین از اطفال کارگر را سروی نموده اید و نتیجه بدست آمده شما این است که 40 فیصد اطفال از کار خود راضی هستند در صورت که حاشیه خطای شما 5% باشد پس شما گفته میتوانید که 35 الی 45 فیصد مجموع اطفال نسبت به کار خویش رضایت دارند.

به یاد داشته باشید به هر اندازه که بخواهید حاشیه خطا در سروی تان کم باشد پس به نمونه بزرگتر ضرورت دارید و هراندازه که نمونه کوچکتر انتخاب نموده باشد حاشیه خطا بزرگتر میباشد.

2. سطح اطمینان (Confidence Level)

سطح اطمینان معمولاً در حدود 90%، 95% و 99% تعیین میشود. سطح اطمینان برای ما این را بازگو میکند که به کدام اندازه بالای نتیجه بدست آمده خود اطمینان داریم.

به تعقیب مثال قبلی اگر سطح اطمینان 95% باشد پس گفته میتوانیم که با 95 فیصد اطمینان در حدود 35% الی 45% اطفال کارگر زیر سن 18 سال نسبت به کار خود رضایت دارند.

بعد از تصمیم در مورد حاشیه خطا و سطح اطمینان شما میتوانید که اندازه نمونه خود را تعیین کنید، جدول ذیل بیانگر تعیین حجم نمونه به اساس جمعیت مورد نظر با درنظرداشت درجه های مختلف سطح اطمینان و حاشیه خطا میباشد.

Confidence Interval = 99% سطح اطمینان			Confidence Interval = 95% سطح اطمینان			
حاشیه خطا Margin of error (MoE)			حاشیه خطا Margin of error (MoE)			
1%	2.5%	5%	1%	2.5%	5%	(Population size) تعداد جمعیت
99	96	87	99	94	80	100
485	421	285	475	377	217	500
943	727	399	906	606	278	1,000
6,239	2,098	622	4,899	1,332	370	10,000
14,227	2,585	659	8,762	1,513	383	100,000
16,055	2,640	663	9,423	1,532	384	500,000
16,317	2,647	663	9,512	1,534	384	1,000,000

شکل: ۲-۱ تعیین حجم نمونه به اساس جمعیت، حاشیه خطا و سطح اطمینان

از جدول فوق چنین نتیجه میگیریم که برای یک جمعیت که در حدود 500,000 نفر را در بر داشته باشد باید که نمونه به اندازه 384 نفر انتخاب کنیم (با درنظرداشت سطح اطمینان 95% و حاشیه خطا 5%).

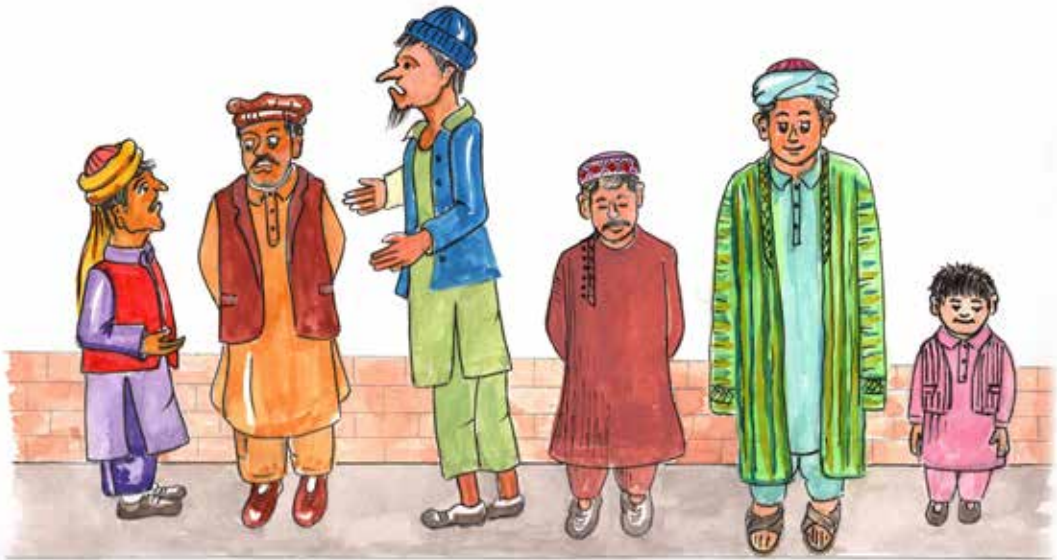
در صورت که جمعیت شما بیشتر از 1,000,000 باشد برای تعیین حجم نمونه به لینک ذیل مراجعه نمایید.

<https://www.checkmarket.com/sample-size-calculator/>

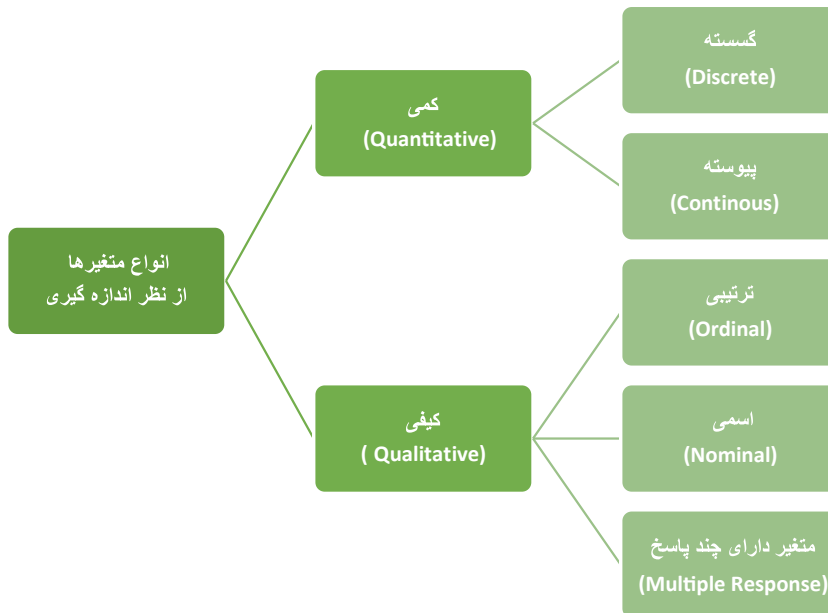
متغیر و انواع آن (Variable and Types of Variable)

متغیر چیست؟

یک مشخصه یا معلومات که در حال تغییر بوده و قیمت های متفاوت را به خود گرفته میتواند. اندازه قد افراد یک جامعه میتواند یک متغیر باشد اندازه قد هر فرد میتواند نظر به فرد دیگر متفاوت باشد. در شکل پائین اندازه قد افراد یک مثال خوب برای نشان دادن متغیر است به دلیل اینکه قد تمام افراد معمولاً یکسان نیست و نظر به هر فرد میتواند تغییر کند.



شکل: ۳-۱ افراد با اندازه قد های متفاوت



1. **متغیر کیفی (Qualitative):** عبارت از متغیر است که تغیر در آن در اندازه ی آن نبوده بلکه تغیر در نوع آن می باشد و انجام دادن عملیه های ریاضی با این متغیر ها معنی ندارد.

مانند: جنسیت که به مرد و زن تقسیم شده و تغیر تنها در نوع جنسیت است نه این که مرد یا زن بودن جنسیت از نگاه عددی دارای کدام مفهوم باشد.

2. **متغیر کمی (Quantitative):** عبارت از متغیر می باشد که تغیرات در اندازه آن صورت می گیرد و انجام دادن عملیه های حسابی با این متغیرها ممکن و با مفهوم می باشد و یا متغیرهای که از اعداد تشکیل شده است در مقابل متغیر کیفی مانند متغیر جنسیت که مرد و زن می باشد از حروف تشکیل شده است.

مانند: سن، عاید، تعداد کتاب که میتواند گفت 23 عدد کتاب که عدد 23 دارای یک مفهوم است برای انجام دادن عملی های حسابی

انواع متغیرهای کمی

- **متغیر پیوسته (Continuous):** نوعی از متغیر است که ارزش های نا محدود را در بر میگیرد و بین دو واحد آن هر نقطه یا ارزشی را میتوان انتخاب کرد.
- به عنوان مثال، وزن یک متغیر پیوسته است
- **متغیر گسسته (Discrete):** متغیری است که ارزش های محدود را در بر میگیرد.
- مانند: ماه های سال از 1 الی 12 است. و رنگهای رنگین کمان که از 1 الی 7 رنگ است.

انواع متغیرهای کیفی

- **متغیر ترتیبی (Ordinal):** در این نوع متغیر ترتیب پاسخ ها بسیار با اهمیت و دارای مفهوم می باشد. مانند: سطح تحصیل که از کم شروع به طرف سطح بلندتر ادامه میابد.

1 = دوره ابتدایی 1 = بسیار ناراحت

2 = دوره متوسطه 2 = ناراحت

3 = دوره لیسه 3 = خوشحال

- **متغیر اسمی (Nominal):** ترتیب پاسخ ها در این نوع متغیر، دارای اهمیت و مفهوم نمی باشد، مانند: جنسیت، حالت مدنی

1 = مرد 1 = متاهل

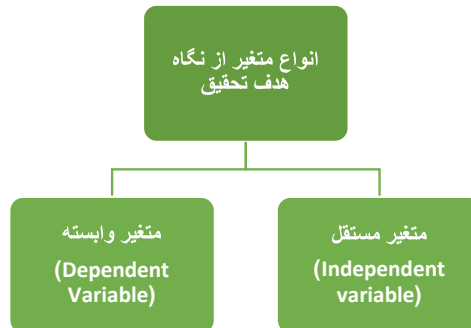
2 = زن 2 = مجرد

- **متغیر دارای چند پاسخ (Multiple Response Variable):** نوعی از متغیر است که دارای چند گزینه می باشد. به طور مثال: شما کدام وسایل ذیل را برای بدست آوردن اطلاعات و معلومات استفاده میکنید؟

1 = رادیو 3 = شبکه های اجتماعی 5 = روزنامه

2 = تلویزیون 4 = مسجد 6 = شورای قریه

انواع متغیر از نگاه هدف تحقیق



- **متغیر وابسته (Dependent variable):** طوری که از نام آن پیداست، عبارت از متغیری است که قیمت ها یا نتیجه آن وابسته به قیمت های متغیر مستقل یا عوامل دیگر است.
- **متغیر مستقل (Independent variable):** این متغیر طوری که از نام آن پیداست به گونه مستقل قیمت گرفته میتواند و قیمت های آن وابسته به متغیر دیگر یا عوامل دیگر نیست.

به طور مثال اگر ما بخواهیم یک نوع آزمایش را در مورد این که بالای سطح خوشحالی کدام عوامل تاثیر گذار می باشد انجام دهیم. در اینجا سطح خوشحالی عبارت از متغیر وابسته بوده و مجموع عوامل دیگر که بالای سطح خوشحالی تاثیر گذار اند عبارت از متغیر های مستقل می باشند.

فریکونسی (Frequency)

در علم احصائیه فریکونسی یک دیتا یا معلومات مشخص به معنی این است که این ارزش چند بار تکرار شده است.

به گونه مثال، اگر در یکی از آمریت های وزارت صحت عامه 4 نفر از طبقه ذکور و 6 نفر از طبقه اناث باشد، در نتیجه فریکونسی برای طبقه ذکور 4 و برای اناث 6 می باشد.

فریکونسی همیشه بوسیله حرف انگلیسی (f) نشان داده میشود.

جدول فریکونسی

یک جدول فریکونسی عبارت از بیان معلومات احصائیوی به شکل ابتدایی اما ارزشمند می باشد. ساختار یک جدول فریکونسی به گونه ی می باشد که گزینه های یک سوال به ترتیب نوشته شده و فریکونسی هر کدام از گزینه ها شده ذکر میشود.

مثال:

شغل پاسخ دهنده گان سروی، 2015

وظیفه	فریکونسی	فیصدی
شاغل	3,946	41
متقاعد	79	0.8
خانم خانه	4,073	42.6
متعلم	627	6.6
بیکار	837	8.7
مجموعه	9,562	100

شکل: ۱-۴ جدول فریکونسی شغل پاسخ دهنده گان



تمرین

1. راجع به علم احصائیه معلومات دهید؟
2. جمعیت و نمونه چه است و برای هر کدام دو مثال بنویسید؟
3. چرا از نمونه استفاده میکنیم؟
4. با درنظر داشت 99% سطح اطمینان و 1% حاشیه خطا برای انجام تحقیق از یک جمعیت که تعداد شان به 100,000 خانم میرسد حد اقل به چه اندازه نمونه ضرورت است؟
5. انواع متغیر از نگاه اندازه گیری را نام بگیرید؟
6. برای متغیر اسمی و پیوسته مثال دهید؟
7. متغیر وابسته و مستقل را تعریف و برای هر کدام دو مثال دهید؟

معیارهای گرایش مرکزی (Measures of Central Tendency)

معرفی

برای مطالعه یک جمعیت و به خصوص یک مشخصه خاص آن جمعیت مانند سن، بررسی تمام افراد در یک جمعیت به گونه انفرادی مشکل به نظر میرسد که به یک نتیجه در مورد سن افراد برسیم، یعنی با مشاهده هر کدام از سن افراد نمیتوان گفت که آیا جمعیت مورد مطالعه نظر به سال گذشته چه تغییرات را در مشخصه سن افراد شاهد بوده است. اما بهتر خواهد بود تا یک سن را دریافته و به اساس آن در مورد جمعیت بتوان اظهار نظر کرد. یکی از مثال ها میتواند این باشد که اوسط سن فارغین پوهنتون طبی در سال 1385، 29 سال بوده اما در سال 1394 به 26 سال کاهش پیدا کرده است. در این مثال اوسط سن به نماینده گی از تمام سن محصلین به بحث گرفته میشود. زیرا یاد آوری سن هر یک از محصلین بی مفهوم بوده ما را به نتیجه نمی رساند.

یکی از معیارهایی که با دریافت آن میتوان به نماینده گی از جمعیت صحبت کرد، عبارت از معیارهای گرایش مرکزی میباشد.

تعریف

معیار گرایش مرکزی عبارت از معیار است که مرکز دیتا را نشان داده و دیگر ارزش‌ها در یک مجموعه ازدیتاها در اطراف آن قرار دارند. میتوان برای فهم بهتر موضوع به دو فرد که در شکل پائین در مرکز سایر افراد ایستاده اند و موسیقی می نوازند و سایر افراد در اطراف این دو نفر که اتن ملی را انجام میدهند از مثالی برای معیار های گرایش مرکزی یاد آور شد که شاخصهای هستند برای نشان دادن مرکز دیتا .



شکل: ۲-۱ معیارهای گرایش مرکزی

مثال: اوسط یکی از معیارهای گرایش مرکزی بوده، فرضاً اوسط سن کارمندان وزارت مالیه در سال 1392، 22 سال بوده اما در سال 1393 این اوسط به 25 سال افزایش پیدا کرده، این افزایش به این معنی است که سن کارمندان وزارت مالیه به صورت اوسط نظر به سال 1392 در سال 1393 افزایش پیدا کرده است.

در این اثر مود، میدیان و اوسط به معرفی گرفته میشوند.

مود (Mode): مود ساده ترین شاخص گرایش مرکزی است، مود عددی است که دارای بیشترین فریکونسی است، یا عددی که بیشتر از اعداد دیگر تکرار شده است. مود از طریق مشاهده ی عددی که بیشتر تکرار گردیده است تعیین می گردد. در شکل پائین گل ارغوانی مود است زیرا نسبت به دیگر گل ها بیشتر در تصویر تکرار شده است.



شکل: ۲-۲ نمایش مود در تصویر

مود در مطالعات بازاریابی حایز اهمیت میباشد، به عنوان مثال یک مدیر بازاریابی میتواند با استفاده از مود دریابد که کدام رنگ لباس بیشترین تقاضا را در بازار دارا میباشد.

به گونه مثال معلومات به دست آمده از یک دیتا ست وزارت معارف به کمک مود دریافته که متعلمین بیشتر علاقمند به خواندن کدام کتاب درسی هستند.

کتاب های مورد علاقه شاگردان لیسه های شهر کابل، 1394

کتاب	فریکونسی (به نفر)
هندسه	8,555
دری	12,333
تاریخ	30,000
ریاضی	9,990

شکل: ۲-۳ جدول کتاب های مورد علاقه شاگردان برای مطالعه

از جدول داده شده در بالا این موضوع واضح میگردد که بیشترین شاگردان (30,000 شاگرد) در لیسه های شهر کابل علاقه مند به مطالعه کتاب تاریخ هستند.

مثال: در اعداد ذیل اندازه های 8 جوړه بوت میباشد که در یک روز از یک دوکان خریداری شده، مود عبارت از اندازه 44 میباشد زیرا نسبت به دیگر اندازه ها زیادتر تکرار شده است.

44, 45, 44, 41, 44, 43, 42, 40

مثال: در اعداد ذیل مود عبارت از اعداد 3 و 6 میباشد زیرا نسبت به دیگر اعداد زیاد تر تکرار گردیده و دارای فریکونسی مساوی اند.

2 6 7 6 1 3 6 3 7 5 3 1

همچنان در صورت که دیتا ما طبقه بندی شده باشد. هرطبقه و گروپ که بزرگ ترین فریکونسی یا تکرار را دارا باشد، همان طبقه یا گروپ، مود است.

میانه (Median): میدیان نقطه ی وسط یک سلسله از ارزش ها یا اعداد میباشد، که سلسله اعداد را به دو قسمت تقسیم کرده، 50% اعداد بزرگتر از میدیان و 50% دیگر کوچکتر از میدیان و خود میدیان در وسط اعداد قرار میگیرد. بنابراین میدیان نقطه وسطی است یعنی عددی است که سلسله اعداد را به دو قسمت مساوی تقسیم می کند.

به عبارت دگر اگر یک سلسله ی اعداد را به ترتیب لست کنیم، میانه (میدیان) همان عددی است که در وسط این سلسله اعداد قرار می گیرد.

میدیان زمانی منحصیث معیار گرایش مرکزی مورد استفاده قرار میگیرد که سلسله اعداد دارای یک عدد فوق العاده کوچک یا بزرگ باشد. در مثال پائین یک درخت که در وسط سایر درختها قرار دارد دیده میشود که ارتفاع آن نسبت به سایر درختان بسیار بلند می باشد در این حالت باید برای دریافت مرکز دیتا از میدیان استفاده شود.



شکل: ۲-۴ مورد استفاده میدیان

در اعداد ذیل میدیان عبارت از عدد (4.5) میباشد

8 6 3 1 2 4 5 7

مرحله اول : ترتیب اعداد به شکل نزولی یا صعودی.

1 2 3 4 5 6 7 8

مرحله دوم : چون مجموع مشاهدات ما جفت است بنأ حاصل جمع دو عدد که در وسط قرار دارد را تقسیم بر 2 میکنیم که عدد بدست آمده عبارت از میدیان میباشد. $4.5 = (4+5)/2$

نوت : اگر مجموع مشاهدات ما طاق باشد عدد که در وسط مشاهدات ما (بعد از ترتیب نمودن اعداد) قرار دارد عبارت از میدیان میباشد.

مثال: در سلسله اعداد ذیل میدیان را دریابید؟

8 6 11 18 10 22 12

مرحله اول اعداد را ترتیب میکنیم از کم به زیاد (و یا از زیاد به کم) و عددی که در وسط سلسله اعداد قرار

گیرد عبارت از میدان می باشد.

به طور مثال در مثال قبلی میدان عبارت از عدد 11 میباشد.

22-18-12 11 10-8-6



نکته: میدان زمانی مورد استفاده قرار میگیرد که سلسله اعداد دارای یک **عدد فوق العاده کوچک** یا **بزرگ** باشد.

مثال: در پائین مقدار تیل که ۱۰ موتر در وزارت مخابرات در یک ماه مصرف می رسانند داده شده و به صورت مشخص موتر با شماره سریال ۲۳۳ نظر به سایر موترها به مقدار قابل ملاحظه ای تیل بیشتر مصرف میکند که در این حالت بهتر است از میدان برای دریافت مرکز دیتا استفاده گردد.

مصرف تیل وسایط نقلیه در وزارت ترانسپورت - 31 حمل 1393

مقدار تیل مصرف شده (لیتر)	شماره موتر
330	225
455	226
322	227
544	228
444	229
456	230
654	231
233	232
2,300	233
267	234

شکل: ۵-۲ جدول مصرف تیل موترها

اوسط (Mean): مشهورترین شاخص گرایش مرکزی اوسط است. اوسط از طریق جمع کردن تمام اعداد و تقسیم حاصل جمع آن بر تعداد تمام مشاهدات به دست می آید. در شکل پائین سه مرد با اندازه ریش های متفاوت رسامی شده و شخصی که در میان دو مرد دیگر ایستاده بلندی ریش وی نسبت به دو فرد دیگر در وسط قرار دارد. و مثالی از اوسط میباشد.



شکل: ۶-۲ نمایش اوسط توسط ریش های افراد با اندازه های متفاوت

در اعداد ذیل اوسط عبارت از عدد (3.3) میباشد.

2 7 5 4 1 2 3 5 3 1

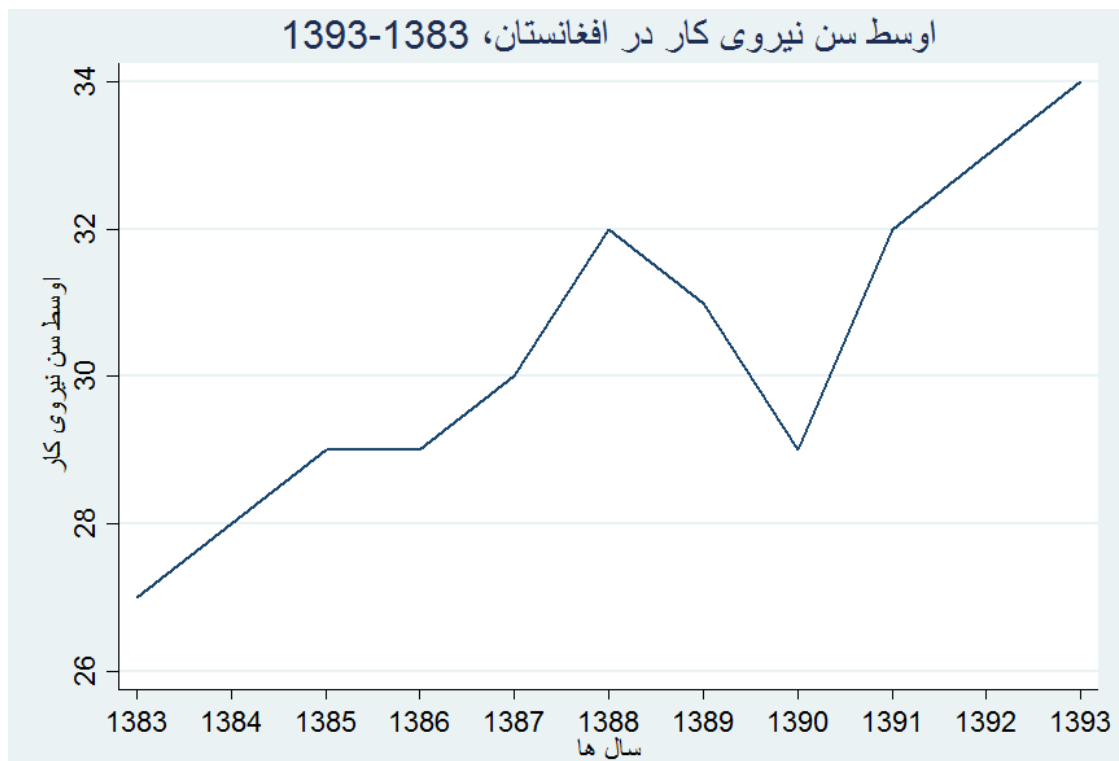
مرحله اول : تمام نمرات را جمع میکنیم .

$$33 = 2 + 7 + 5 + 4 + 1 + 2 + 3 + 5 + 3 + 1$$

مرحله دوم : حاصل جمع تمام نمرات را تقسیم تعداد مشاهدات میکنیم

$$3.3 = 10 / 33$$

مثال: از اوسط میتوان استفاده های مختلف کرد، مثال ذیل اوسط سن نیروی کار در افغانستان را در جریان 10 سال نشان میدهد که اوسط سن نیروی کاری از 27 سال در سال 1383 به 34 سال در سال 1393 رسیده است.



شکل: ۷-۲ گراف اوسط سن نیروی کار در ۱۰ سال

مثال: مصارف دو ریاست وزارت تجارت (پلان و پالیسی و انکشاف سکتور خصوصی) برای ۱۲ ماه سال در دسترس است، توسط اوسط میتوان دریافت که به صورت اوسط مصارف سالانه کدام ریاست بیشتر است:

ماه ها	پلان و پالیسی (افغانی)	انکشاف سکتور خصوصی (افغانی)
1	15,000	12,000
2	12,300	13,200
3	13,300	14,500
4	11,000	11,000
5	10,000	12,020
6	12,300	14,000
7	11,000	15,000
8	14,000	10,000
9	11,500	8,500
10	10,000	12,300
11	15,000	16,000

12	13,000	15,500
اوسط	12,835	12,367

شکل: ۸-۲ جدول مصارف دو ریاست وزارت تجارت

حالا نظر به جدول در بالا بعد از 12 ماه معلوم گردیده که مصارف ماهانه ریاست پلان و پالیسی به صورت اوسط نسبت به انکشاف سکتور خصوصی بیشتر است.

به عنوان مثال در سلسله اعداد پائین عواید ماهانه 5 آمریت ریاست زراعت ولایت لوگر داده شده در این میان عواید یکی از آمریت ها با 4 آمریت دیگر تفاوت زیادی دارد و در این خواهد بود تا از میدیان استفاده گردد.

نام آمریت	عواید ماهانه
اداری	1,000
تحقیق	1,200
تهیه و تدارکات	1,400
کنترول و ارزیابی	800
مالی	170,000

شکل: ۹-۲ جدول عواید ماهانه ۵ آمریت ریاست زراعت ولایت لوگر

مثلا اگر بخواهیم دریابیم که به طور اوسط ریاست زراعت ولایت لوگر به چه اندازه عواید دارد به این منظور از اوسط (mean) استفاده می کنیم. در این صورت عواید تمام آمریت ها را با هم جمع میکنیم و تقسیم بر تعداد آمریت ها میکنیم که 34,880 میشود.

با وجود این که نتیجه گیری ما اشتباه نمیشود، اما میتواند گمراه کننده باشد و ما را به راه اشتباه راهنمایی کند و موجب اتخاذ تصمیم اشتباه گردد زیرا 34,880 نسبت به عواید 4 آمریت (اداری، تحقیق، تهیه و تدارکات و کنترول و ارزیابی) بلند میباشد و نسبت به عواید آمریت مالی به مراتب کمتر میباشد.

پس در این جا بهتر این خواهد بود که از میدیان استفاده کنیم. میدیان بعد از محاسبه 1,400 بدست می آید.



تمرین

1. معیارهای گرایش مرکزی چیست؟
2. مود چیست و یک مثال دهید؟
3. در کدام موقع برای پیدا نمودن مرکز دیتا باید از میدان استفاده کنیم؟
4. میدان را توضیح دهید و در سلسله اعداد ذیل میدان را دریابید؟
12, 23, 13, 18, 25, 15
5. عواید پنج شرکت هوایی به اساس معلومات وزارت مالیه در سال 1393 قرار جدول ذیل می باشد، برای دریافت مرکز دیتا باید از کدام یک از معیارهای گرایش مرکزی استفاده کرد؟

عواید به میلیون افغانی	نام شرکت هوایی
330	امارات
0.5	کام ایر
245	صافی
340	آریانا
430	افق شرق

شکل: ۲-۱۰ عواید ۵ شرکت هوایی در افغانستان

بخش سوم

معیارهای پراگندگی (Measures of Dispersion)

معرفی

با مطالعه معیارهای گرایش مرکزی میتوان مرکز یک سلسله اعداد یا یک توزیع را دریافت، اما معیارهای گرایش مرکزی توانایی این را ندارند تا در مورد اینکه به کدام اندازه اعداد در دو طرف مرکز اعداد پراکنده هستند معلومات دهد. در یک سلسله اعداد یا یک توزیع تمام اعداد با هم مساوی نیستند. به معنی اینکه شاید بعضی اعداد از مرکز بسیار تفاوت داشته و یا بعضی کمتر تفاوت از مرکز را نشان دهند. این تفاوت بنام انحراف یاد شده که توسط معیارهای پراگندگی نشان داده میشود. انحراف کمتر به معنی این است که در یک سلسله اعداد تفاوت کمتر هر یک از اعداد از اوسط وجود دارد. از سوی دیگر انحراف بیشتر به معنی این است در یک سلسله اعداد تفاوت بین هر یک اعداد (مشاهدات) از اوسط زیاد میباشد و کمتر مشابه هم هستند.

مثال ذیل را در نظر داشته باشید:

مقایسه نمرات دو شاگرد، لیسه حبیبیه 1394

مضامین	نمرات احمد	نمرات فهیم
انگلیسی	85	68
تاریخ	90	75
دری	80	65
جغرافیه	25	67
ریاضی	65	70
مجموعه	345	345
اوسط	69	69

شکل: ۱-۳ جدول مقایسه نمرات دو شاگرد لیسه حبیبیه

در جدول بالا نمرات دو تن از شاگردان لیسه حبیبیه در 5 مضمون مختلف داده شده است، که مجموعه هر دوی آنها 345 و در نتیجه اوسط نمرات آنها مساوی به 69 می باشد. حقیقت این است که احمد در یکی از مضامین نمره 25 را گرفته ناکام مانده حال آنکه هر دوی آنها دارای مجموع نمرات و اوسط مساوی میباشند. از جدول بالا

میتوان چنین نتیجه گرفت که نمرات فهم یکی از دیگری کمتر انحراف دارد و برعکس نمرات احمد دارای انحراف بیشتر است و انحراف کمتر مطلوب است.

معیارهای که در اینجا ما به منظور دانستن و آشنایی با انحراف مورد مطالعه قرار میدهیم عبارت اند از: دامنه تغییرات (Range)، انحراف (Variance) و انحراف معیاری (Standard Deviation).

دامنه تغییرات (Range): یکی از ابتدایی ترین معیارهای پراکندگی عبارت از دامنه تغییرات می باشد که تفاوت میان کوچکترین و بزرگترین عدد را در یک سلسله اعداد نشان میدهد. دامنه تغییرات در شکل که در پائین رسامی شده عبارت از تفاوت بین ارتفاع قد بلندترین پسر که کلاه زرد به سر داشته و قد کوتاه ترین پسر که پیراهن سرخ به تن دارد دریافت شده میتواند.



شکل: ۳-۲ دامنه تغییرات از کوتاه ترین تا بلندترین قد

فرمول دامنه تغییرات عبارت از:

دامنه تغییرات = بزرگترین عدد - کوچکترین عدد

$$\text{Range} = L - S$$

L = largest value

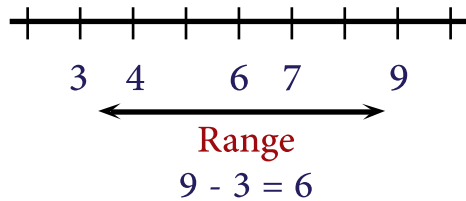
S = smallest value

مثال: در سلسله اعداد ذیل دامنه تغییرات را دریابید؟

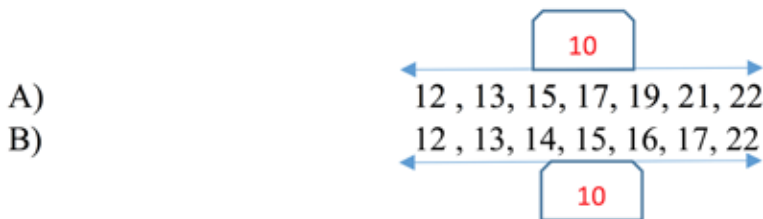
3, 4, 6, 7, 9

$$\text{Range} = L - S$$

$$\text{Range} = 9 - 3 = 6$$



با وجود که دامنه تغییرات یکی از معیارات پراگندگی است اما منحنیث یکی از بهترین معیار های پراگندگی به شمار می رود. زیرا دامنه تغییرات تنها با کوچکترین عدد و بزرگترین عدد یک سلسله اعداد سرو کار دارد. به دو مثال ذیل توجه کنید با وجود که پراگندگی در سلسله اعداد اولی نسبت به سلسله اعداد دومی کمتر است اما باز هم هر دو دیتا ست دارای دامنه تغییرات یکسان اند.



انحراف معیاری (Standard Deviation)

هر یک انحراف معیار جمعیت و انحراف معیار نمونه دارای فرمول های جداگانه می باشد.

انحراف معیار جمعیت:

$$\sigma = \sqrt{\frac{1}{N} \sum (x_i - \mu)^2}$$

انحراف معیار نمونه :

$$s = \sqrt{\frac{1}{N-1} \sum (x_i - \bar{x})^2}$$

انحراف معیاری عبارت از یک معیاری برای اندازه گیری پراگندگی اعداد از اوسط آنها است.

سمبول که انحراف معیاری به وسیله آن نشان داده میشود:

الف: σ (سیگما حرف لاتین) میباشد برای نشان دادن انحراف معیاری یک جمعیت استفاده میشود.

ب: s، برای نشان دادن انحراف معیاری یک نمونه مورد استفاده قرار می گیرد.

فورمول برای دریافت انحراف معیاری بسیار ساده است؛ و عبارت از جذر دوم انحراف می باشد. شاید حالا برای شما این سوال مطرح گردد که انحراف چیست؟ به همین خاطر در مرحله اول انحراف را به مطالعه میگیریم.

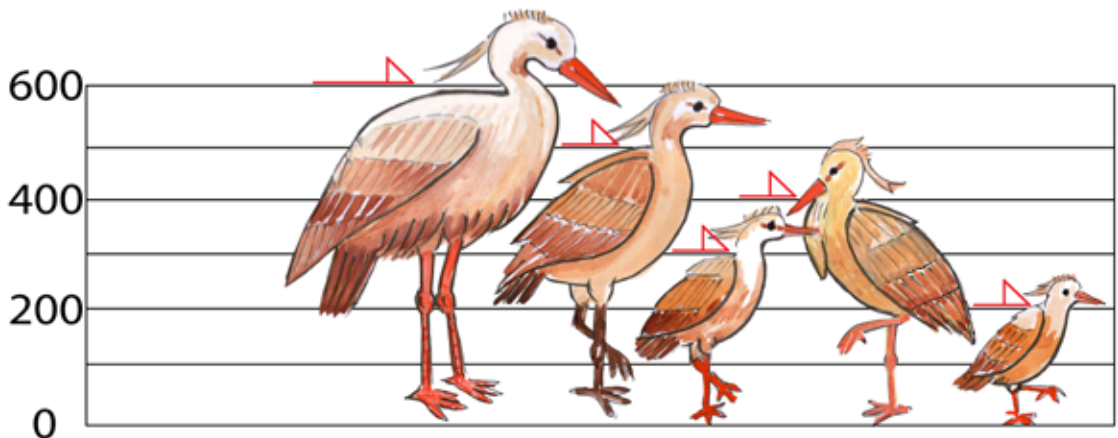
انحراف (Variance)

عبارت از اوسط مربعات تفاوت ها از اوسط یک سلسله اعداد یا یک توزیع می باشد. یا به عبارت دگر انحراف یا وریانس عبارت از مجموع مربعات تفاوت های یک سلسله ی اعداد از اوسط آن اعداد، تقسیم بر تعداد آن سلسله عددی است.

نگران نباشد شاید تعاریف بالا کمی پیچیده به نظر برسد اما برای فهم ساده تر انحراف به مراحل ذیل توجه کنید .

- در مرحله اول اوسط یک سلسله اعداد را دریابید.
- در مرحله بعد از هر عدد، اوسط بدست آمده را تفریق نموده و نتیجه آن را مربع نمایید (مربعات تفاوت ها از اوسط).
- مجموع آنها (مربعات تفاوت از اوسط) را تقسیم بر تعداد آنها میکنیم.
- اما به یاد داشته باشید در صورت پیدا نمودن انحراف نمونه باید که مجموع آنها را تقسیم (n-1) کنید.

مثال: برای فهم بیشتر به تصویر ذیل که قد (به میلی متر) چند حیوان (مرغابی ها) داده شده توجه کنید.



قد هر یک (تا شانه): 600 میلی متر، 500 میلی متر، 300 میلی متر، 400 میلی متر و 200 میلی متر میباشد.

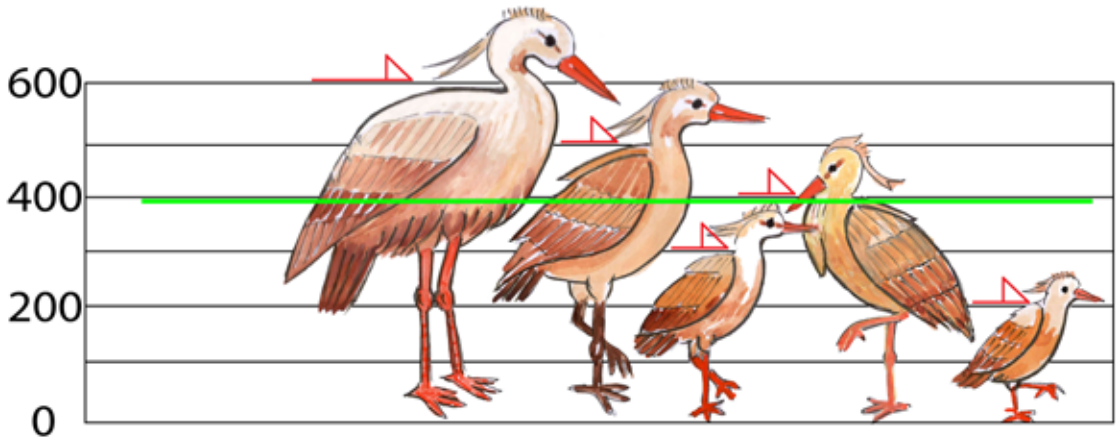
حالا اوسط، انحراف و انحراف معیاری قدهای داده شده در بالا را دریابید.

در مرحله اول اوسط را پیدا میکنیم.

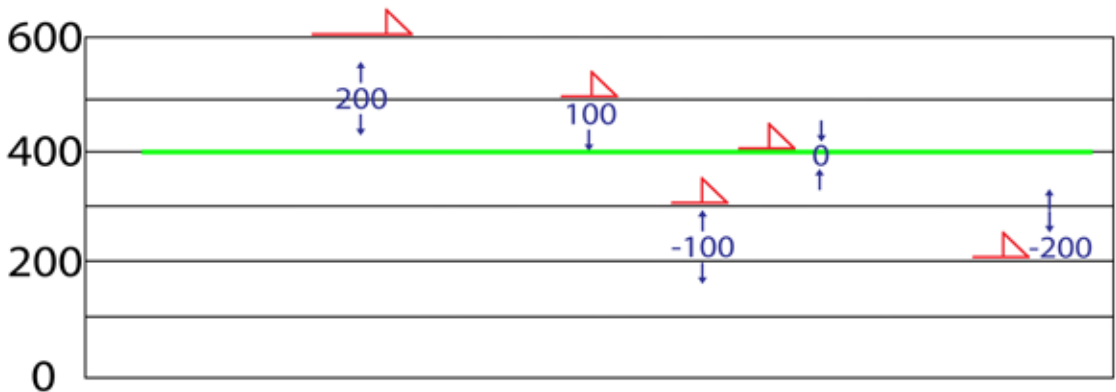
پاسخ:

$$\text{اوسط} = \text{Mean} = \frac{600+500+300+400+200}{5} = \frac{2000}{5} = 400$$

دیدیم که اوسط قدها عبارت از 400 میلی متر می باشد، حالا اوسط را توسط خط سبز نشانی میکنیم.



در این مرحله تفاوت هریک از قدها را از اوسط پیدا میکنیم.





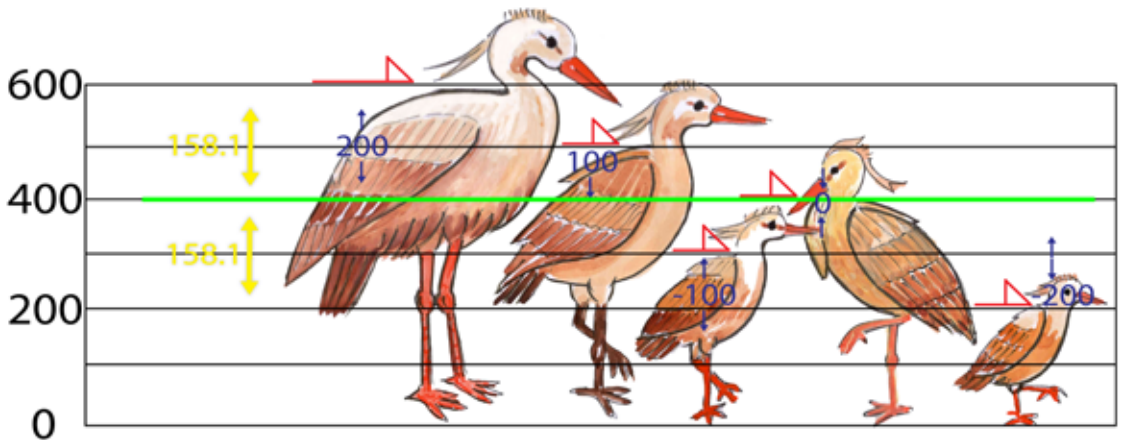
نکته: انحراف به سمبول s^2 نشان داده می باشد.

$$\text{انحراف} = \text{Variance} = \frac{200^2 + 100^2 + (-100)^2 + 0^2 + (-200)^2}{5 - 1} = \frac{100,000}{4} = 25,000$$

در نتیجه انحراف مساوی به 25,000 می باشد.

و انحراف معیاری عبارت از جذر مربع انحراف می باشد.

$$\text{انحراف معیاری: } sd = \sqrt{25,000} = 158,1$$



در نتیجه با استفاده از انحراف معیاری، گفته می‌توانیم که کدام یک از مشاهدات (قد‌ها) نورمال، کدام یک بسیار بلند و کدام آنها بسیار کوتاه می باشد.

مثال: یک نمونه دارای 5 مشاهده، که معلومات در مورد معاشات کارمندان دو وزارت خانه دولتی می باشد.

معاشات پنج نفر کارمند در دو وزارت خانه، حمل 1395

نام کارمندان	تجارت و صنایع (افغانی)	مخابرات و تکنالوژی (افغانی)
احمد	100,000	2,200
رشید	500	300
فرید	2,000	450
رحیم	40,000	1,000
حامد	3,200	2,500
(Std. Dev.) انحراف معیاری	42,921	1,008

شکل: ۳-۳ جدول معاشات ۵ کارمند وزارت تجارت و صنایع و مخابرات و تکنالوژی

حالا گفته می‌توانیم که معاشات کارمندان وزارت تجارت و صنایع نسبت به معاشات کارمندان وزارت مخابرات و تکنالوژی پراکنده تر می باشند. یا به عبارت دیگر معاشات کارمندان وزارت تجارت نسبت به وزارت مخابرات و تکنالوژی بیشتر پراکنده است.



تمرین

1. ضرورت دانستن معیارهای پراکندگی را توضیح داده و انواع آن را نام بگیرید؟
2. چرا دامنه تغییرات نمیتواند یک معیار خوب برای نشان دادن پراکندگی یک سلسله ی اعداد باشد؟
3. جدول ذیل نشان دهنده سن 11 نفر از مرد ها و زنان ولایت کابل میباشد، با استفاده از انحراف معیاری دریابید که پراکندگی در میان کدام گروه بیشتر میباشد؟

سن مرد ها	سن خانم ها
18	18
24	15
26	28
22	79
19	33
45	17
35	40
32	30
33	37
40	22
60	60

شکل: ۳-۴ جدول سن ۱۱ نفر مرد و زن ولایت کابل

4. معاش مامورین دیپارتمنت های تدارکات و دیپارتمنت تحقیق وزارت اقتصاد که هر یک دارای 11 نفر کارمند اند به ترتیب ذیل میباشد، با استفاده از انحراف معیاری دریابید که پراکندگی در معاشات کدام دیپارتمنت زیاد است؟

معاشات کارمندان دیپارتمنت های تدارکات و تحقیق وزارت اقتصاد، در ماه حمل 1395

شماره	نام کارمندان	تدارکات (افغانی)	تحقیق (افغانی)
1	احمد	12,000	20,000
2	فرهاد	15,000	20,000
3	نصرت	10,000	21,000
4	صادق	10,000	19,000
5	احمد شفیق	9,000	19,500
6	محد زمان	13,000	21,000
7	محمد امان	17,000	20,500
8	احمد فواد	12,000	19,300
9	احمد جمشید	12,000	21,000
10	احمد رشاد	15,000	20,000
11	محمد مسعود	9,000	21,500

شکل: ۳-۵ جدول معاشات ۱۱ نفر کارمند دو ریاست وزارت اقتصاد

بخش چهارم

فرضیه (Hypothesis) چیست؟

تصور کنید که شما فردا یک امتحان دارید. شما شب قبل تا ناوقت در بیرون از منزل هستید و تلویزیون می بینید. شما میدانید که اگر روز قبل از امتحان درس بخوانید در نتیجه نمره بهتر را در امتحان فردا خواهید گرفت. به نظر شما نتیجه امتحان فردای شما چگونه خواهد بود با فرض اینکه شب قبل از امتحان تا ناوقت شب در بیرون از منزل بودید؟



زمانیکه شما به این سوال پاسخ میدهید در واقع شما یک فرضیه را شکل میدهید. یک فرضیه عبارت از پیش بینی قابل اندازه گیری و مشخص می باشد. و فرضیه عبارت از ارایه علمی پیش بینی شما می باشد، توسط یک فرضیه شما پیش بینی میکنید که چه اتفاق خواهد افتاد.

در حالت مشخص فرضیه، شما میتوانید به شکل ذیل باشید: " اگر من امشب درس نخوانم، پس نمره پائین تر را در امتحان فردا خواهم گرفت."

در شکل ذیل پائین مثال خوبی برای فهم بهتر فرضیه داده شده است دو فرد هستند که هر کدام آنها با دیدن دود از فاصله دورتر فرضیات متفاوت دارند. یکی از آنها فکر میکند که در قریه آتش سوزی شده و فرد دیگر به این باور است که گویا دیگران برای آماده ساختن غذا روشن شده است.

شکل: ۴-۱ نمایش فرضیه

هدف یک فرضیه

معمولاً دانشمندان در یک بخش مشخص علوم، بعد از مشاهدات و اندیشیدن به نتایج میرسند که تا قبل از انجام آزمایش‌ها و ثابت شدن یا نشدن آن یک فرضیه است. فرضیه، حدس و گمان است که هنوز به اثبات نرسیده و با استفاده از احصاییه و تست‌های احصائیوی میتواند یک فرضیه به اثبات برسد یا اینکه به اثبات نرسد.

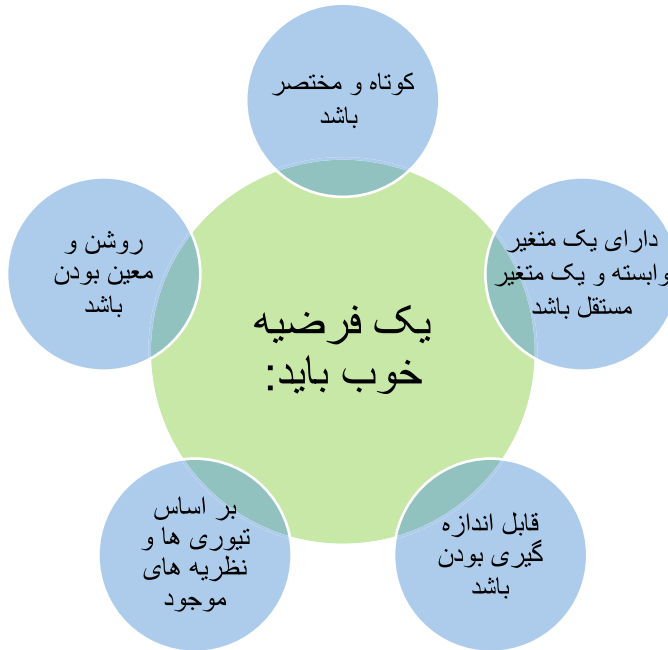
هدف از استفاده فرضیه در یک تحقیق این است که رابطه بین دو پدیده یا به زبان تحقیق رابطه بین دو متغیر، یک متغیر وابسته و دیگری متغیر مستقل را بررسی و مطالعه کند. یک فرضیه به ما کمک می‌کند تا اینکه چه نتیجه را در یک تحقیق باید به بررسی بگیریم و تمام تمرکز تحقیق بالای آن می‌باشد.

تغییرات در متغیر وابسته، بستگی به تغییرات در متغیر مستقل دارد. یعنی متغیر وابسته در پی تغییرات در متغیر آزاد از خود تغییر نشان می‌دهد. لازم به یاد آوری است که متغیر وابسته را بنام متغیر نتیجه (Outcome Variable) نیز یاد میکنند. و در مثال بالا نتیجه امتحان عبارت از متغیر وابسته است زیرا نتیجه آن وابسته به متغیر مستقل است. در مثال بالا متغیر مستقل عبارت از درس خواندن در شب قبل از امتحان می‌باشد.

باز هم میتوان از مثال بالا استفاده کرد و به این نکته اشاره کرد که یک فرضیه باید قابل آزمایش باشد، چونکه نمره امتحان شما قابل سنجش است. شما میتوانید نمرات امتحان تان را که برای آن درس نخوانده اید و نمرات امتحان را که شب قبل برای آن درس خوانده اید مقایسه کنید.

ویژگی‌های یک فرضیه خوب (Characteristics of a Good Hypothesis)

- دارای یک متغیر وابسته و یک متغیر مستقل باشد.
- قابل اندازه‌گیری بودن باشد.
- روشن و معین بودن باشد.
- بیان بر اساس تیوری‌ها و نظریه‌های موجود.
- کوتاه و مختصر باشد.



انواع فرضیه: (Types of Hypothesis)

1. فرضیه تحقیق یا متناوب (H_A) یا (H_1) - Alternative Hypothesis

2. فرضیه صفر (H_0) - Null Hypothesis

منظور از فرضیه تحقیق، فرضیه ای است که در آن به بیان تفاوت بین متغیرها (رابطه) پرداخته می شود، ولی فرضیه صفر بیان گر عدم وجود رابطه بین متغیرها است.

مثال:

فرضیه تحقیق: افزایش فساد اداری باعث این میشود که رشد اقتصادی در یک اقتصاد کاهش پیدا کند.

فرضیه صفر: میان سطح فساد اداری و میزان رشد اقتصادی هیچ رابطه ی وجود ندارد.

انکشاف فرضیه (Developing Hypothesis)

در هنگام انکشاف یا ساختن فرضیه نکات مهم که در پائین به آن اشاره میشود را در نظر گرفت:

- فرضیه باید به صورت یک جمله خبری، روشن و بدون ابهام بیان شود. چون بیان فرضیه به شکل خبری قابل فهم تر است.

مثال: شاگردانیکه باهوش تر اند فرصت بیشتر برای ادامه تحصیل دارند.

- فرضیه باید رابطه بین یک یا چند متغیر را بیان کند. رابطه معمولاً به شکل مثبت و مستقیم بیان

میشود هر چند در بعضی مواقع رابطه به شکل معکوس نیز بیان میشود. مثال که در بالا به آن اشاره شد نشان دهنده رابطه مستقیم و مثبت است.

- فرضیه باید قابل آزمون باشد. فرضیه باید طوری انتخاب شود که بتوان با استفاده از دیتا و مشاهدات آنها را رد یا تایید کرد. در صورت تایید فرضیه بتوان نتایجی را از آن پیش بینی کرد.
- در نوشته های علمی کوتاه بودن و واضح بودن همیشه دارای اهمیت می باشد. در نتیجه یک فرضیه باید کوتاه و واضح باشد.

فرضیه نامناسب: اثر تدریس به شیوه سخنرانی و پرسش و پاسخ در یادگیری شاگردان کمتر از روش شاگرد محوری است.

در فرضیه بالا اختصار نوشته کردن مراعات نشده است.

فرضیه مناسب: تاثیر تدریس به روش سخنرانی در یادگیری شاگردان کمتر از روش آزمایشی است.

مثال اول:

در یک آزمایش که در لابراتوار ادویه جات وزارت صحت عامه بالای یک داروی جدید صورت گرفته:

فرضیه صفر (H_0): ادویه که جدید انکشاف و تولید شده است، از نگاه تاثیرگذاری با ادویه که در گذشته ساخته شده تفاوت ندارد.

فرضیه تحقیق (H_A): ادویه که جدید تولید شده، تاثیر گذاری آن متفاوت تر از ادویه همانند است که در گذشته موجود بوده است.

مثال دوم:

فرضیه صفر (H_0): روشن بودن چراغ های اتاق در شب و بیدار ماندن یک شاگرد تا دیروقت در شب باهمدیگر ارتباط ندارند.

فرضیه تحقیق (H_A): روشن ماندن چراغ های اتاق در شب، سبب بیدار ماندن شاگردان تا دیروقت در شب می شود.



تمرین

1. فرضیه چیست در باره آن معلومات دهید؟
2. در مورد کارهای روزمره وظیفه تان فکر نموده و نظر به تصور خودتان یک فرضیه بسازید؟
3. یک فرضیه بسازید و هم در فرضیه خویش متغیر مستقل و وابسته را نشان داده و قابل اندازه گیری بودن متغیرها را مشخص سازید؟
4. یک فرضیه متناوب و یک فرضیه صفری را انکشاف دهید؟
5. هدف یک فرضیه چیست با ذکر یک مثال توضیح دهید؟

آزمون های احصائیوی

آزمون همبستگی (Correlation)

اصطلاح همبستگی رابطه بین دو متغیر را نشان میدهد. به عنوان مثال زمانیکه والدین فرزندان شان را تشویق به درس خواندن بیشتر می نمایند،

در حقیقت آنها بر این باور هستند که تلاش بیشتر فرزندان شان باعث افزایش نمرات آنها در امتحانات میشود. به عنوان مثال این باور نزد مردم وجود دارد که میان افزایش سن افراد و فشار خون شان رابطه وجود دارد.

تعریف: همبستگی یک تحلیل احصائیوی است که هدف آن مطالعه رابطه بین دو متغیر می باشد، به عنوان مثال رابطه بین قد پدر و پسر، بارش باران و حاصلات زراعتی، معاش یک کارمند و کیفیت انجام کار.

تصویر که رسامی شده است میتواند نشان دهنده رابطه بین ثروت و خوشحالی باشد. یکی از افراد که ثروتمند تر است خوشحال تر بوده و شخص دیگر ثروتمند نیست ناراضی به نظر می رسد. و هدف اساسی از آزمون همبستگی نشان دادن رابطه بین متغیر ها است که در این شکل ثروت یک متغیر و سطح رضایت از زندگی یا خوشحالی متغیر دیگر می باشد.



شکل: ۵-۱ ارتباط همبستگی میان ثروت و خوشحالی

تحلیل احصائیوی همبستگی به اندازه گیری شدت و چگونگی موجودیت رابطه بین دو متغیر می پردازد. لازم به یادآوری می باشد که همبستگی رابطه علت و معلول را نشان نمیدهد. به طور مثال در انگلستان معلومات در مورد اشخاصی که مبتلا به سرطان بودند جمع آوری شد و در نتیجه میان استفاده از سگرت و سرطان شش رابطه همبستگی مثبت بدست آمد. اما این ارتباط همبستگی مثبت به معنای رابطه علت معلولی میان این دو متغیر نیست شاید یک عامل سومی دیگر مانند آلودگی هوا باعث افزایش مرض سرطان شده باشد.

همبستگی مثبت و منفی (Negative & Positive Correlation): همبستگی مثبت و منفی دو متغیر وابسته به جهت حرکت دو متغیر است.

اگر هر دو متغیر هم جهت حرکت کنند به این معنی که با افزایش یک متغیر، متغیر دیگر نیز افزایش پیدا کند و یا اینکه با کاهش یک متغیر، متغیر دیگر نیز کاهش پیدا کند در این صورت رابطه بین دو متغیر

مثبت و یا مستقیم گفته میشود. از مثال همبستگی مثبت میتوان به رابطه مستقیم بین افزایش قد و وزن، حاصلات زراعتی و بارش باران، نمرات مضمون ریاضی و احصائیه اشاره کرد.

اگر دو متغیر خلاف جهت همدیگر حرکت کنند، به این معنی که با افزایش یک متغیر، متغیر دیگر کاهش

پیدا کند و یا اینکه با کاهش یک متغیر متغیر دیگر افزایش پیدا کند، در این صورت **ارتباط همبستگی منفی** یا **غیر مستقیم** بین دو متغیر وجود دارد. مانند مقدار محصولات زراعتی در بازار و قیمت آنها، در بازار از مثال همبستگی منفی می باشد.

ضریب همبستگی (Correlation Coefficient) شدت رابطه بین دو متغیر توسط ضریب همبستگی نشان داده میشود و سمبول ضریب همبستگی عبارت از (r) بوده که قیمت -1 الی $+1$ را گرفته می تواند. در هنگام تفسیر ضریب همبستگی نکات ذیل دارای اهمیت می باشد:

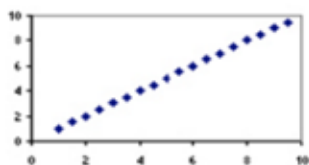
الف: هرگاه ضریب همبستگی $r = +1$ در اینصورت ارتباط کامل و مثبت بین دو متغیر وجود دارد.

ب: هرگاه ضریب همبستگی $r = -1$ باشد در اینصورت ارتباط همبستگی کامل و منفی بین دو متغیر وجود دارد.

ج: هرگاه ضریب همبستگی $r = 0$ در این صورت هیچ ارتباط همبستگی بین دو متغیر وجود ندارد.

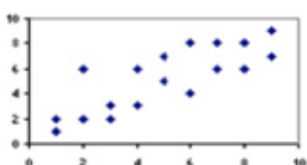
د: هرگاه قیمت مطلقه ضریب همبستگی بزرگتر از 0.5 ($r > 0.5$) باشد در این صورت ارتباط همبستگی بین دو متغیر قوی می باشد.

در شکل پائین انواع مختلف ارتباط همبستگی بین متغیر ها را دیده میتوانید.



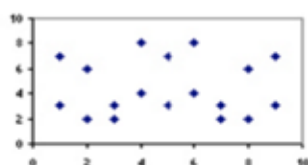
$(r = 1.0)$

ارتباط همبستگی کامل مثبت



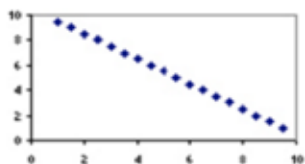
$(r = 0.80)$

ارتباط همبستگی قوی مثبت



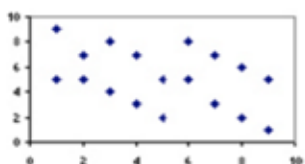
$(r = 0)$

ارتباط همبستگی صفر است



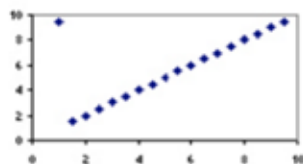
$(r = -1.0)$

ارتباط همبستگی منفی کامل



$(r = -0.43)$

ارتباط همبستگی منفی ضعیف



$(r = 0.71)$

ارتباط همبستگی قوی مثبت

فرمول ارتباط همبستگی:

$$r = \frac{\sum XY}{\sqrt{\sum X^2 \sum Y^2}} \quad , \quad X = x - \bar{x} \quad , \quad Y = y - \bar{y}$$

x عبارت از هر مشاهده جداگانه متغیر مستقل یا

y عبارت از هر مشاهده جداگانه متغیر وابسته یا می باشد.

مثال: ریاست زراعت ولایت پروان برای بررسی اینکه تاثیر افزایش کود اصلاح شده زراعتی و حاصلات زمین چگونه رابطه دارد، از آزمون همبستگی استفاده کرده و نتایج آن در پائین به تفسیر گرفته شده:

حاصلات زمین به تن $y =$

مقدار کود زراعتی به کیلوگرام $x =$

معلومات در مورد کود اصلاح شده زراعتی و حاصلات یک قطعه زمین در ولایت پروان

XY	y^2	$Y = y - \bar{y}$ $X = x - 68$	x^2	$X = x - \bar{x}$ $X = x - 67$	y	x
6	4	-2	9	-3	66	64
2	1	-1	4	-2	67	65
3	9	-3	1	-1	65	66
0	0	0	0	0	68	67
2	4	2	1	1	70	68
0	0	0	4	2	68	69
12	16	4	9	3	72	70
25	34	0	28	0	476	469

شکل: ۲-۵ معلومات در مورد کود اصلاح شده زراعتی و حاصلات زمین

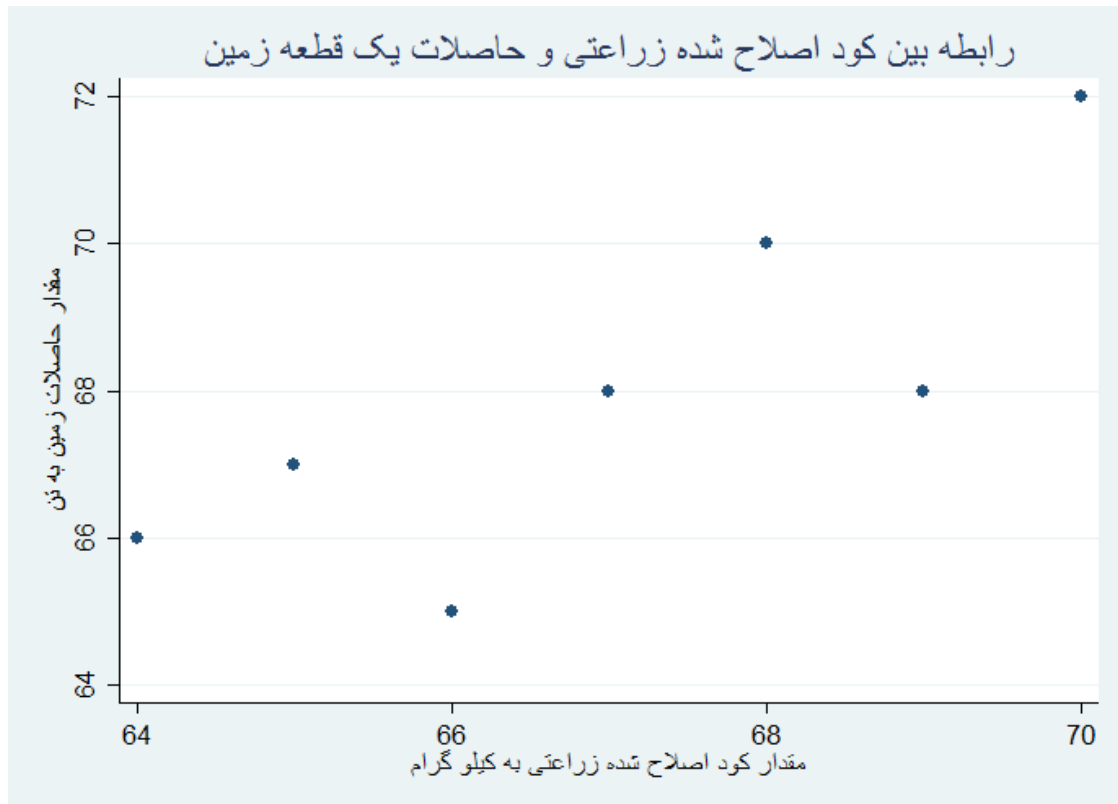
فرمول:

$$r = \frac{\sum XY}{\sqrt{\sum X^2 \sum Y^2}} \quad , \quad X = x - \bar{x} \quad , \quad Y = y - \bar{y}$$

$$\bar{x} = \frac{469}{7} = 67 \quad , \quad \bar{y} = \frac{476}{7} = 68$$

$$r = \frac{\sum XY}{\sqrt{\sum X^2 \sum Y^2}} = \frac{25}{\sqrt{28 * 34}} = \frac{25}{\sqrt{952}} = \frac{25}{30.85} = 0.81$$

تفسیر ضریب همبستگی (Correlation Coefficient): در این مثال بعد از سنجش ضریب همبستگی رابطه بین کود اصلاح شده زراعتی و حاصلات زمین مثبت بدست آمده، یعنی با افزایش کود اصلاح شده زراعتی حاصلات زمین نیز افزایش پیدا میکند. چون ارزش (r) علامه آن مثبت و ارزش آن بزرگتر از 0.5، یعنی $+0.8$ می باشد، علامه مثبت آن نشانه دهنده رابطه مستقیم و ارزش مطلقه آن نشان دهنده قوت و شدت رابطه می باشد.



شکل: ۳-۵ رابطه همبستگی میان کود اصلاح شده زراعتی و حاصلات زمین زراعتی

مثال: معلومات در مورد درجه حرارت و فروشات شرکت آیسکریم داده شده است، با استفاده از ارتباط همبستگی بگویید که چگونه رابطه بین فروشات آیسکریم و درجه حرارت وجود دارد؟

معلومات در مورد درجه حرارت و فروشات شرکت آیسکریم

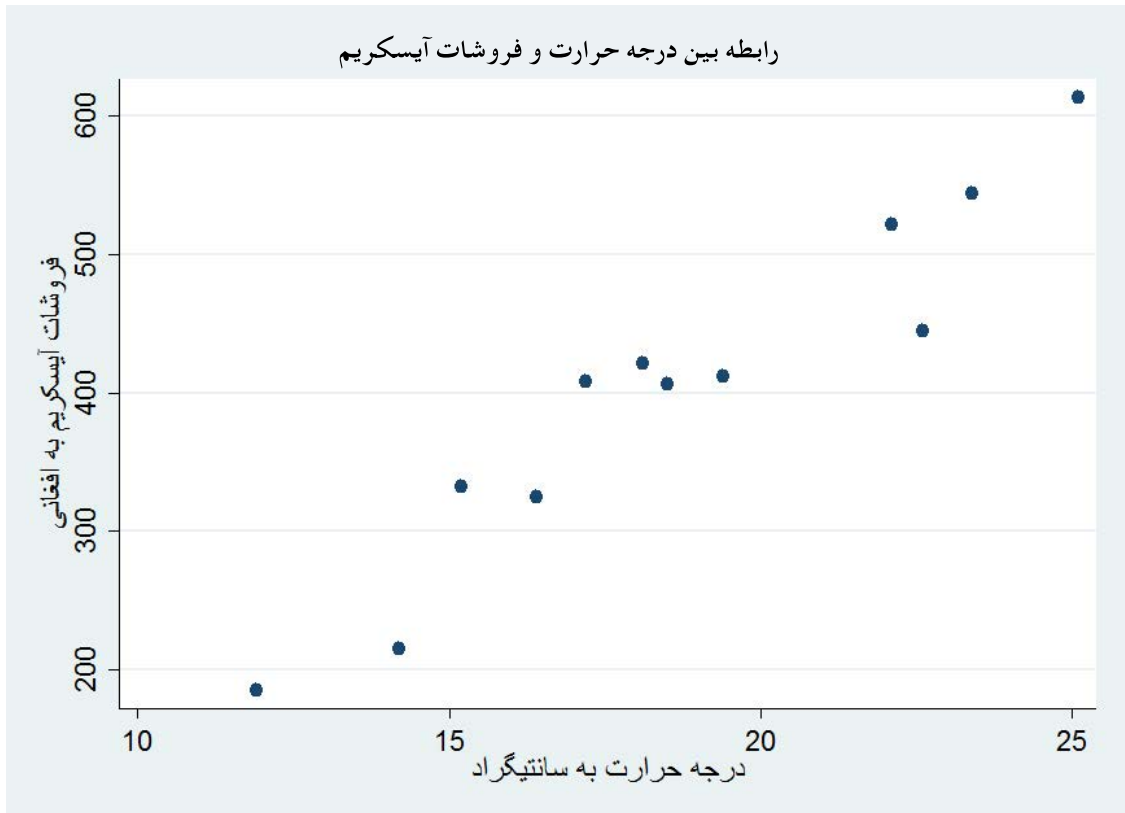
فروشات آیسکریم y = به افغانی	درجه حرارت x = به سانتیگراد
215	14.2
325	16.4
185	11.9
332	15.2
406	18.5

حل:

$$r = \frac{\sum XY}{\sqrt{\sum X^2 \sum Y^2}} = \frac{5,325}{\sqrt{177 * 174.757}} = 0.9575$$

شکل: ۵-۵ رابطه همبستگی میان درجه حرارت و فروشات آیسکریم

تفسیر ضریب همبستگی: در اینجا ضریب همبستگی عبارت از $(r=+0.95)$ می باشد که علامه مثبت آن نشان دهنده رابطه مستقیم بین درجه حرارت و فروشات آیسکریم دارد، یعنی به هر اندازه که درجه هوا گرم تر میشود در نتیجه فروشات آیسکریم نیز افزایش پیدا می کند، و چون ارزش ضریب همبستگی بزرگتر از 0.5 می باشد نشان دهنده رابطه قوی بین دو متغیر است.



شکل: ۵-۶ رابطه همبستگی میان درجه حرارت و فروشات آیسکریم

سطح معنی داری (P-Value): در هنگام مطالعه ارتباط همبستگی سطح معنی داری دارای اهمیت زیادی می باشد. سطح معنی داری به سمبول (P-Value) نشان داده شده و سطح معنی داری نشان دهنده این موضوع است که ضریب همبستگی به دست آمده دارای اعتبار می باشد و تصادفی و از روی چانس بدست نیامده است.

قیمت که (P-Value) میگیرد اگر کوچکتر از 0.05 ($P\text{-Value} < 0.05$) باشد، در این صورت رابطه همبستگی بین دو متغیر از نگاه احصائیوی معنی دار می باشد.

رگرسیون (Regression)

در بحث گذشته به مطالعه ارتباط همبستگی پرداختیم که برای بررسی چگونگی ارتباط بین دو متغیر بود، حالا به مطالعه رگرسیون می پردازیم که برای مطالعه موجودیت رابطه بین دو متغیر یا بیشتر از آن می پردازد و تفاوت آن با ارتباط همبستگی عبارت از این است که در ارتباط همبستگی متغیرها به متغیرهای وابسته و مستقل تقسیم بندی نمی شود، حال آنکه در رگرسیون متغیرها به دو گروه تقسیم میشوند، یکی متغیر مستقل و دیگری وابسته و رگرسیون به مطالعه رابطه بین متغیرهای مستقل و وابسته می پردازد.

در رگرسیون متغیرها به دو دسته تقسیم میشوند، یکی متغیرهای مستقل (Independent Variables) که مستقلانه قیمت میگیرند و دیگر متغیرهای وابسته (Dependent Variables) که قیمت آنها وابسته به قیمت متغیرهای مستقل است.

رگرسیون که یک تحلیل احصائی می باشد برای ساختن مدل های تحلیلی مورد استفاده قرار می گیرد. در هنگام ساختن مدل ها، متغیرها به یک متغیر وابسته و حداقل یک متغیر مستقل تقسیم میشوند. بعد از ساختن مدل و انجام تحلیل رگرسیون معلوم می گردد که تغییرات در متغیر مستقل چگونه باعث تغییر در متغیر وابسته میشود.

در اینجا از انواع مختلف رگرسیون تنها دو نوع آن را به مطالعه می گیریم.

1. رگرسیون حداقل مربعات (OLS regression or Ordinary Least Square)

2. رگرسیون پروبیت (Probit Regression)



نکته: رگرسیون یکی از مباحث احصائی است که در مورد آن کتاب و معلومات بسیار موجود است و در اینجا تنها مباحث که برای ما در هنگام استفاده برنامه ستاتا (Stata) دارای اهمیت می باشد مورد مطالعه قرار می گیرد.

رگرسیون به طریقه حداقل مربعات (Ordinary Least Square or OLS Regression) برای آشنایی بیشتر با رگرسیون، در ابتدا یک مدل تحلیلی ساخته و متغیرها را به دو دسته مستقل و وابسته تقسیم می کنیم. به مثال ذیل توجه کنید!



معادله رگرسیون به شکل ذیل می باشد:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + u$$

در معادله رگرسیون، (x_1, x_2, \dots, x_n) عبارت از متغیرهای مستقل می باشد که قیمت های مستقل را به خود میگیرند و (y) عبارت از متغیر وابسته می باشد که قیمت های آن وابسته به قیمت متغیرهای مستقل می باشد.

تحلیل رگرسیون در واقع تاثیرات متغیرهای مستقل را بالای متغیر وابسته طوری مطالعه میکند که تاثیرات متغیرهای مستقل دیگر ثابت بماند. به عبارت دیگر، ضریب های B برای هر متغیر مستقل، تاثیر خالص همان متغیر مستقل بالای متغیر وابسته می باشد. برای محاسبه B یکی از طریقه ها طریقه حداقل مربعات می باشد. در برنامه (Stata) توسط فرمان (reg) می باشد که در بخش های قبلی در مورد آن به تفصیل معلومات داده شده است که مطالعه آن قبل از مطالعه این بخش ضروری می باشد.

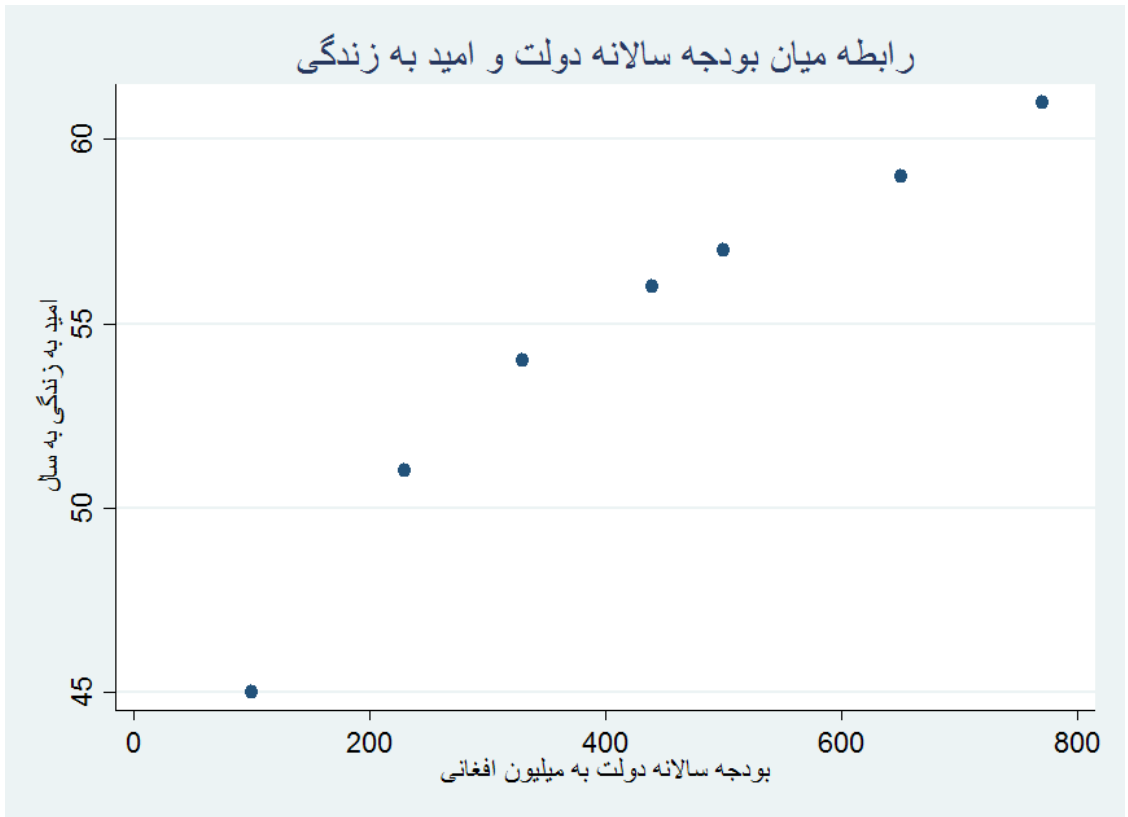
فرضیه های رگرسیون خطی (OLS Regression Assumptions)

در این بخش به تشریح فرضیه های رگرسیون خطی می پردازیم. بخاطر داشته باشید که ابتدا باید فرضیه رگرسیون خطی در نظر گرفته شده و بعد از آن فرمان رگرسیون انجام شود و از نتایج به دست آمده تفسیر رابطه بین متغیرها صورت گیرد.

1. فرضیه رابطه خطی بین متغیر وابسته و متغیرهای غیر وابسته (Linearity):

در رگرسیون خطی چند متغیره یکی از فرضیه ها این است که رابطه بین متغیرهای مستقل و وابسته باید خطی

باشد. و برای این کار میتوان از نمودار پراکندگی (Scatter-plot) استفاده کرد.



شکل: ۷-۵ رابطه میان بودجه خدمات صحتی و امید به زندگی

3. تغییر پذیری در تمام متغیرها (variability in all variables)

این موضوع بسیار حایز اهمیت است که هر یک از متغیرها باید تغییرپذیر باشند. زیرا اگر تغییری از ابتدا تا انتها قیمت یکسان و تغییر نا پذیر داشته باشد، در این صورت متغیر گفته نمی شود و جا دادن چنین یک متغیر در رگرسیون هیچ معنایی نخواهد داشت. پس تغییر پذیر بودن قیمت ها یا مشاهدات موجود در یک متغیر یکی از اساسی ترین شرط ها برای رگرسیون است.

مثال: اگر عاید ماهانه مد نظر باشد، باید چند عاید ماهانه با ارزشهای متفاوت داشته باشیم تا اینکه در مدل شامل گردد.

(100 افغانی، 230 افغانی، 450 افغانی و) این عاید چون دارای چند عاید با ارزش های متفاوت است، میتواند در مدل شامل گردد.

(100 افغانی، 100 افغانی، 100 افغانی) این عاید ماهانه به عنوان یک متغیر در مدل شامل شده نمیتواند، چون

که تمام ارزش های آن 100 افغانی می باشد و تغییر را نشان نمیدهد و یا واریانس ندارد.

4. نمونه تصادفی (Random Sample)

نمونه که از جمعیت گرفته میشود و از آن متغیرها برای ساختن مدل رگرسیون استفاده میشود، باید به شکل تصادفی انتخاب شده باشد، و معلومات جمع آوری گردیده باشد.

5. مولتی کولینریتی "هم خطی چندگانه" (Multicollinearity)

مولتی کولینریتی به سخن ساده زمانی وجود دارد که دو و یا بیشتر از دو متغیر مستقل با هم رابطه قوی همبستگی داشته باشند. برای آزمون این فرضیه در پروگرام ستاتا (Stata) میتوانیم از فرمان (pwcorr) استفاده کنیم، در صورتی که ضریب همبستگی بزرگتر از 0.8 باشد، باید یکی از دو متغیر مستقل از مدل رگرسیون بیرون کشیده شود.

مثال: در این مثال دو متحول مستقل (قد و سن) را مدنظر میگیریم

در این مدل میخواهیم بررسی کنیم که آیا رابطه بین این دو متحول مستقل وجود دارد یا خیر. برای انجام این کار از آزمون همبستگی استفاده میکنیم و رابطه همبستگی بین دو متحول مستقل قد و سن را بررسی میکنیم. بعد از آزمون همبستگی بین دو متغیر به این نتیجه می رسیم که ارتباط همبستگی قوی میان این دو متغیر وجود داشته است و ضریب همبستگی ($r=0.894$) بزرگتر از 0.8 است که در نتیجه یکی از این متغیرهای مستقل را از مدل رگرسیون ناگزیر بیرون می سازیم.

`pwcorr height age, sig`

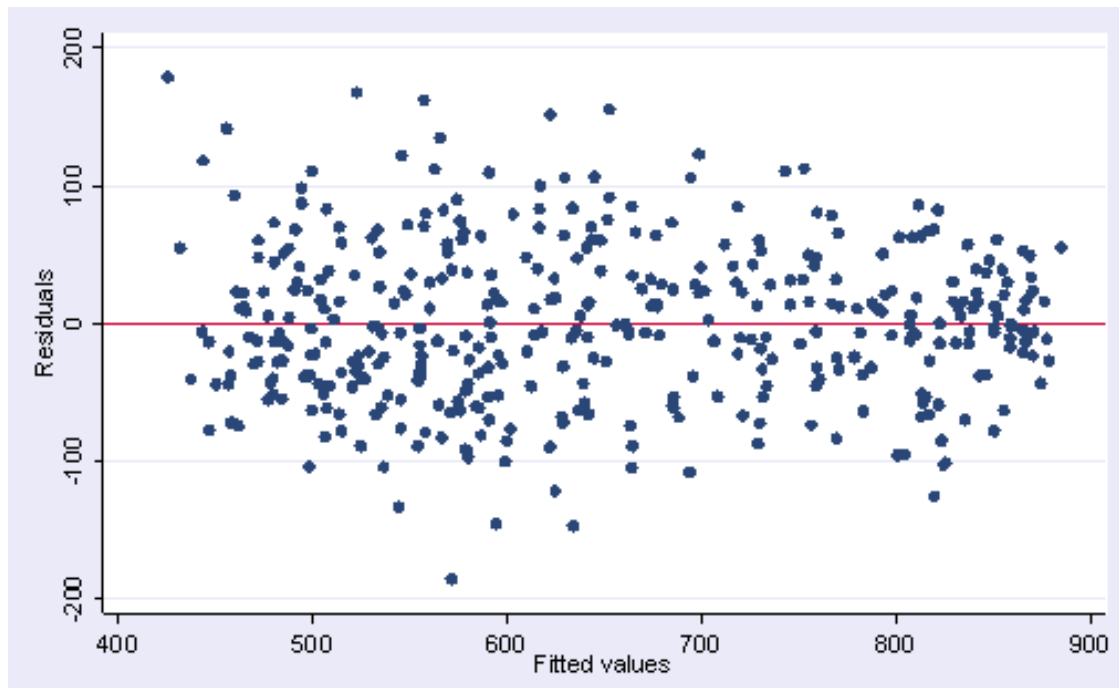
	height	age
height	1.0000	
age	0.8944 0.0066	1.0000

6. فرضیه برابری انحراف ها (Homoscedasticity):

یکی دیگر از فرضیه های رگرسیون طریقه حداقل مربعات عبارت از فرضیه برابری وریانس خطاها می باشد. خطاها عبارت از تفاوت بین قیمت های واقعی (y) و قیمت های تخمین شده آن به صورت اوسط به اساس قیمت های (x) می باشد.

برای انجام این آزمون در برنامه ستاتا از فرمان ذیل استفاده میکنیم. البته به خاطر داشته باشید که در مرحله اول باید رگرسیون اجرا گردد و بعد از آن فرمان برای انجام آزمون برابری واریانس خطاها انجام شود.

rvfplot, yline(0)



شکل: ۹-۵ فرضیه برابری انحراف خطاها

در شکل بالا دیده میشود که نقطه ها در اطراف (بالا و پائین) خط سرخ با واریانس تقریبی مساوی پراکنده می باشند و فرضیه درست است. هرگاه نقطه ها در اطراف خط با واریانس تقریبی مشابه به حرکت نباشد، در این صورت مشکل نابرابری انحراف خطاها (Heteroskedasticity) پدید می آید.

داده های خارج از محدوده (Outlier) یکی از دلایل پدید آمدن مشکل انحراف خطاها می باشد که نتیجه رگرسیون چند متغیره خطی را تحت تاثیر قرار میدهد. یک مثال میتواند که برای فهم موضوع کمک کند.

به گونه مثال وزارت زراعت معلومات را در ارتباط به محصولات زراعتی پنج ولایت کشور در سال 1393 به نشر رسانده اما محصولات زراعتی یک ولایت نسبت به ولایت های دیگر تفاوت بسیار دارد و از حالت معمول متفاوت است. در جدول ذیل به محصولات زراعتی ولایت لوگر نگاه کنید.

محصولات گندم پنج ولایت افغانستان در سال 1393

نام ولایت	محصولات به تن
کابل	200
غزنی	230
میدان وردک	130
لوگر	38,000
پنجشیر	234

شکل: ۵-۸ داده های خارج از محدوده

7. متغیر وابسته پیوسته (Continuous Dependent Variable)

در مدل رگرسیون به طریقه ای حد اقل مربعات که ترتیب می گردد، باید متغیر وابسته یک متغیر گسسته باشد. از مثال متغیر پیوسته میتوان به قد درآمد نمرات امتحان فروشات اشاره کرد.

اما متغیرهای مستقل میتوانند که پیوسته یا دارای کتگوری (Categorical Variable) باشد، مثال های متغیر پیوسته در بالا ذکر گردید و مثال های متغیر دارای کتگوری عبارت اند از: جنسیت (مرد یا زن)، ملیت (تاجیک، هزاره، پشتون و ...)، محل زندگی (شهر یا دهات).



نکته: در اینجا به چند مثال رگرسیون پرداخته و موضوع مهم تفسیر نتیجه رگرسیون می باشد.

مثال: در این مثال اول یک رگرسیون ساده یعنی متشکل از دو متغیر، یکی وابسته و دیگری مستقل می باشد، به عنوان مثال میخواهم بدانیم که آیا با افزایش سطح تحصیل رضایت مردم از دموکراسی افزایش پیدا می کند یا خیر؟

برای رسیدن به این نتیجه مراحل ذیل ضروری می باشد:

1. دانستن کودها و لیبل های متغیرها
2. انجام فرمان رگرسیون
3. تفسیر نتیجه رگرسیون

1. تحلیل کودهای متغیرها:

متغیر (q32) نشان دهنده سطح رضایت مردم از دموکراسی می باشد که کودهای آن از 1 الی 4 می باشد. که از کود 1 شروع شده به طرف کود 2 رفته و به همین ترتیب تا کود 4 نشاندهنده رضایت کمتر از دموکراسی می

باشد.

متغیر (d10) دارای کدهای 1 الی 13 می باشد و با هر واحد از کد 1 الی 13، نشان دهنده تحصیل بیشتر می باشد.

```
. labellist q32 d10 m6b d1
Q32:
      1 Very satisfied
      2 Somewhat satisfied
      3 Somewhat dissatisfied
      4 Very dissatisfied
     98 Refused
     99 Don't Know

D10:
      1 Never went to a school
      2 Informal schooling at home or at a literacy class
      4 Primary School, incomplete (classes 1 to 5)
      5 Primary School, complete (finished class 6)
      6 Secondary education, incomplete (classes 7 to 8)
      7 Secondary education, complete (finished class 9)
      8 High School incomplete (classes 10-11)
      9 High School complete (finished class 12)
     10 14th grade incomplete (class 13)
     11 14th grade complete (finished class 14)
     12 University education incomplete (have no degree diploma)
     13 University education complete (have degree diploma)
     98 Refused (vol.)
     99 Don't know (vol.)

M6B:
      1 Rural
      2 Urban

D1:
      1 Male
      2 Female
```

2. اجرای فرمان رگرسیون در برنامه (Stata): اجرای فرمان رگرسیون به شکل ذیل می باشد،

```
reg [dependent variable] [independent variables]
```

```
reg q32 d10 if q32<98
```

3. تفسیر نتیجه فرمان رگرسیون بعد از اجرای آن در برنامه (Stata)

. reg q32 d10 if q32<98						
Source	SS	df	MS	Number of obs	=	9,468
Model	.073044591	1	.073044591	F(1, 9466)	=	0.09
Residual	7823.77349	9,466	.826513151	Prob > F	=	0.7663
				R-squared	=	0.0000
				Adj R-squared	=	-0.0001
Total	7823.84654	9,467	.826433562	Root MSE	=	.90913

q32	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
d10	.0007893	.0026551	0.30	0.766	-.0044152	.0059938
_cons	2.399013	.013428	178.66	0.000	2.372691	2.425335

بعد از اجرای فرمان رگرسیون برنامه ستاتا نتیجه که در بالا ملاحظه میکند را نشان میدهد. حالا نکات که برای تحلیل ما دارای اهمیت می باشد در اینجا به ترتیب به تفسیر میگیریم.

1. **ضریب رگرسیون (Regression Coefficient):** در مدل های رگرسیون ساده تنها یک ضریب رگرسیون وجود دارد.

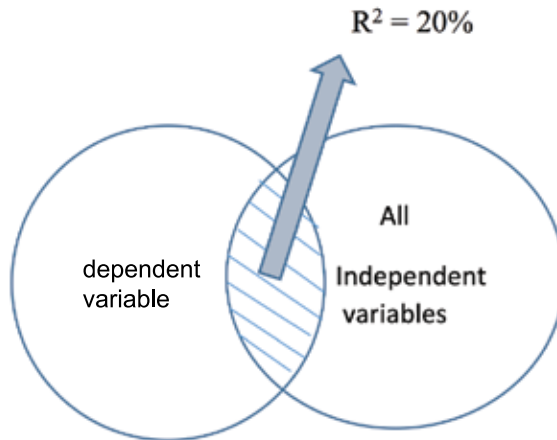
این ضریب نشان دهنده آن است که با یک واحد یا یک سال افزایش در سطح تحصیل به اندازه B) $=0.00078$ واحد رضایت از دموکراسی کاهش پیدا می کند. یا به تفسیر ساده تر شهروندان تحصیل کرده نسبت به افراد که تحصیل کمتر دارند رضایت کمتر از دموکراسی نشان داده اند.

2. **$P>|t|$ یا P-value** همانگونه که در مبحث ارتباط همبستگی مطالعه کردیم نشان دهنده سطح معنی داری می باشد و اگر قیمت آن از 0.05 کوچکتر باشد رابطه بدست آمده معنی دار می باشد. اما در این مثال قیمت آن بزرگتر از 0.05 است، پس در نتیجه رابطه بدست آمده معنی دار نیست و نیاز به تفسیر آن نیست اما در بالا تنها بخاطر کمک برای یاد گیری نتیجه رگرسیون قیمت B تفسیر شده است.

3. **R-squared:** نام دیگر آن (Coefficient Determination) ضریب تشخیص می باشد، که نشان دهنده توانایی مدل در توضیح دادن تغییرات در متغیر وابسته می باشد. و قیمت آن در اینجا عبارت از (0.000) می باشد به معنی اینکه تقریباً این مدل به اندازه 0% توانایی تغییرات در متغیر وابسته (رضایت از دموکراسی) را توسط متغیر مستقل (سطح تحصیل) دارا می باشد.



نکته: با اینکه ضریب تشخیص در مثال گذشته تقریباً 0% بود، اما اگر ضریب تشخیص 20% باشد، میتوان آنرا به شکل ذیل توسط یک شکل گرافیکی به منظور فهم بهتر نشان داد.



مثال 2. در این مدل که در اینجا به آن می پردازیم برعلاوه تحصیل متغیرهای مستقل دیگر را نیز در مدل شامل می سازیم تا رابطه آنها را نیز با متغیر وابسته (رضایت از دموکراسی) به مطالعه بگیریم. مانند: سن (d2)، جنسیت (d1)، عاید ماهانه (d18a) و شهری یا دهاتی بودن (m6b).

برای دانستن روابط بین متغیرهای مستقل و وابسته در این مدل باید مانند قبل سه مرحله را اجرا کنیم.

1. دانستن کودها و لیبل های متغیرها

2. انجام فرمان رگرسیون

3. تفسیر نتیجه رگرسیون

1. دانستن کود و لیبل متغیرها

- در متغیر (m6b) با حرکت از کود 1 به 2 یعنی نظر شهری های نسبت به یک موضوع، در این مثال رضایت شهری ها از دموکراسی.
- (d1) متغیر است که دارای دو کود 1 و 2 می باشد و حرکت از کود 1 به 2، یعنی نظر زنان در مورد یک موضوع، در اینجا رضایت از دموکراسی.

Q32:

- 1 Very satisfied
- 2 Somewhat satisfied
- 3 Somewhat dissatisfied
- 4 Very dissatisfied
- 98 Refused
- 99 Don't Know

D10:

- 1 Never went to a school
- 2 Informal schooling at home or at a literacy class
- 4 Primary School, incomplete (classes 1 to 5)
- 5 Primary School, complete (finished class 6)
- 6 Secondary education, incomplete (classes 7 to 8)
- 7 Secondary education, complete (finished class 9)
- 8 High School incomplete (classes 10-11)
- 9 High School complete (finished class 12)
- 10 14th grade incomplete (class 13)
- 11 14th grade complete (finished class 14)
- 12 University education incomplete (have no degree diploma)
- 13 University education complete (have degree diploma)
- 98 Refused (vol.)
- 99 Don't know (vol.)

M6B:

- 1 Rural
- 2 Urban

D1:

- 1 Male
- 2 Female

2. انجام فرمان رگرسیون در برنامه (Stata)

reg [dependent variable] [independent variables]

reg q32 d10 d2 m6b d1 d18a if q32<98 & d18a>99

3. تفسیر نتیجه فرمان رگرسیون بعد از اجرای آن در برنامه (Stata)

```
. reg q32 d10 d2 m6b d1 d18a if q32<98 & d18a>99
```

Source	SS	df	MS	Number of obs	=	5,746
				F(5, 5740)	=	4.52
Model	19.2228333	5	3.84456665	Prob > F	=	0.0004
Residual	4880.17588	5,740	.850204857	R-squared	=	0.0039
				Adj R-squared	=	0.0031
Total	4899.39871	5,745	.852810916	Root MSE	=	.92207

q32	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
d10	-.0037196	.0036565	-1.02	0.309	-.0108878	.0034486
d2	.0015388	.0009915	1.55	0.121	-.000405	.0034826
m6b	.0410946	.0318054	1.29	0.196	-.021256	.1034452
d1	-.1044162	.0260869	-4.00	0.000	-.155564	-.053276
d18a	-7.95e-08	3.57e-07	-0.22	0.824	-7.80e-07	6.21e-07
_cons	2.481562	.0697972	35.55	0.000	2.344733	2.61839

1. تفسیر رابطه بین سن و رضایت از دموکراسی، یعنی با یک واحد (یک سال) افزایش در سن به اندازه 0.00153) واحد رضایت از دموکراسی کمتر می شود. تفسیر رابطه بین شهری بودن و دهاتی بودن، در اینجا رابطه بین این دو متغیر مستقیم می باشد، یعنی با افزایش هر واحد در متغیر شهری یا دهاتی بودن (m6b) باعث افزایش در متغیر (q32) رضایت از دموکراسی به اندازه $B=0.04109$ واحد میشود، یعنی شهری ها بیشتر از کسانی که در روستاها زندگی میکنند، از دموکراسی رضایت دارند.

2. تفسیر رابطه بین جنسیت و رضایت از دموکراسی، در اینجا رابطه بین این دو متغیر غیر مستقیم می باشد، یعنی یک واحد افزایش در متغیر جنسیت (d1) باعث کاهش در متغیر (q32) رضایت از دموکراسی به اندازه $B=0.1044$ واحد میشود، یعنی زنان بیشتر از مردان از دموکراسی در کشور رضایت دارند.

3. R-squared: نام دیگر آن (Coefficient Determination) ضریب تشخیص می باشد، که نشان دهنده توانای مدل در توضیح دادن تغییرات در متغیر وابسته می باشد. و قیمت آن در اینجا عبارت از (0.0039) می باشد به معنی اینکه تقریباً این مدل به اندازه 0.39% توانایی تغییرات در متغیر وابسته (رضایت از دموکراسی) را توسط متغیرهای مستقل (سطح تحصیل، جنسیت، سن، عاید ماهانه و شهری و دهاتی) دارا می باشد.

4. در این مدل رابطه بین هیچ یک از متغیرهای مستقل با متغیر وابسته معنی دار نمی باشد به استثنای رابطه بین جنسیت و رضایت از دموکراسی که قیمت $P\text{-value} = 0.000$ از 0.05 کوچکتر می باشد، و متبانی آنها از 0.05 بزرگتر می باشند.

پروبیست رگرسیون (Probit Regression)

از دو نوع رگرسیون که مورد مطالعه ما قرار دارد، ابتدا رگرسیون به طریقه حداقل مربعات را به بررسی گرفتیم. در این جا رگرسیون پروبیست را به تشریح می گیریم. این نوع رگرسیون زمانی که متغیر وابسته در یک تحلیل رگرسیون دارای دو قیمت (Dummy variable) باشد مورد استفاده قرار می گیرد و رگرسیون به طریقه حداقل مربعات در این حالت استفاده نمی شود.

تفاوت بین رگرسیون به طریقه حداقل مربعات و پروبیست

رگرسیون حداقل مربعات	رگرسیون پروبیست
هنگامی استفاده میشود که متغیر وابسته یک متغیر پیوسته (Continuous)	هنگامی استفاده میشود که متغیر وابسته دارای دو قیمت 0 و 1 باشد.

شکل: ۵-۱۰ تفاوت میان رگرسیون به طریقه خطی و رگرسیون پروبیست

مثال: با استفاده از دیتای سروی مردم افغانستان در سال 2015، بنیاد آسیا می خواهیم یک مدل تحلیلی ساخته و رابطه بین جنسیت (d1)، سن (d2)، دهاتی و شهری بودن (m6b) و سطح تحصیل (d10) را با نظر مردم افغانستان در مورد اینکه کشور به کدام در حرکت است را به بررسی می گیریم. برای دریافت این رابطه سه مرحله را مانند رگرسیون حداقل مربعات انجام میدهیم:

1. دانستن کودها و لیبل های متغیرها

2. انجام فرمان رگرسیون

3. تفسیر نتیجه رگرسیون

1. دانستن کود و لیبل متغیرها

کود سایر متغیرها را از مثال های گذشته میدانید، در اینجا تنها کود متغیر وابسته (q1) که دارای دو کود 0 و 1 می باشد را مورد مطالعه قرار میدهیم، البته بخاطر داشته باشید که قبل از انجام فرمان رگرسیون پروبیست باید کودهای متغیر وابسته را به 0 و 1 تبدیل کنید که این موضوع به تفسیر گرفته شده در فصل های قبلی و مطالعه آن ضروری می باشد.

```
. labellist q1
direction:
    0 Right direction
    1 Wrong Direction
```

2. انجام فرمان پروبیست رگرسیون

`probit [dependent variable] [independent variables]`

`probit q1 d1 m6b d10 d2`

3. تفسیر نتیجه فرمان رگرسیون بعد از اجرای آن در برنامه ستاتا

در تفسیر نتیجه رگرسیون پروبیت یک تفاوت وجود دارد، در بحث مربوط به رگرسیون به طریقه حداقل مربعات شما آموختید که در هنگام تفسیر نتیجه رگرسیون، یک واحد تغییر در متغیر مستقل با ثابت نگه داشته تاثیرات دیگر متغیرهای مستقل به چه تعداد واحد تغییرات (کاهش یا افزایش) در متغیر وابسته روفا می گردید. اما در تفسیر نتیجه رگرسیون پروبیت، این گونه بیان میشود که با یک واحد تغییر در متغیر مستقل با ثابت نگه داشتن تاثیرات دیگر متغیرهای مستقل، متغیر وابسته چند فیصد احتمال وقوع آن کاهش یا افزایش پیدا می کند. برای فهم بیشتر به تفسیر ذیل توجه کنید:

1. در اینجا به تفسیر رابطه بین جنسیت و نظر مردم در مورد اینکه افغانستان به کدام سمت در حرکت است می پردازیم، چونکه علامه ضریب رگرسیون ($B = -0.06463$) منفی می باشد، رابطه بین متغیر مستقل جنسیت ($d1$) و متغیر وابسته ($q1$) منفی بوده، یعنی زن ها نسبت به مردان 6.4 درصد بیشتر احتمال دارد تا بگویند که کشور به سمت درست در حرکت است.

. probit q1 d1 m6b d10 d2

```
Iteration 0: log likelihood = -6025.1992
Iteration 1: log likelihood = -5995.5281
Iteration 2: log likelihood = -5995.5221
Iteration 3: log likelihood = -5995.5221
```

Probit regression	Number of obs	=	8,978
	LR chi2(4)	=	59.35
	Prob > chi2	=	0.0000
Log likelihood = -5995.5221	Pseudo R2	=	0.0049

direction	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
d1	-.0646357	.0286046	-2.26	0.024	-.1206998	-.0085717
m6b	.2383293	.0359654	6.63	0.000	.1678383	.3088202
d10	-.0053066	.0042001	-1.26	0.206	-.0135387	.0029254
d2	.0030824	.0010977	2.81	0.005	.0009309	.005234
_cons	-.0068433	.077997	-0.09	0.930	-.1597145	.146028

آزمون میانگین دو جمعیت T-test

محققین در رشته های مختلف از t-test برای مقایسه اوسط دو گروه استفاده میکنند.

مثلاً شما میخواهید تفاوت اوسط سن را در دو ولایت افغانستان دریابید و شما از هر ولایت درمورد سن 3 نفر پرسان میکنید و هر یک به ترتیب ذیل جواب میدهند.

سن سه نفر از باشندگان ولایت کابل و مزار شریف

(n) تعداد مشاهدات	کابل	مزار شریف
1	20	29
2	28	18
3	36	22
اوسط	28	23

شکل: ۵-۱۰ سن سه نفر از باشندگان دو ولایت کابل و مزار شریف

به نظر ساده دیده میشود که اوسط سن در شهر کابل (28 سال) نسبت به شهر مزار شریف (23 سال) بزرگتر است اما این تفاوت در اوسط سن دو ولایت میتواند که از نظر احصایوی منطقی نباشد، دلیل آن این است که تنها 3 نفر نمیتواند نماینده گی از تمام شهروندان این ولایت کند یعنی این تفاوت سنی میتواند به اساس چانس باشد.

اما اگر شما از هر شهر (کابل و مزار شریف) در مورد سن 300 باشندگان پرسید و هنوز هم یک تفاوت آشکارا در اوسط سنی این دو گروه مشاهده شود پس کمتر احتمال دارد که این تفاوت بر اساس چانس اتفاق بیافتد و به همین ترتیب به هر مقدار که تعداد نمونه ما زیاد باشد چانس اینکه نتیجه بدست آمده ما به اساس چانس نباشد زیادتراست.

به صورت عموم دو نوع T-test وجود دارد.

1. T-test یکطرفه

2. T-test دو طرفه

آزمون میانگین یک طرفه One sided T-test

در صورت که ما دیتای یک گروه را داشته باشیم و از گروه دومی را تنها اوسط آنرا داشته باشیم، بناً برای مقایسه اوسط دو گروه از T-test یکطرفه استفاده میکنیم.

ساختار انجام فرمان T-test در برنامه ستاتا به شکل ذیل میباشد.

```
ttest var1== [mean_var2]
```

در اینجا برای اجرای مثال عملی از دیتای سروی مردم افغانستان SAP استفاده میکنیم که شما میتوانید آنرا از لینک ذیل دریافت کنید.

<http://asiafoundation.org/where-we-work/afghanistan/survey>

یکی از سوالات که در این سروی از شندگان پرسان گردیده است عبارت سن آنها میباشد حالا اگر بخواهیم اوسط سن این گروه را با اوسط سن یک گروه دیگر که به صورت فرضی آنرا 35 سال در نظر گرفتیم مقایسه

کنیم در صفحه فرمان Command Window فرمان ttest نوشته به تعقیب آن متغیر را که نشان دهنده سن است (d2) را نوشته کرده و بعد اوسط سن جمعیت دوم را که به صورت تخمینی 35 فرض نمودیم نوشته میکنیم. نتیجه فرمان اجرا شده به شکل ذیل میباشد.

```
. ttest d2==35
```

One-sample t test

Variable	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
d2	9,586	34.6039	.1299777	12.72587	34.34912	34.85869

```
mean = mean(d2)                                t = -3.0474
Ho: mean = 35                                degrees of freedom = 9585
```

```
Ha: mean < 35                                Ha: mean != 35                                Ha: mean > 35
Pr(T < t) = 0.0012                        Pr(|T| > |t|) = 0.0023                        Pr(T > t) = 0.9988
```

1. فرضیه H_0 بیانگر این موضوع است که اوسط سن گروه اولی مساوی به 35 (اوسط سن گروه دومی که به صورت تخمینی فرض نمودیم) میباشد.

2. فرضیه H_A بیانگر این موضوع است که اوسط سن گروه اولی مساوی به 35 (اوسط سن گروه دومی که به صورت تخمینی فرض نمودیم) نمیشود.

اما از روی P-value قضاوت نموده میتوانیم که کدام فرضیه را تایید کنیم و کدام فرضیه را رد کنیم.

در صورت که P-value بزرگتر از 0.05 باشد پس گفته میتوانیم که فرضیه H_0 درست بوده و فرضیه H_A را رد میکنیم اما در صورت که P-value کوچکتر از 0.05 باشد پس فرضیه H_A را تایید کرده و فرضیه H_0 را رد میکنیم.

آزمون میانگین دو طرفه Two sided T-test

در صورت که ما معلومات در مورد هر دو گروه داشته باشیم پس از T-test دو طرفه استفاده میکنیم.

به طور مثال در سال 2015 یکی از سوالات در سروی مردم افغانستان این بود که در مورد سن هر یک از پاسخ دهنده گان (مردان و زنان) سوال شده، حالا اگر بخواهیم بدانیم که اوسط سن مردان و زنان باهم برابر بوده پس از T-test دو طرفه استفاده میکنیم.

ساختار فرمان T-test دو طرفه در برنامه ستاتا به شکل ذیل میباشد.

```
ttest continuous variable = by (category variable)
```

```
ttest d2, by (d1)
```

و نتیجه فرمان به شکل ذیل در صفحه ستاتا نمایش داده میشود.


```
. ttest d2 , by(d1)
```

Two-sample t test with equal variances

Group	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
Male	4,847	36.07427	.1946841	13.55399	35.6926	36.45594
Female	4,739	33.10002	.1689367	11.62967	32.76883	33.43122
combined	9,586	34.6039	.1299777	12.72587	34.34912	34.85869
diff		2.974252	.2582042		2.468117	3.480386

```
diff = mean(Male) - mean(Female)                                t = 11.5190
Ho: diff = 0                                                    degrees of freedom = 9584

Ha: diff < 0                                                    Ha: diff != 0                                                    Ha: diff > 0
Pr(T < t) = 1.0000        Pr(|T| > |t|) = 0.0000        Pr(T > t) = 0.0000
```

در اینجا به چند نکته ذیل باید توجه شود.

1. فرضیه H_0 بیانگر این موضوع است که اوسط سن بین مردان و زنان یکسان می باشد.

2. برعکس فرضیه H_A بیانگر این موضوع است که اوسط سن مردان و زنان یکسان نمی باشد.

در مورد این که بدانیم کدام فرضیه ما درست می باشد و کدام فرضیه ما اشتباه است از روی P-value میتوانیم قضاوت کنیم.

در صورت که p-value بزرگتر از 0.05 باشد پس فرضیه H_0 را تایید میکنیم، یعنی اوسط سن مردان و خانم ها یکبرابر است، و فرضیه H_A را رد میکنیم.

اما در صورت که P-value کوچکتر از 0.05 باشد پس فرضیه H_A را تایید میکنیم، یعنی اوسط سن مردان و زنان یکسان نیست، و فرضیه H_0 را رد میکنیم.

3. نوع دیگر ttest دو طرفه شباهت زیاد با هم دارد، در این دو نوع، شما دو گروه یا متغیر را مستقیماً با همدیگر مقایسه میکنید. مثلاً در تحقیقات طبی وقتی یک دانشمند یک دوا را کشف میکند، آنرا در قدم اول روی بعضی حیوانات آزمایشگاهی آزمایش میکند. درین صورت، بخاطر دریافت تاثیر این دوا، دو گروه میباشند که یکی دوا را گرفته و دیگری نگرفته (RCT). یک ttest به این دانشمند کمک میکند تا بداند که آیا این دوا تاثیر داشته یا نداشته است، اگر داشته است این تاثیر قابل اعتبار است یا خیر.

اما به هر صورت نتیجه فرمان ttest به شکل بالا می باشد که باز هم نظر به قیمت p-value تصمیم میگیریم که فرضیه H_0 و یا فرضیه H_A را قبول کنیم.

• به یاد داشته باشید که در هنگام راپور باید قیمت های t و df (degree of freedom) را همراه با p-value گزارش دهید.



تمرین

1. در مورد آزمون همبستگی معلومات دهید؟
2. همبستگی مثبت و منفی را تعریف و برای هر یک مثال دهید؟
3. از فرمول همبستگی استفاده کنید و رابطه بین دو متغیر کمی که روزانه در محیط کار یا درسی با آنها روبرو هستید، را دریابید و چگونگی رابطه بین دو متغیر را به صورت مفصل تشریح کنید؟
4. درمورد رگرسیون خطی تحقیق نمائید و معلومات خود را به صورت مفصل بنویسید؟
5. حالا بعد از مطالعه رگرسیون خطی بگویید که کدام زمان و چگونه از آن استفاده میکنید؟
6. از دیتاست سروی مردم افغانستان در سال 2015 استفاده کنید و یک مدل تحلیلی بسازید و معلومات ذیل را تکمیل کنید؟
 الف: فرضیه صفر و فرضیه جاگزین بسازید.
 ب: راجع به اندازه نمونه معلومات دهید.
 ج: رابطه بین متغیر وابسته و مستقل هر یک را به صورت جداگانه تفسیر دهید با ذکر (Coefficients and P-values)
 د: در مورد (Coefficient determinant) معلومات دهید و در این مدل آن را به تفسیر بگیرید.
7. چه زمانی از رگرسیون پروبیت استفاده میکنید؟
8. مدل تحلیلی برای رگرسیون پروبیت انکشاف دهید و معلومات ذیل را بعد از انجام فرمان رگرسیون تکمیل کنید؟
 الف: فرضیه صفر و متناوب بسازید.
 ب: راجع به اندازه نمونه معلومات دهید.
 ج: رابطه بین متغیر وابسته و مستقل هر یک را به صورت جداگانه تفسیر دهید با ذکر (Coefficients and P-values)
 د: در مورد (Coefficient determinant) معلومات دهید و در این مدل آن را به تفسیر بگیرید.

References

- Acock, Alan C. (2014). A Gentle Guide to Stata. Texas: Stata Press Publication
- Dougherty C. (2016). Elements of Econometrics. London: University of London
- Gujarati, Damodar N. (2004). Basic Econometrics, Fourth Edition. New York: The McGraw-Hill Companies
- Hansen, Bruce E. (2015). Econometrics. Wisconsin: John Wiley & Sons
- Kerns, Jay G. (2011). Introduction to Probability and Statistics. Vienna: The McGraw-Hill Companies
- Reagle, Derrick & Salvatore, Dominick. (2001). Statistics And Econometrics, Second Edition. New York: The McGraw-Hill Companies
- Rumsey, Deborah. (2010). Statistics for Dummies. Indiana: Wiley Publishing. Inc.
- Spiegel, Murray R. (1999). Theory and Problems of Statistics. New York: The McGraw-Hill Companies
- Suseela, Tmt. N. & Sundaram Gnana G. (2005). Statistics for Higher Secondary. Tamilnadu: TEXTBOOK CORPORATION
- Wooldridge, Jeffery M. (2012). Introductory Econometrics A Modern Approach, Fifth Edition. Michigan: South-West, Cengage Learning

انگیزه اساسی برای تهیه این کتاب در نخست ادامه مواد درسی برای کورس تحلیل دیتا توسط برنامه ستاتا و ضمیمه برای درک بهتر مفاهیمی که در رهنمای برنامه ستاتا تحت عنوان "تحلیل دیتا توسط برنامه ستاتا" است، می باشد.

بنیاد آسیا از اداره انکشاف بین المللی ایالات متحده امریکا (USAID)، وزارت امور خارجه و تجارت آسترالیا (DFAT)، و انجمن همکاری های بین المللی های آلمان (GIZ) بخاطر حمایت شان از برنامه های ارتقای ظرفیت تحقیقاتی ابراز سپاس و امتنان مینماید.