

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

راهنمای طلایی  
تسلی طلایی  
پیشگام طلایی

انتشارات طلایی  
پویندگان دانشگاه



[www.bookgolden.com](http://www.bookgolden.com)

عنوان درس :

# آمار مقدماتي

راهنمای طلایی  
تست طلایی  
پیشک طلایی

انتشارات طلایی  
پویندگان دانشگاه



[www.bookgolden.com](http://www.bookgolden.com)

# فصل اول

مفاهيم اساسي علم آمار

© راجه‌ای طلایی  
کتاب طلایی  
پیشگام طلایی

انتشارات طلایی  
پویندگان دانشگاه



[www.bookgolden.com](http://www.bookgolden.com)

# مقدمه :

واژه **statistics** که به فارسی آن را آمار ترجمه کرده اند از کلمه لاتین **status** یا کلمه ایتالیایی **statista** و یا کلمه آلمانی **statistik** که همگی به معنی دولت هستند استخراج شده است.

آمار را به معنی روشهایی برای جمع آوری، تنظیم و تجزیه و تحلیل اطلاعات عددی درباره موضوعی تعریف نموده اند.

# مفهوم آمار توصیفی و آمار استنباطی

توصیفی: روش‌های تنظیم ، خلاصه سازی ، رسم

نمودار ،

جدول و پردازش اطلاعات

} آمار

استنباطی : روش‌های جمع آوری ، تجزیه و تحلیل

اطلاعات ،

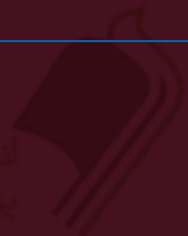
تفسیر و تعمیم اطلاعات

موضوع علم آمار را می‌توان از چند دیدگاه مختلف تقسیم  
بندی کرد . آمار در نخستین مرحله به دو بخش توصیفی  
و آمار استنباطی تفکیک می‌شود .

وقتی با مجموعه‌ای از اطلاعات روبرو هستیم ابتدا باید تصمیم بگیریم که  
چگونه این اطلاعات را خلاصه کنیم تا بهتر بتوانیم از آنها استفاده کنیم.

در واقع آمار توصیفی بخشی از علم آمار است که به جمع آوری ، خلاصه کردن،نمایش و پردازش اطلاعات می پردازد.

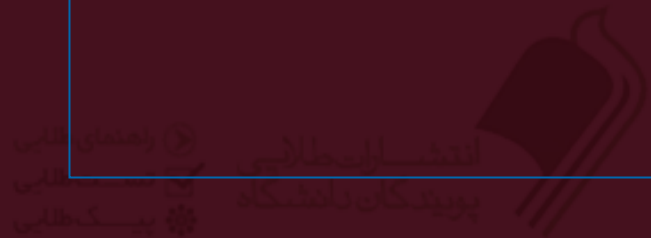
در این روش پس از جمع آوری اطلاعات،تهیه جدولهای خلاصه و رسم نمودارهای آماری به یافته های مناسب برای مطالعه تحقیق می رسد



گرچه آمار توصیفی بخش مهم از آمار است و هنوز هم  
استفاده گسترده ای از آن مخصوصاً برای توصیف ویژگیهای  
مهم داده‌ها

برای به عمل می آید

اما در سالهای اخیر توجه بیشتر به آمار استنباطی ، یعنی روشهایی است  
که به کمک آنها می توان اطلاعات موجود در مجموعه‌ای محدود از  
داده‌ها را به مجموعه‌ای بزرگتر که داده‌ها از آن بدست آمده‌اند تعمیم



# مفاهيم بنيادي آمار

انجام پژوهشهاي آماري مستلزم مشخص شدن  
برخي مفاهيم اساسي مي باشد كه در اين قسمت  
مطالعه و تشریح اين مفاهيم اوليه و اساسي  
در آمار مي پردازيم كه عبارتند از:  
جمعيت ، نمونه ، متغير و داده



## جمعیت :

مجموعه ای از افراد یا اشیاء که می‌خواهیم یک‌یاچند نفر خصوصیت را در مورد آنها مطالعه کنیم را جمعیت یا جامعه آماری می‌نامند.

مثلاً َّ جمعیت دانشجویان دانشگاه پیام نور ، از لحاظ سن، جمعیت کارگران يك کارخانه از نظر میزان دستمزد یا جمعیت نوزدانی که سال گذشته بدنیا آمده‌اند از نظر مصرف شیر خشک.

براي انجام هر کار آماری باید آن جمعیت و خصوصیت مورد مطالعه بدون هر گونه انجام قبلاً مشخص می شوند .  
جمعیت دو نوع است جمعیت متناهی که تعداد اعضاء تشکیل دهنده معین و محدود است

به طوری که اگر آنها را تک تک شناسایی کنیم و کنار بگذاریم تمام می شوند مانند جمعیت اتومبیلها شهر تهران و جمعیت نامتناهی عناصر تشکیل دهنده آن نامعین و پایان ناپذیر می باشند.

مانند جمعیت ستارگان آسمان یا جمعیت دانه های یک مزرعه گندم که هر یک از عناصر جمعیت را یک فرد می نامند و واژه فرد تنها به انسان اطلاق نمی شود.

تعداد عناصر جمعیت را اندازه جمعیت گویند و آن را با  $N$  نمایش می دهند.

مطالعه تک تک افراد جمعیت به علت هزینه زیاد، کمی وقت و نداشتن امکانات کافی، اغلب مقدور نمی باشد .  
بنابراین قسمتی از جمعیت را به جای تمام آن برای مطالعه در نظر می گیرند.

**نمونه :**

نمونه بخشی از جامعه آماری است که طبق اصول و ضوابطی معین انتخاب می شود و مطالعه نمونه به جای مطالعه، کل جمعیت توصیه می شود .



تعداد عناصر نمونه را اندازه نمونه گویند و با  $n$  نمایش می‌دهند.

انتخاب نمونه طبق روشها، اصول و ضوابطی معین که به روشهای نمونه‌گیری معروفند صورت می‌پذیرد.

## متغیر :

خصوصیت یا صفتی مانند گروه خونی، هوش، قد، وزن، محل سکونت، طول عمر و ... که از هر فرد به فرد دیگر و از هر شخص به شخص دیگر در جمعیت آماری تغییر می‌کند را متغیر می‌نامند و اغلب با  $X$  نشان داده می‌شود.



## متغیرها

الف) متغیرهای کمی: صفاتی هستند که قابل اندازه گیری یا شمارش هستند مانند سن طول عمر، وزن، درآمد، هزینه، و...

ب) متغیرهای کیفی: صفاتی هستند که واحد نداشتند و قابل اندازه گیری نیستند مانند جنس، مرغوبیت، مهارت، گروه خونی و...



همچنین متغیرها را بر اساس نقشی که بر عهده دارند می توانند به دودسته زیر تقسیم شوند.

الف) متغیرهای مستقل : این متغیرها تحت کنترل محقق بوده و تغییرات آنها، دستخوش تغییرات سایر متغیرها نباشد .  
متغیرها

ب) متغیرهای وابسته : تغییرات این گونه متغیرها بسته به تغییرات

متغیر یا متغیرهای مستقل می باشد.

## داده :

در مطالعه خصوصیتی در جمعیت، متغیر مورد نظر را با مقیاسی مناسب اندازه می گیرند در این صورت مجموعه ای از اعداد حاصل می شود

که به آنها داده می گویند. در واقع اطلاعات عددی درباره موضوعی (خصوصیت) را داده می گویند . داده ها دو نوع اند

الف) داده های گسسته: داده هایی که ماهیت اعشاری

ندارند

مانند تعداد دانشجویان يك

كلاس، تعداد مدارس يك شهر و ....

داده ها

ب) داده های پیوسته: داده هایی که ماهیت اعشاری

دارند

پارامتر و آماره

مطالعه جامعه برای بدست آوردن برخی از شاخصهاست. این شاخصها، چنانچه با

اندازه گیری تمامی عناصر جامعه آماری (سرشماری) بدست آمده شده باشند آنها را پارامتر و اگر با استفاده از نمونه که بخشی از جامعه است بدست آمده باشند آنها را آماره می گویند .

اگر متوسط درآمد کارکنان دولت با استفاده از

اندازه گیری درآمد کلیه کارکنان دولت محاسبه شود آن را

پارامتر و اگر با استفاده از اندازه گیری درآمد نمونه ای

از کارکنان بدست آید آن را آماره می گویند .



# مراحل پژوهشي علم آمار

الف) تعیین هدف

ب) جمع آوري داده‌ها

ج) تجزيه و تحليل داده‌ها

د) نتیجه گيري



اندازه گيري و مقياس سازي

اطلاق يك عدد حقيقي به متغير  $x$  ، طبق قاعده‌اي مشخص را  
اندازه گيري يا مقياس سازي  $x$  و اين قاعده را مقياس گویند.  
استیونس در مقاله‌اي (به سال 1946 میلادي) چهار نوع  
مختلف از مقیاس‌هاي اندازه‌گيري را مشخص کرده است .  
این مقیاسها عبارتند از :

الف) مقياس اسمي :

هر گاه عدد مورد نظر تنها براي تمیيز افراد يا اشیاء  
بكار رود و نتوان آنرا براي مقایسه و چهار عمل اصلي  
حساب بكار برد از مقياس اسمي استفاده مي‌شود.

(ب) مقیاس ترتیبی :

مقیاس اسمی هر گاه صعودی باشد را مقیاس ترتیبی می گویند.

(ج) مقیاس فاصله ای :

وقتی که مقیاس همه خصوصیات يك مقیاس ترتیبی را دارا باشد و علاوه بر آن فاصله بین هر دو عدد بر روی مقیاس میزان مشخصی داشته باشد.

در آن صورت به مقیاسی دست یافته ایم به مراتب قوی تر که آن را مقیاس فاصله ای می نامند در این مقیاس نسبت بی معنی است و نسبت فواصل ثابت می ماند در اندازه گیری با این مقیاس صفر به معنای هیچ نمی باشد و صرفاً جنبه قراردادی و اختیاری دارد

د) مقیاس نسبی :

به مقیاسی که علاوه بر همه خصوصیات مقیاس  
فاصله‌ای دارای نقطه صفر واقعی نیز هست مقیاس  
نسبی می‌گویند در این مقیاس ، نسبت هر دو نقطه روی  
مقیاس از واحد  
اندازه گیری مستقل است.



نماد جمعبندي ( )

در مطالعه آمار همه جا با عمل جمع کردن داده ها سر  
و کار داریم براي اجتناب از نوشتن مکرر علامت (+)  
شکل بزرگ حرف يوناني سيگما ( ) را با عنوان  
نمادي براي عمل جمع کردن به کار مي بریم



www.bookgolden.com

مثال :

$$x_i = x_1 + x_2 + \dots + x_n$$

4

$$x_i = x_2 + x_3 + x_4$$

5

$$(x_i - 5) = (x_3 - 5) + (x_4 - 5) + (x_5 - 5)$$

$i=3$

و



# برخي از خواص اصلي

$$\sum_{i=1}^n b x_i = b \sum_{i=1}^n x_i$$

(الف)

$$\sum_{i=1}^n (b x_i + a) = b \sum_{i=1}^n x_i + n a$$

(ب)

$$\sum_{i=1}^n (x_i - a)^2 = \sum_{i=1}^n x_i^2 - 2a \sum_{i=1}^n x_i + n a^2$$

(ج)

# فصل دوم

## جدولهای فراوانی

راهنمای طلایی  
تسلی طلایی  
پیشگامی

انتشارات طلایی  
پویندگان دانشگاه



[www.bookgolden.com](http://www.bookgolden.com)



## جدول فراواني :

نمایش داده‌ها را با نظم خاص در چند سطر و ستون يك  
جدول آماری می‌گویند يك جدول آماری باید به گونه‌ای  
تنظیم شود که به كمك آن بتوان تا حدودي پاره‌ای  
اطلاعات نهفته در  
داده‌ها را در آن مشاهده کرد .



مثال: فرض کنید پزشکی 20 بیمار قلبی دارد که گروه خونی آنها به صورت زیر است.

AB, AB, O, AB, B, B, A, A, O, O, O, A, A, A, AB, O, A, B, O  
O

می توان اطلاعات فوق را در جدول زیر خلاصه نمود.

گروه خونی  $f_i$  (فراوانی)

6

A

3

B

4

AB

هر گاه  $n$  چیز از  $k$  نوع به ترتیب با تعداد  $f_1, f_2, \dots, f_k$  تشکیل شده باشند

این تعداد را فراوانی ها و  $\frac{f_1}{n}, \frac{f_2}{n}, \dots, \frac{f_k}{n}$  را فراوانی نسبی گویند و با  $r_i$  نمایش می دهند و

فراوانی جمعی

$$F_i = \sum_{j=1}^i f_j$$

فراوانی نسبی جمعی نمایش

$$R_i = \sum_{j=1}^i r_j = \frac{F_i}{n} \text{ می دهند.}$$

بنابراین می توان جدول فراوانی گروه خونی 20 بیمار قلبی را به صورت زیر تکمیل نمود.

**X: گروه خونی**

	$f_i$	$r_i$	$F_i$	$R_i$
$A$	6	$\frac{6}{20}$	6	$\frac{6}{20}$
$B$	3	$\frac{3}{20}$	9	$\frac{9}{20}$
$AB$	4	$\frac{4}{20}$	13	$\frac{13}{20}$
$O$	7	$\frac{7}{20}$	20	1

**n = 20**

**1**

# الگوریتم تنظیم جدولهای فراوانی

برای طبقه‌بندی داده‌ها الگوریتم غیرحقیقی در چهار مرحله ارائه می‌گردد

مرحله اول : داده‌ها را جمع‌آوری شده را از کوچک به بزرگ مرتب و دامنه تغییرات را محاسبه می‌کنیم .

$$\text{دامنه تغییرات} = R = x_{(n)} - x_{(1)} = \text{Max}x_i - \text{Min}x_i$$

مرحله دوم: تعداد طبقات (  $k$  ) را مشخص می‌کنیم دو روش برای تعیین  $k$  پیشنهاد می‌شود .

الف) فرمول زیر توسط استورجس برای تعیین  $k$  بدست آورده شده است .

که در آن  $n$  تعداد کل داده ها است .  $k = \lceil \log_2 n \rceil$

ب)  $k$  کوچکترین عدد صحیحی است که به ازای آن برای اولین بار نامساوی غیر اکید  $2^k \geq n$  برقرار است .

## مرحله سوم :

محاسبه می کنیم

$$\omega = \frac{R}{K}$$

هر گاه در تعیین طول طبقات حاصل تقسیم دامنه تغییرات بر تعداد طبقات عدد اعشاری بدست آمد عدد مذکور را به نزدیکترین واحد به بالا رند می کنیم .

مرحله چهارم : سازماندهی داده ها را آغاز می کنیم حد پایین طبقه اول را با کوچکترین مشاهده شروع کرده و حدود طبقات را با فاصله  $\omega$  تشکیل می دهیم

مثال : مؤسسه‌اي يك نوع باطري براي اتومبيلهاي خود استفاده

مي‌کند از مسئول اين مؤسسه طول عمر 40 باطري مستعمل را بر حسب سال سؤال کرده و داده‌هاي زير را بدست مي‌آوريم .

4/1 3/5 4/5 3/2 3/7 3 2/6 3/4 1/6 3/1 3/2

2/2

3/1 4/7 3/6 2/5 4/3 3/6 3/4 2/9 3/3 3/9 3/1

3/8

3/1 3/7 4/4 3/2 4/1 1/9 3/4 4/7 3/8 3/2 2/6

3/9 3 4/2 3/5 3/3





در این صورت  $n = 40$  و تعداد طبقات از رابطه  $k = \lceil \frac{n}{3.22 \log 40} \rceil$

محاسبه می شود بنابراین  $k = \lceil \frac{40}{3.22} \rceil = 13$  بنابراین طول هر طبقه عبارت است از :

$$\omega = \frac{R}{K} = \frac{\max(x_i) - \min(x_i)}{k} = \frac{4.7 - 1/6}{13} = 0.33 \quad 0.6$$

حدود واقعی طبقات	حدود طبقات	$f_i$ (فراوانی)	$x_i$ (مرکز (نماینده) طبقات)
1/55- 2/25	1/6- 2/2	2	1/9
2/25- 2/85	2/2- 2/8	4	2/5
2/85- 3/45	2/8- 3/4	13	3/1
3/45- 3/95	3/4- 4/0	13	3/7
3/95- 4/55	4/0- 4/6	6	4/3
4/55- 5/15	4/6- 5/2	2	4/9

# فصل سوم

## نمودارهاي آماری

© راجه‌های طلایی  
کتاب طلایی  
پیشگامی

انتشارات طلایی  
پویندگان دانشگاه



[www.bookgolden.com](http://www.bookgolden.com)

## نمودارهاي آماری

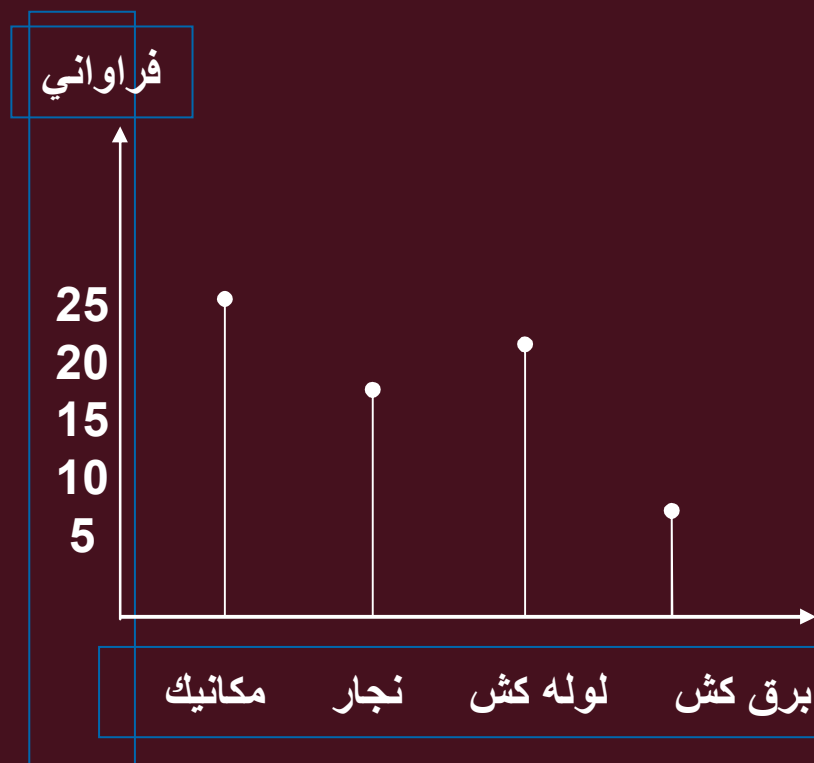
نمایش داده ها را، طبق قراردادهای خاص به صورت هندسی ، يك نمودار آماری می گویند.

### نمودار میله ای

در این نمودار صفت کیفی را به صورت نقاطی روی محور ( X ها ) مشخص می کنیم و سپس از نقاط حاصل، خطهای عمودی بر محور رسم می کنیم که ارتفاع آن برابر فراوانی یا فراوانی نسبی متغیر مربوطه می باشد.

مثال: توزيع فراواني کارگران يك کارخانه بر حسب تخصص داده شده است نمودارميله اي آن عبارت است از:

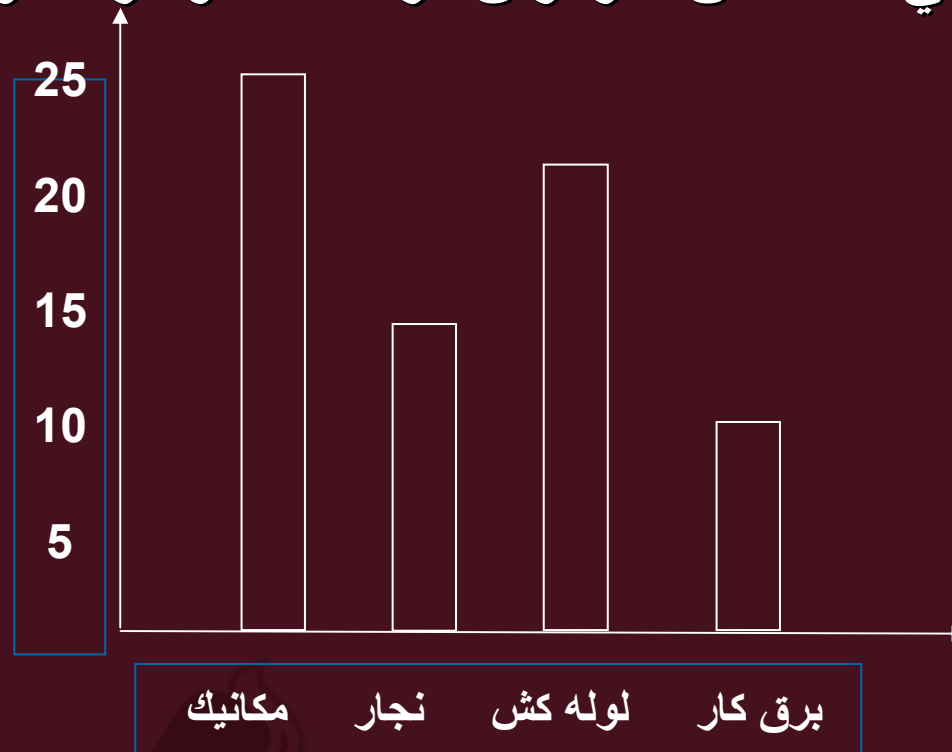
تخصص	فراواني
مکانیک	25
نجار	16
لوله کش	20
برق کش	9



## نمودار ستونی :

در نمودار میله ای، هرگاه میله ها را به صورت مستطیل نشان دهیم نمودار ستونی حاصل می شود.

مثال: در توزیع فراوانی تخصص کارگران کارخانه نمودار ستونی عبارت است از

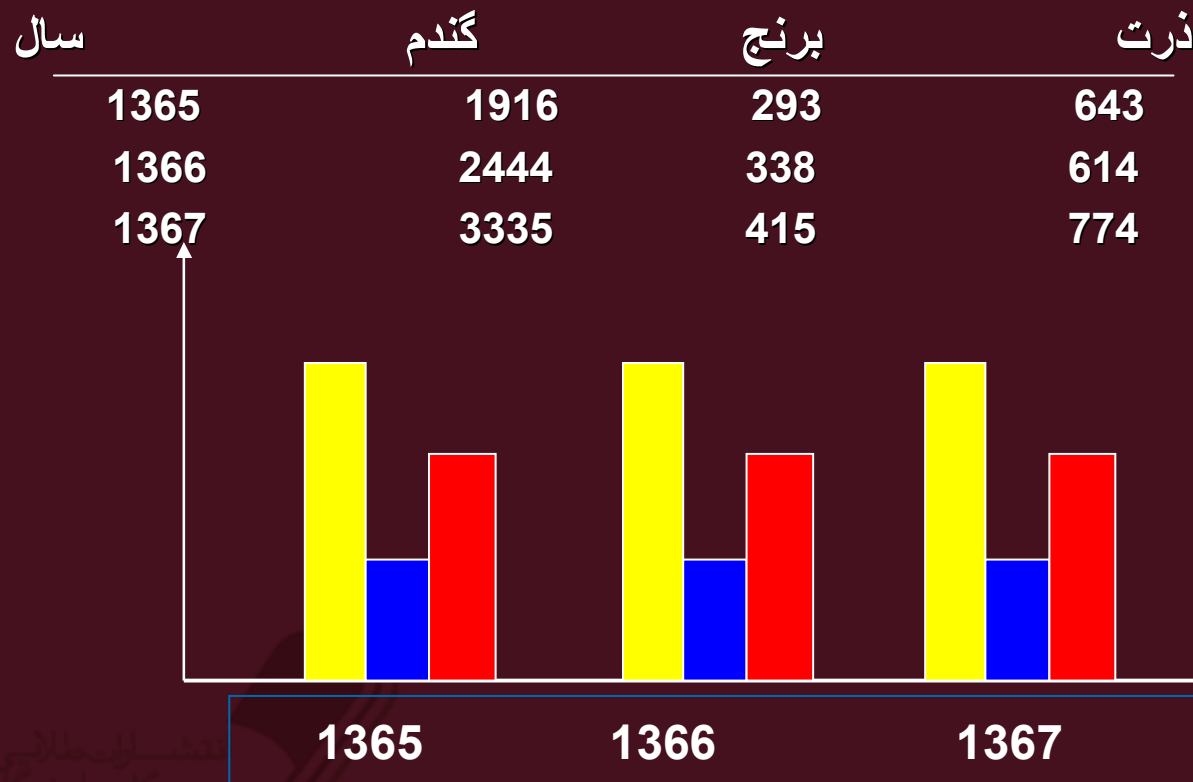


## نمودار میله ای چند تایی :

اگر داده ها را بر مبنای صفات کیفی طبقه بندی کنیم و برای مقایسه آنها از دو یا چند صفت استفاده کنیم در آن صورت از نمودار میله ای

### چند تایی استفاده می کنیم .

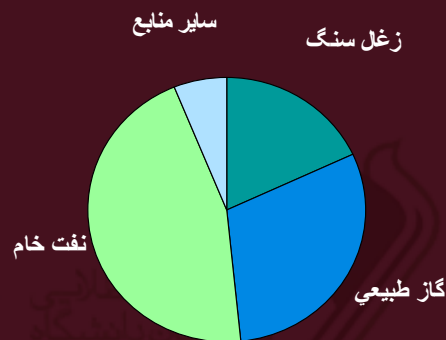
مثال : مقدار محصول گندم ، برنج و ذرت بر حسب هزار تن در سالهای 1365 و 1366 و 1367 در یک استان به شرح زیر است.



این نمودار برای نمایش داده های کیفی مورد استفاده قرار می گیرد دایره به قطاعهایی متناسب با فراوانی نسبی هر طبقه تقسیم می شود اندازه هر قطاع از رابطه ی  $S_i = r_i \times 360^\circ$  بر حسب درجه بدست می آید.

مثال: وضعیت مصرف انرژی بر حسب صد هزار B.T.U از هر نوع در سال 1354 داده شده است نمودار دایره ای مربوط به تغییرات نوع انرژی عبارت است از :

سال	گاز طبیعی	زغال سنگ	نفت خام	سایر منابع	مجموع
1354	105	62	158	21	346



گاز طبیعی

$$s = \frac{105}{346} \times 360 = 109$$

زغال سنگ

$$s = \frac{62}{346} \times 360 = 65$$

نفت خام

$$s = \frac{158}{346} \times 360 = 164$$

سایر منابع

$$s = \frac{21}{346} \times 360 = 22$$



## نمودار شاخه و برگ :

برای تعبیه نمودار شاخه و برگ ارقام مشاهدات را که دستکم دو رقمی هستند را به دو بخش شاخه و برگ تقسیم می‌کنیم شاخه شامل يك یا چند رقم اولیه و برگ شامل ارقام باقیمانده ، وقتی که مجموعه شاخه ها انتخاب شد این اعداد در سمت چپ صفحه نوشته می شود و در کنار هر شاخه ، تمام برگهای متناظر با مقادیر داده ها به ترتیب مشاهده ثبت می شود.

مثال: نمودار شاخه و برگ تعداد دانشجویان 15 درس دانشکده عبارت است از :

شاخه	برگ
0	9
1	5
2	0 2 2 4 4 4 8
3	0 0 4 6 7
4	2

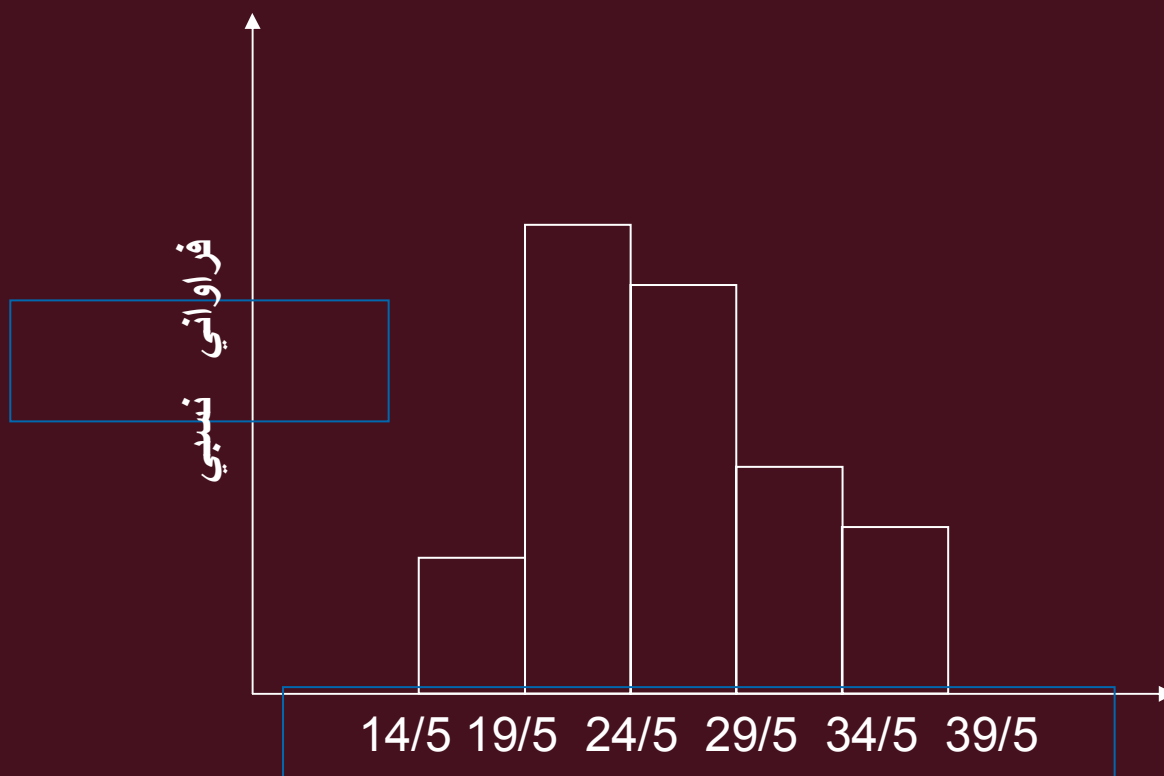


## نمودار هیستوگرام :

در توزیعهای فراوانی که دارای فاصله طبقات مساوی هستند هیستوگرام نموداری است مرکب از مستطیل هایی که عرض آن برابر فاصله طبقه و ارتفاع آن برابر فراوانی یا فراوانی نسبی هر طبقه است. در نمودار هیستوگرام ، محور افقی دستگاه مختصات با حدود طبقات و محور عمودی با فراوانی یا فراوانی نسبی مدرج می گردد.

مثال : نمودار هیستوگرام توزیع سن ازدواج عبارت است از :

سن ازدواج	$f_i$	$r_i$
14 / 5 – 19 / 5	18	0 / 09
19 / 5 – 24 / 5	74	0 / 37
24 / 5 – 29 / 5	62	0 / 31
29 / 5 – 34 / 5	26	0 / 13
34 / 5 – 39 / 5	20	0 / 1
مجموع	n= 200	



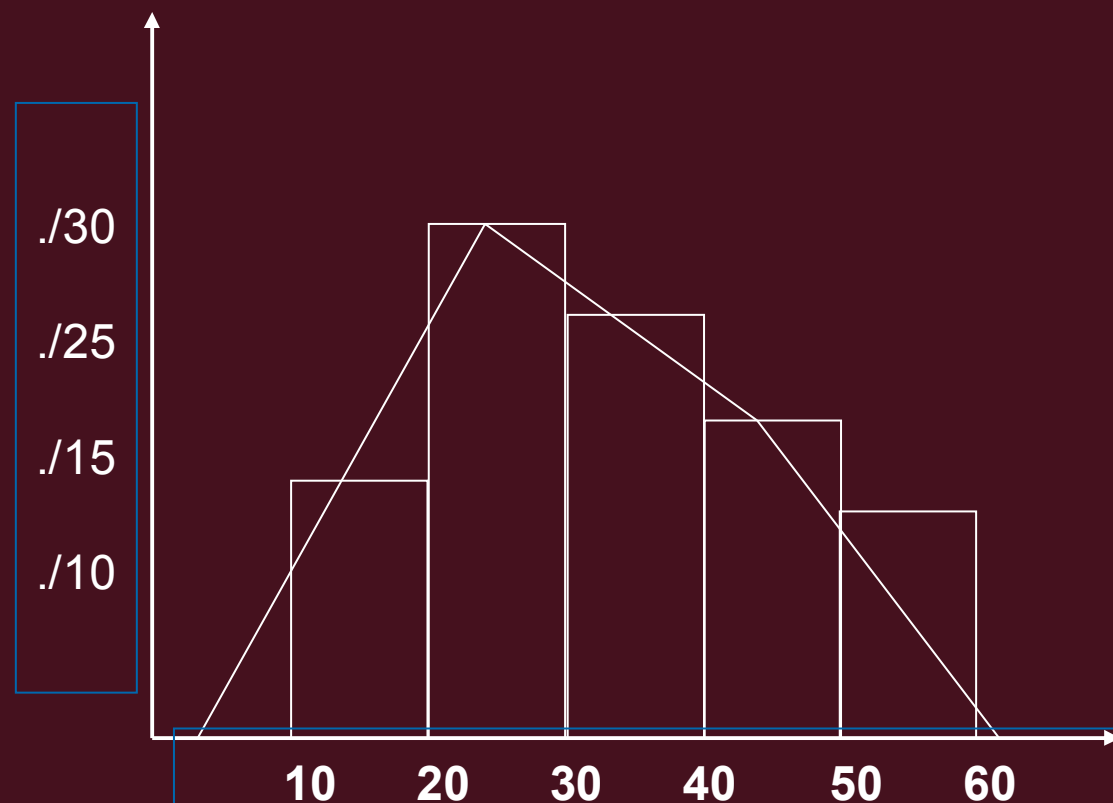
## نمودار چند ضلعي (چند بر فراواني)

اگر نقطه‌اي وسط قاعده هاي بالايي مستطيل‌هاي هيس‌توگرام  
ونقطه‌هاي وسط رده‌هايي را كه بلافاصله در دو انتهاي  
هيس‌توگرام بوده و داراي فراواني صفر هستند به هم بپيوندیم يك  
خط شكسته به دست مي آید كه آن را چند بر يا چند ضلعي  
فراواني مي گويند .



مثال: شرکت ایران دارو داراي 100 کارمند است که در آمد ماهیانه آنها در جدول زیر آمده است.

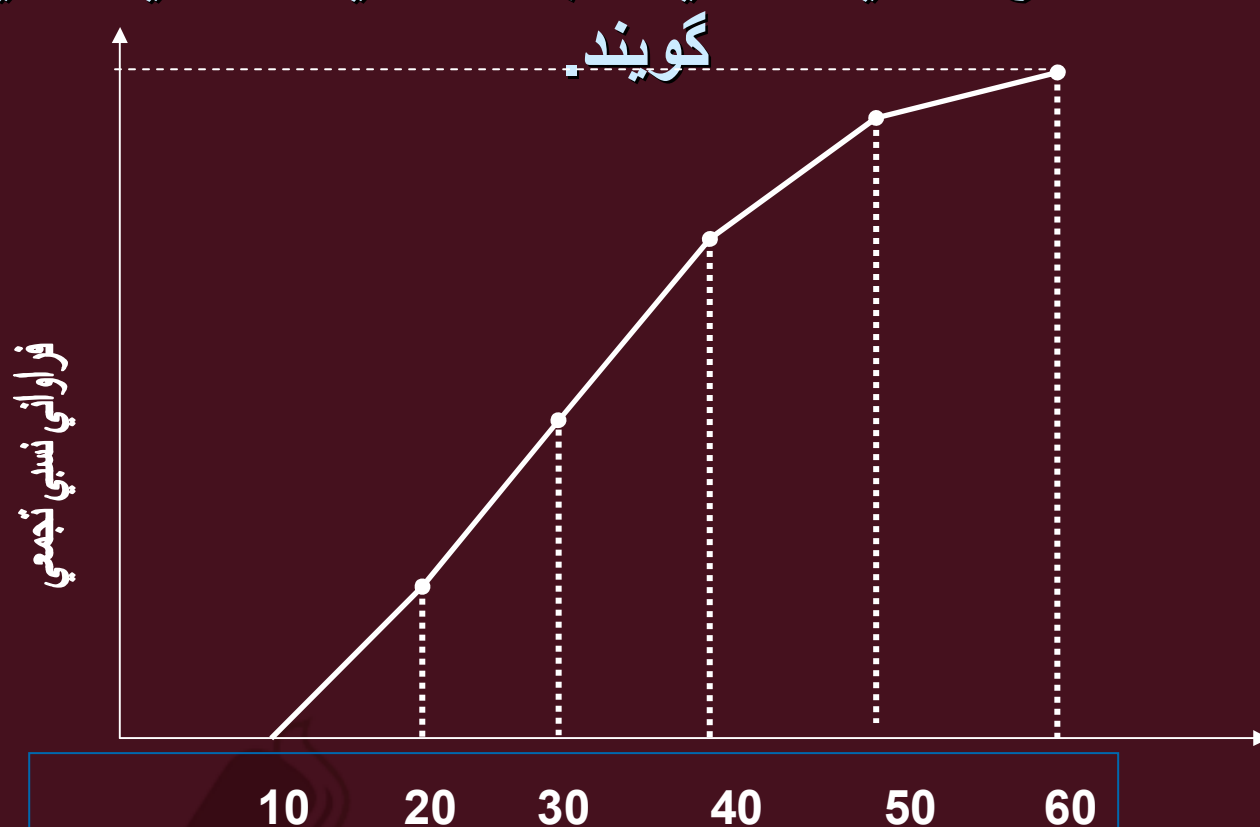
در آمد	$f_i$
10 – 20	15
20 – 30	30
30 – 40	25
40 – 50	20
50 – 60	10
مجموع	$n = 100$



نمودار هیستوگرام و چند ضلعي مربوط به درآمد ماهیانه

## نمودار توزیع نسبی تجمعی

اگر نقاطی را که طول آنها مرز طبقه‌ها و عرض آنها فراوانی نسبی تجمعی تا آن مرز باشد به هم بپیوندیم يك خط شکسته بدست می‌آید که آن را نمودار توزیع نسبی تجمعی یا چند ضلعی فراوانی نسبی تجمعی



نمودار توزیع نسبی تجمعی مربوط به در آمد ماهیانه

# منحنی‌های فراوانی و فراوانی تجمعی

هرگاه تعداد داده‌ها (  $n$  ) زیاد باشد تعداد طبقات طبق

$$k = \left[ 1 + \frac{1}{5} \log n \right]$$

زیاد می‌شود و در نتیجه طول

$$\left( w = \frac{R}{K} \right)$$

طبقات

کوچک می‌شود در این صورت چند ضلعی

فراوانی و چند ضلعی فراوانی تجمعی دارای اضلاع

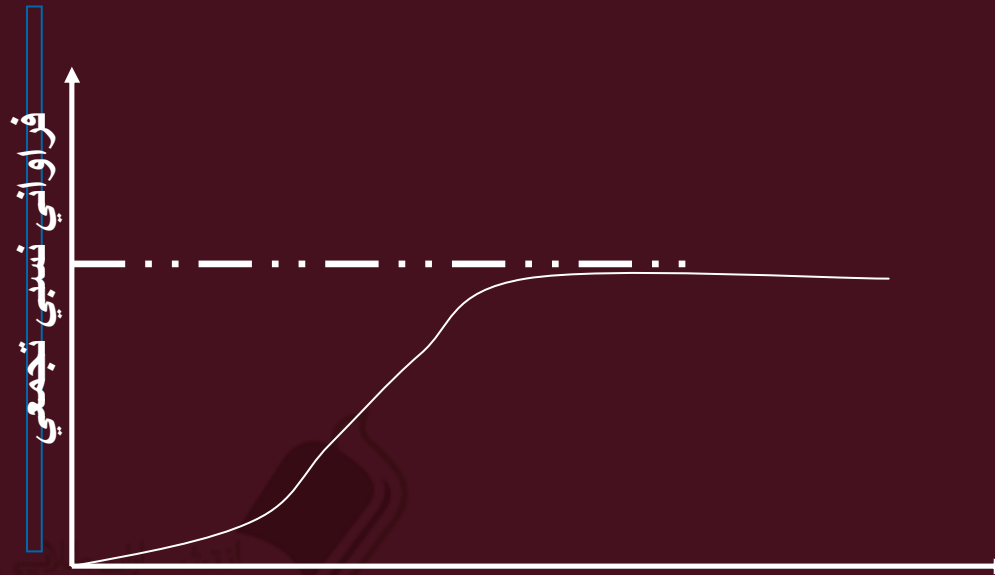
کوچک و زیادی خواهد شد و می‌توان بر آنها منحنی‌هایی

منطبق کرد که به ترتیب منحنی فراوانی تجمعی نامیده

می‌شوند .



# منحني فراواني



منحني فراواني نسبي تجمعي

# فصل چهارم

## معیارهای گرایش به مرکز

راهنمای طلایی  
تسبیح طلایی  
پیشگامی

انتشارات طلایی  
پویندگان دانشگاه



[www.bookgolden.com](http://www.bookgolden.com)

# معیارهای گرایش به مرکز

در آمار همواره سعی می‌شود تا اطلاعات نهفته در کل داده‌ها را به صورت يك یا چند عدد معقول در آورد این اعداد معیارهای گرایش به مرکز نامیده می‌شود. از مهمترین معیارهای گرایش به مرکز عبارتند از :

الف) میانگین

ب) نما

ج) میانه



## الف ( میانگین :

اصلی‌ترین و مورد استفاده‌ترین معیار گرایش به مرکز میانگین است .  
 $x_1, x_2, \dots, x_k$

فرض کنید  $n$  داده به صورت  $f_1, f_2, \dots, f_k$  با فراوانی‌ها

باشند در این صورت

$$\bar{x} = \frac{\text{تعداد کل داده‌ها}}{\text{مجموع کل داده‌ها}} = \frac{\sum_{i=1}^k f_i x_i}{n}$$

را میانگین داده‌ها می‌گویند

مثال : میانگین اعداد 1، 3، 5، 7، 9، 11

$$\bar{x} = \frac{1+3+5+7+9+11}{6} = 6$$

مثال : میانگین حقوق کارمندان اداره‌ای بر حسب واحد ده‌هزار تومان در جدول فراوانی زیر عبارت است از :



حقوق	$f_i$	$x_i$	$f_i x_i$	$\bar{x} = \frac{\sum_{i=1}^k f_i x_i}{n} = \frac{360}{40} = 59$
20-40	5	30	150	
40-60	17	50	850	
60-80	13	70	910	
80-100	5	90	450	
	n=40		$\sum f_i x_i = 2360$	

مثال : میانگین تعداد فرزندان 20 خانواده در  
جدول فراوانی زیر عبارت است از :

$x_i$	$f_i$	$f_i x_i$
0	5	0
1	3	3
2	4	8
3	6	18
4	2	8
$n = 20$		37

$$\bar{x} = \frac{\sum_{i=1}^k f_i \times x_i}{n} = \frac{37}{20} = 1.85$$



## میانگین جامعه:

هر گاه اطلاعات جمع آوری شده از جامعه ای به حجم  $N$  باشد در این صورت میانگین بدست آمده میانگین جامعه بوده و با نشان می دهیم که در آن :

$$\mu = \frac{\sum_{i=1}^N f_i x_i}{N}$$





## میانگین وزنی

در برخی مطالعات یا اندازه گیریها ممکن است کمیتهای بدست آمده دارای وزن یا ضریب خاصی باشند در این حالت میانگین وزنی محاسبه می شود . که فرمول آن عبارت است از:

$$\bar{x}_w = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

که در آن  $w_i$  وزن داده  $i$  ام  $(x_i)$  می باشد

مثال : دانشجویی در یک درس 3 واحدی نمره 14، در درس 2 واحدی 12،

در درس 4 واحدی 15 و در 3 واحدی دیگر 11 گرفته است.  
معدل وی و یا به بیان آماری میانگین نمرات وی چقدر است؟

$x_i$	$w_i$	$w_i x_i$
14	3	42
12	2	24
15	4	60
11	3	33

$$\bar{x}_w = \frac{\sum w_i x_i}{\sum w_i} = \frac{159}{12} = 13.25$$

مجموع

$$\sum w_i = 12 \quad \sum w_i x_i = 159$$

مثال : در يك شهر سه روزنامه محلي منتشر مي شود. 18 در صد خانواده‌هاي ساكن اين شهر به هيچيك از روزنامه ها مشترك نيستند . 61 در صد آنها يكي از روزنامه‌ها ، 17 درصد دو روزنامه و 4 درصد ديگر هر سه روزنامه مشتركند .

متوسط اشتراك خانواده هاي اين شهر به اين سه روزنامه چقدر است؟

$x_i$	$w_i$	$w_i x_i$
0	18	0
1	61	61
2	17	34
3	4	12

$$\bar{x}_w = \frac{\sum w_i x_i}{\sum w_i} = \frac{107}{100} = 1.07$$

$$\sum w_i = 100 \quad \sum w_i x_i = 107$$

# میانگین هندسی

در صورتی که داده‌ها به صورت نسبت ، نرخ رشد یا از این قبیل باشند از میانگین هندسی استفاده می‌شود که فرمول آن عبارت است از :

$$G = \sqrt[n]{x_1^{f_1} \times x_2^{f_2} \times \dots x_k^{f_k}}$$

در کارهای اقتصادی و جمعیت‌شناسی نیز از این میانگین استفاده می‌کنند.



مثال : قیمت يك کالا در چهار سال گذشته به ترتيب 100 و 15 و 300 و 800 ريال بوده است میانگین قیمت هر سال به سال قبل را تعیین کنید .

$$x_1 = \frac{150}{100}$$

$$x_2 = \frac{300}{150}$$

$$G = \sqrt[3]{\frac{150}{100} \times \frac{300}{150} \times \frac{800}{300}} = \sqrt[3]{8} = 2$$

$$x_3 = \frac{800}{300}$$



مثال : اگر سرمایه شرکت‌ها پس از یکسال 21 درصد افزایش یابد و سال بعد 44 درصد نسبت به سال قبل افزایش یابد متوسط میزان افزایش سرمایه را محاسبه کنید.

پس از يك سال ، سرمایه نسبت به سال قبل  $1/21$  برابر و در سال بعد نسبت به سال قبل  $1/44$  برابر شده است ، لذا يعني میزان افزایش متوسط سرمایه 32 درصد بوده است

$$G = \sqrt[2]{1/21 \times 1/44} = 1/32$$

يعني میزان افزایش متوسط سرمایه 32 درصد بوده است



## میانگین توافقی ( هارمونیک )

چنانچه داده‌ها به صورت میزان تغییرات بیان شده باشند مانند سرعت ، شتاب ، و ... میانگین توافقی مناسب است که فرمول آن عبارت است از :

$$H = \frac{n}{\sum_{i=1}^k \frac{f_i}{x_i}}$$

این میانگین در عینک سنجی و مطالعه برق بکار می‌رود و گاهی نیز به میانگین همسازه معروف است .



مثال : اتومبيلي سه قطعه از يك جاده با طولهاي يكسان را با  
سرعتهاي 75 و 80 و 95 كيلومتر در ساعت طي مي كند متوسط  
سرعت اين اتومبيل در اين سه قسمت عبارت است از :

$$n=3$$

$$x_1 = 75$$

$$x_2 = 80$$

$$x_3 = 95$$

$$H = \frac{3}{\frac{1}{75} + \frac{1}{80} + \frac{1}{95}} = 82/51$$





مثال : جدول توزیع فراوانی مصرف برق 100 خانوار به صورت زیر داده شده است میانگین هارمونیک را حساب کنید

حدود طبقات	$f_i$	$x_i$
60-62	5	61
62-64	18	63
64-66	42	65
66-68	27	67
68-70	8	69

N = 100

$$H = \frac{n}{\sum \frac{f_i}{x_i}} =$$

$$\frac{5}{\frac{5}{61} + \frac{18}{63} + \frac{42}{65} + \frac{27}{67} + \frac{8}{69}} = 68.03$$

## میانگین رتبه دو:

این میانگین اغلب در محاسبات کارهای فیزیکی مانند سیستم‌های قدرت و لثاژها و جریان‌ها کاربرد دارد و به صورت:

$$M_2 = \left( \frac{1}{n} \sum_{i=1}^k f_i x_i^2 \right)^{\frac{1}{2}}$$

تعریف می‌شود درحقیقت این

میانگین جذر میانگین حسابی  $x_1^2, x_2^2, \dots, x_k^2$  می‌باشد.

مثال : براي داده هاي

5، 7، 7، 1، 3، 7، 8، 3، 12، 11، 15، 11، 16،  
19، 15 میانگین رتبه دو عبارت است از :

$$M_2 = \left( \frac{(5)^2 + 3 \times (7)^2 + (1)^2 + 2 \times (3)^2 + (8)^2 + (12)^2 + 2 \times (11)^2 + 2 \times (15)^2 + (16)^2 + (19)^2}{15} \right)^{\frac{1}{2}} = 10/6$$



[www.bookgolden.com](http://www.bookgolden.com)

# میانگین میانگین‌ها

فرض کنید  $k$  نمونه با حجم‌های  $n_1, n_2, \dots, n_k$  با میانگین‌های  $x_1, x_2, \dots, x_k$  باشند در این صورت میانگین میانگین‌ها عبارت است از :

$$\bar{x} = \frac{\sum_{i=1}^k n_i \bar{x}_i}{\sum_{i=1}^k n_i}$$

مثال : میانگین مزد 10 کارگر 350 هزار تومان و میانگین مزد 15 کارگر دیگر 250 هزار تومان می باشد میانگین مزد این 25 کارگر را حساب کنید .

$$\bar{x} = \frac{\sum_{i=1}^k n_i \bar{x}_i}{\sum_{i=1}^k n_i} = \frac{10 \times 350 + 15 \times 250}{25} = 290$$

$$n_1 = 10$$

$$\bar{x}_1 = 350$$

$$n_2 = 15$$

$$\bar{x}_2 = 250$$



ب) نما :

داده‌ای که فراوانی آن نسبت به سایر داده‌ها بیشتر باشد مد یا نما نامیده می‌شود و با  $M$  نشان داده می‌شود .

مثال : نمای داده‌های 3 و 1 و 3 و 1 و 3 و 3 و 2 و 4 و 5 و 3 و 3 و 1

عبارت است از  $M=3$  مثال : جدول زیر تعداد واکسیناسیون انجام شده علیه انواع بیماری‌ها ، در یک سال را نشان می‌دهد مد یا نما عبارت است از :

نوع بیماری	فلج اطفال	سرخک	سل	دیفتري
فراوانی بر	11	3	8	4

فلج اطفال =  $M$

مثال : توزیع فراوانی تعداد افراد خانوار به صورت زیر داده شده است

تعداد افراد خانوار	1	2	3	4	5	6	7
فراوانی	5	10	20	10	40	20	10

در این صورت  $M=5$  چون خانواده‌های 5 نفری بیشترین فراوانی را دارند

مثال : مد یا نما در مجموعه داده‌های 2 و 1 و 3 و 5 و 2 و 5 و 3 و 1 و 3 و 1 وجود ندارد زیرا فراوانی تمام داده‌ها یکسان می‌باشد .

محاسبه نما برای داده‌های طبقه بندی شده :

ابتدا نماینده‌ی طبقه‌ای را که دارای بیشترین فراوانی است و به طبقه‌ی نما معروف است را به عنوان نما اختیار می‌کنند برای دقت بیشتر می‌توان نما را از فرمول

$$M = L_M + \frac{d_1}{d_1 + d_2} \cdot w$$

را بدست آورد .

که در آن :

حد پایین طبقه‌ی نما

اختلاف فراوانی طبقه‌ی نما با طبقه‌ی قبلی

اختلاف فراوانی طبقه‌ی نما با طبقه بعدی

طول طبقه

$L_M$

$d_1$

$d_2$

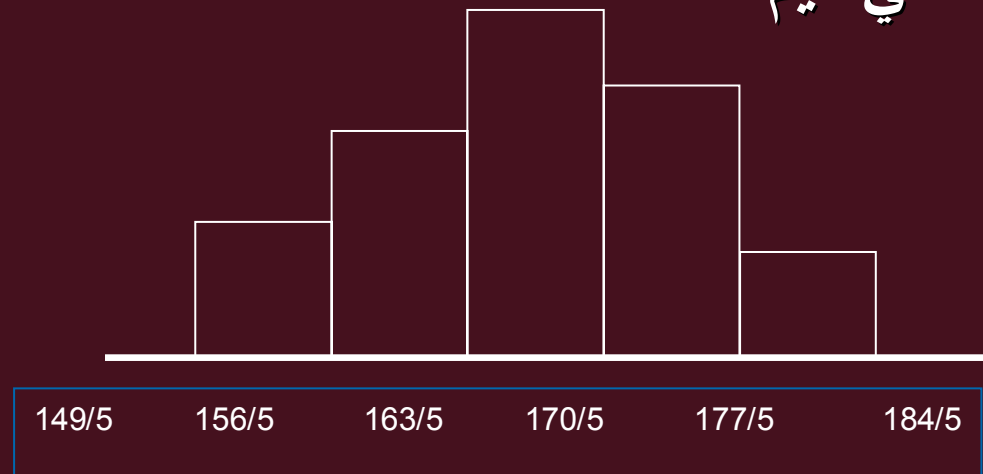
$w$



مثال : در جدول فراوانی داده شده مقدار نما را بیابید.

حدود طبقه	فراوانی
149/5–156/5	15
156/5–163/5	2
163/5–170/5	30
170/5–177/5	25
177/5–184/5	10

مقداری که دارای بیشترین فراوانی است بعنوان طبقه نما انتخاب می‌کنیم



بعنوان مقدار تقریبی نما

$$\frac{163/5 + 170/5}{2} = 167 \text{ یعنی}$$

و

مقدار دقیق آن عبارت است از :

$$M = L_M + \frac{d_1}{d_1 + d_2} \cdot w$$

$$M = 163/5 + \frac{10}{10+5} \times 7$$

$$M = 168/17$$



ج ( میانه :

یکی دیگر از معیارهای گرایش به مرکز میانه می باشد و آن را با  $m$  نشان می دهند میانه داده ای است که تقریباً نصف داده ها از آن کوچکتر باشند.

برای بدست آوردن میانه ابتدا داده ها را به صورت غیر نزولی از کوچک به بزرگ مرتب می کنند اگر تعداد داده ها فرد باشد در این صورت داده ای که در وسط قرار دارد میانه محسوب می شود و اگر تعداد داده ها زوج باشد نصف مجموع دو داده ای که در وسط قرار دارند میانه محسوب می شود .

مثال : براي پيدا كردن ميانه داده‌هاي  
20 و 23 و 20 و 19 و 20 و 21 و 27  
آنها را به صورت غير نزولي 27 و 23 و 21 و 20 و 20 و  
20 و 19  
مرتب مي‌كنيم لذا ميانه عدد وسطي است  $m=20$

مثال : ميانه داده‌هاي 2 و 3 و 5 و 4 و 1 و 2 و 3 و 2 مي‌شود  $m$   
 $= 2/5$

اگر داده‌ها گسسته همراه با فراوانی داده باشند در این صورت با تشکیل ستون فراوانی جمعی، داده‌ای که در وسط قرار دارد را بعنوان میانه انتخاب می‌کنند.

$x_i$	12	14	15	17	19
$f_i$	4	8	6	4	3

مثال: جدول توزیع فراوانی نمرات 25 دانشجو به صورت زیر داده شده است.

$x_i$	$f_i$	$F_i$
12	4	4
14	8	12
15	6	18
17	4	22
19	3	25

با تشکیل ستون فراوانی جمعی و با توجه به  $n=24$  ملاحظه می‌شود که داده‌ی وسطی سیزدهمین داده مرتب شده می‌باشد.

$$M = x_{(13)} = 15$$

محاسبه میانه برای داده‌های طبقه بندی شده :

ابتدا طبقه‌ای را که نصف داده‌ها در آن قرار دارد و یا به عبارتی اولین طبقه‌ای که فراوانی تجمعی آن برابر نصف یا بیشتر از نصف داده‌ها باشد و به آن طبقه میانه می‌گویند را مشخص کرده و سپس از فرمول زیر میانه را محاسبه می‌کنند.



$$m = L_m + \frac{n \times 0.5 - F_{m-1}}{f_m} \times w$$

که در آن

$n$  : تعداد کل داده‌ها

$L_m$  : حد پایین طبقه میانه

$f_m$  : فراوانی طبقه میانه

$F_{m-1}$  : فراوانی تجمعی طبقه‌ی ما قبل طبقه میانه

$w$  : طول طبقه



مثال : توزیع فراوانی وزن 100 نفر دانشجو در جدول زیر داده شده است میانه وزن دانشجویان عبارت است از :

حدود طبقات	فراوانی تجمعی فراوانی	
59-62	5	5
62-65	18	23
65-68	42	65
68-71	27	92
71-74	8	100
		$n = 100$

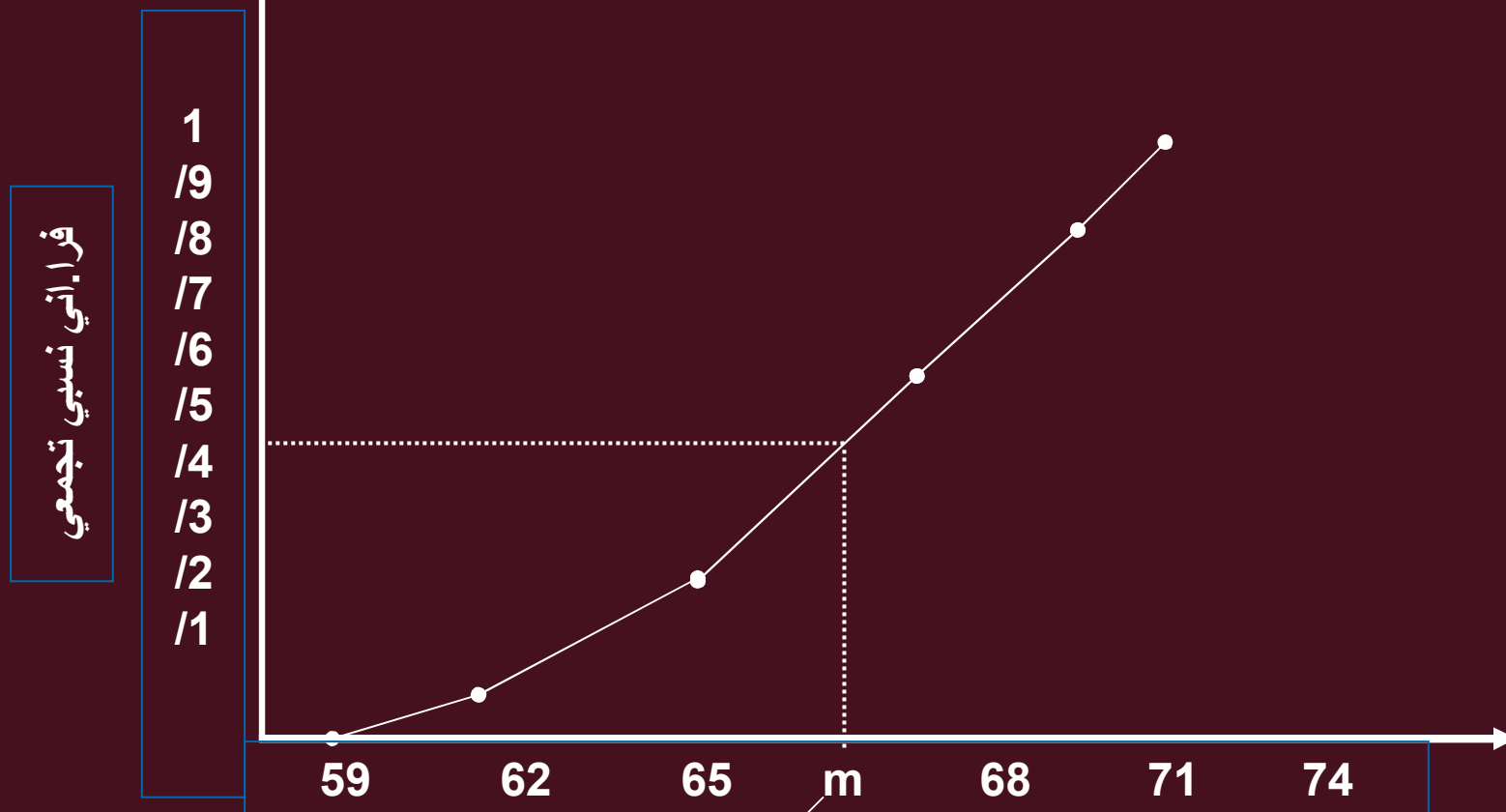
$$m = L_m + \frac{n \times 0.5 - F_{m-1}}{f_m} \times w$$

$$m = 65 + \frac{100 \times 0.5 - 23}{42} \times 3$$

$$m = 66.93$$



مثال: با روش ترسیمی میانه توزیع وزن دانشجویان عبارت است از:



نمودار چند ضلعی نسبی تجمعی وزن دانشجویان

$$m = 66/93$$

خاصیت مهم میانه این است که حاصل جمع قدر مطلق تفاضلهای مقادیر از میانه، از حاصل جمع قدر مطلق تفاضلهای مقادیر از هر عدد دیگر مثل  $c$ ، کوچکتر است یعنی

$$\sum_{i=1}^k f_i |x_i - m| \leq \sum_{i=1}^k f_i |x_i - c|$$

به عبارت دیگر میانه حاصل جمع قدر مطلق تفاضلهای مقادیری از عدد دیگر را به حداقل می رساند.

# فصل پنجم

## چندکها

راهنمای طلایی  
تسبیح طلایی  
پیشک طلایی

انتشارات طلایی  
پویندگان دانشگاه



[www.bookgolden.com](http://www.bookgolden.com)

## چندك ها :

عدد  $Q_p$  را كه در آن  $0 < P < 1$  ، چندك مرتبه  $P$  مي نامند ،  
هرگاه  $100P$  درصد داده ها كوچكتر از آن باشند ، مثلاً  $Q_{./35}$  ،  
چندك مرتبه  $./35$  مي باشد ، هرگاه تقريباً 35 درصد داده ها  
كوچكتر از  $Q_{./35}$  باشد .  
چندك هاي معروف عبارتند از :

### الف) چارك ها

چارك ها كه به ازاي  $./75$  و  $./5$  و  $./25$  بدست مي آيند و آنها را به  
 $Q_1$  و  $Q_2$  و  $Q_3$  نشان مي دهند كه در آن  $Q_2$  همان ميانه مي باشد .

## ب) دهك ها

دهك ها كه به ازاي  $1/9, 1/2, \dots, 1/1$   $P =$  بدست مي آيند و آن ها را با  $D_9, \dots, D_2$  و  $D_1$  نشان مي دهند.

## ج) صدك ها

صدك ها كه بازي  $1/9, 1/2, \dots, 1/1$   $P =$  بدست مي آيند و آنها را با  $P_9, \dots, P_2$  و  $P_1$  نشان مي دهند.



محاسبه چندكها در داده هاي طبقه بندي نشده

فرض كنيد  $n$  داده به صورت غير نزولي  $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$  مرتب شده باشد براي بدست آوردن  $Q_p$  ابتدا  $(n+1)P$  را تشكيل

مي دهيم اگر حاصل ، عددي صحيح مانند  $k$  شود بنابر اين

$$Q = X_{(k)} \text{ در غير اينصورت: } (n+p) = k + r$$

كه در آن  $0 < r < 1$  و قسمت اعشاري  $(n+1)P$  مي باشد در اين

صورت با يك تناسب فاصله اي ساده داريم:

$$Q_p = X_{(k)} + r ( X_{(k+1)} - X_{(k)} )$$



مثال: در داده‌ها زیر ، چارك اول را محاسبه كنيد.

13 و 12 و 9 و 8 و 14 و 16 و 16 و 17 و 14 و 13 و 20 و 16 و 15 و 11 و 10 و 7 و 14 و 6 و 12

ابتدا داده ها را به صورت غير نزولي مرتب مي كنيم

20 و 17 و 16 و 16 و 16 و 15 و 14 و 14 و 14 و 13 و 13 و 12 و 12 و 11 و 10 و 9 و 8 و 7 و 6

در اين صورت چارك اول كه همان صدك 25م مي باشد عبارت

است از:  $Q_1 = Q_{./25}$

$$P = . /25 \quad 19$$

$$(n+1)P = 20 \times . /25 = 4$$

$$n =$$

بنابر اين  $k = 5$

$$Q_{./25} = X_{(5)} = 10$$

و ...

## محاسبه چندك ها در داده هاي طبقه بندي شده

روش محاسبه چندك ها كاملاً شبیه محاسبه میانه است. برای محاسبه چندك  $Q_p$ ، ابتدا عدد  $n.p$  را با ستون فراواني تجمعی مقایسه می‌کنیم اولین طبقه‌ای که فراواني تجمعی آن بزرگتر یا مساوی  $n.p$  باشد را طبقه چندك  $p$ ام نامیده و با استفاده از يك تناسب ساده ریاضی  $Q_p$  را از فرمول زیر محاسبه می‌کنند.

$$Q_p = L_p + \frac{n.p - F_{(p-1)}}{f_p} . w$$



که در آن

$L_p$  : حد پایین طبقه چندانك

$f_p$  : فراواني طبقه چندانك

$F_{p-1}$  : فراواني تجمعي طبقه ماقبل چندانك

$w$  : طول طبقه

مثال: در جدول توزيع فراواني سن ازدواج چارك اول و صدك 60ام عبارتند از:

سن ازدواج	فراواني سن	فراواني تجمعي
14/5 – 19/5	18	18
19/5 – 24/5	74	92
24/5 – 29/5	62	154
29/5 – 34/5	26	180
34/5 – 39/5	20	200

راهنمای طلایی

کتاب تست طلایی

پویندگان دانشگاه

پیشگامی

www.bookgolden.com

$n = 200$

$$Q_1 = Q_{0/25} = 19/5 + \frac{200 \times 0/25 - 18}{74} \times 5 = 21/66$$

$$Q_{60} = Q_{./60} = 24/5 + \frac{20 \times ./6 - 92}{62} \times 5 = 26/76$$

# فصل ششم

## معیارهای پراکندگی

راهنمای طلایی  
تسلی طلایی  
پیشگام طلایی

انتشارات طلایی  
پویندگان دانشگاه



[www.bookgolden.com](http://www.bookgolden.com)

داده‌ها را معمولاً به صورت يك عدد به نام معيار تمرکز خلاصه مي‌کنند و قسمتي از اطلاعات موجود در آنها را در اين عدد منعكس مي‌سازند ولي در تحليلهاي آماري ، معيارهاي مركزي به تنهائي نمي‌توانند ميزان تنوع ، تفاوت و يا به عبارتي ديگر ميزان پراكندگي داده‌ها را بيان كنند.

فرض كنيد دو مجموعه داده B و A عبارتند از :

$A = 1 \text{ و } 2 \text{ و } 3 \text{ و } 4 \text{ و } 6 \text{ و } 8$

$B = 0 \text{ و } 1 \text{ و } 2 \text{ و } 4 \text{ و } 4 \text{ و } 8 \text{ و } 9$

ملاحظه مي شود كه ميانه و نماي هر دو مجموعه ، عدد 4 مي باشد در حالي كه ميزان تجمع و پراكندگي دو مجموعه با هم متفاوت است .

بنابراین معرفی معیارهایی که بتواند بیانگر میزان اختلافها ، تنوعها و میزان پراکندگی مجموعه داده‌ها باشد ضرورت دارد. اینک به معرفی چند معیار مهم برای سنجش پراکندگی می‌پردازیم فرض کنید  $x_1$  و  $x_2$  و ... و  $x_k$  یک سری داده  $n$  تایی با فراوانی‌های  $f_1$  و  $f_2$  و ... و  $f_k$  با میانگین باشند .

$$R = \max(x_i) - \min(x_i) = x_{(n)} - x_{(1)}$$



## الف) دامنه تغييرات

تفاضل كوچكترين داده از بزرگترين داده را كه با  $R$  نشان مي دهند دامنه تغييرات مي گويند .

مثال : حقوق پرداختي ماهيانه كاركنان يك شركت كوچك به شرح زير است .

0 و 6 و 22 و 20 و 17 و 17 و 17 و 16 و

15

دامنه تغييرات  $R = 60 - 15 = 45$  مي باشد .

ب) میانگین قدر مطلق انحرافات :  
 معیار مناسب در سنجش پراکندگی داده‌ها، معیاری است که  
 بتواند میزان انحراف و پراکندگی کل داده‌ها را در مجموعه  
 داده‌ها اندازه بگیرد . پراکندگی و انحراف زمانی مفهوم پیدا می  
 کند که داده ها نسبت به يك مبدأ یا مرکز مقایسه شوند مناسب  
 ترین مرکز برای داده‌ها ، میانگین آنهاست . بنابراین تعریف می  
 کنیم :

مجموع کل قدر مطلق انحرافات  
 =  $\frac{\text{میانگین قدر مطلق انحرافات}}{\text{تعداد کل انحرافات}}$

$$M.A.D = \frac{\sum_{i=1}^n |x_i - \bar{x}|}{n}$$

مثال : میزان مهارت مدیران يك سازمان در برنامه ریزی به صورت زیر است.

45 و 35 و 50 و 65 و 60 و 85 و

$$x_i = 80$$

$$\bar{x} = \frac{\sum x_i}{n} = \frac{420}{7} = 60$$

در این صورت

$x_i$	80	85	60	65	50	35	45
$ x_i - \bar{x} $	20	25	0	5	10	25	15

$$M.A.D = \frac{\sum_{i=1}^n |x_i - \bar{x}|}{n} = \frac{100}{7} = 14.28$$



محاسبه M.A.D در جدولهاي فراواني

اگر جدول فراواني داراي  $K$  طبقه با نماينده  
طبقات

$x_1$  و  $x_2$  و ... و  $x_k$  و فراوانهاي متناظر  $f_1$  و  $f_2$  و ... و  $f_k$   
و  $f_1$  باشد در اين صورت

$$M.A.D = \frac{\sum_{i=1}^k f_i |x_i - \bar{x}|}{n}$$



مثال: در جدول توزیع فراوانی زیر سن 200 مرد زیر 40 سال در اولین ازدواج داده شده است.

سن ازدواج	$f_i$	$x_i$	$f_i x_i$	$ x_i - \bar{x} $	$f_i \times  x_i - \bar{x} $
14/5 - 19/5	18	17	306	8/9	160/2
19/5 - 24/5	74	22	1628	3/9	288/6
24/5 - 29/5	62	27	1674	1/1	682
29/5 - 34/5	26	32	832	6/1	158/6
34/5 - 39/5	20	37	740	11/1	222
مجموع	200		5180		1511/4

$$\bar{x} = \frac{\sum f_i x_i}{n} = \frac{5180}{200} = 25.9$$

$$M.A.D = \frac{\sum f_i |x_i - \bar{x}|}{n} = \frac{1511/4}{200} = 7.557$$

## ج) واریانس :

در اندازه‌گیری مجموع انحرافات و پراکندگی از مرکز داده‌ها، انحرافات منفي متأثر از داده‌هاي کوچکتر از میانگین با انحرافات مثبت خنثي می‌شوند روشي دیگر براي رفع این معضل آن است که انحرافها را مربع کنیم.



در این صورت معیار میانگین مربعات انحرافات از میانگین به عنوان معیار سنجش پراکندگی تعریف می شود. و آن را واریانس می نامیم .

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

با استفاده از جبر مقدماتی می توان فرمول فوق را به صورت خلاصه تر بازنویسی نمود .

$$s^2 = \frac{\sum_{i=1}^n x_i^2}{n} - \bar{x}^2$$

مثال: سود ده ماه يك شركت سرم سازي به اين صورت است.  
4 و 6 و 8 و 2 و 3 و 4 و 7 و 5 و 1

و  $x_i = 0$   
بنابراین

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n} = \frac{60}{10} = 6$$

و نیز به روش خلاصه تر می توان واریانس را محاسبه نمود .

$$s^2 = \frac{\sum_{i=1}^n x_i^2}{n} - \bar{x}^2 = \frac{220}{10} - (4)^2 = 6$$

انحراف استاندارد

$$s = \sqrt{s^2}$$

جذر مثبت واریانس را انحراف استاندارد می گویند

## ■ واریانس جامعه ■

$$\sigma^2 = \frac{\sum_{i=1}^N f_i (x_i - \mu)^2}{N}$$

و گاهی با

$$\sigma_x^2 = \frac{\sum_{i=1}^N f_i (x_i - \mu_x)^2}{N}$$

نشان می دهند

انحراف استاندارد جامعه را با جذر مثبت از واریانس جامعه

می توان بدست آورد

## محاسبه واریانس در جدولهای فراوانی

اگر جدول فراوانی دارای  $K$  طبقه باشد نماینده طبقات را با

$x_1, x_2, \dots, x_k$  و فراوانیهای  $f_1, f_2, \dots, f_k$

نمایش دهند

در این صورت فرمول واریانس را می توان به صورت زیر بازنویسی نمود .

$$s^2 = \frac{\sum_{i=1}^k f_i (x_i - \bar{x})^2}{n}$$

و نیز با استفاده از جبر مقدماتی می توان فرمول فوق را خلاصه تر نمود .

$$s^2 = \frac{\sum_{i=1}^k f_i x_i^2}{n} - \bar{x}^2$$

مثال: توزیع فراوانی سن 200 مرد زیر 40 سال در اولین ازدواج در جدول زیر آمده است .

سن بر حسب سال	$f_i$	$x_i$	$f_i x_i$	$f_i x_i^2$
14/5–19/5	18	17	306	5202
19/5–24/5	74	22	1628	35816
24/5–29/5	62	27	1674	45198
29/5–34/5	26	32	832	26624
34/5–39/5	20	37	740	27380
مجموع	200		5180	140220

$$\bar{x} = \frac{\sum_{i=1}^x f_i x_i}{n} = \frac{5180}{200} = 25/9 \quad \text{سال}$$

$$s^2 = \frac{\sum_{i=1}^k f_i x_i^2}{n} - \bar{x}^2 = \frac{140220}{200} - (25/9)^2 = 30/29$$

$$s = \sqrt{30/29} = 5/5 \quad \text{سال}$$



## (و) دامنه میان چارکي :

در توزیع هایی که دارای تعداد اندکی داده نامعلوم یا پرت ،  
یعني

خیلي کوچک یا خیلی بزرگ در ابتدا و انتهاي منحنی فراواني  
هستند از انحراف چارکي بعنوان شاخص پراکندگی استفاده  
مي شود.

$$IQR = Q_3 - Q_1$$

گاهی اوقات به جاي دامنه میان چارکي از نیمه میان چارکي  
با عنوان

$$SIQR = \frac{Q_3 - Q_1}{2}$$

انحراف چارکي استفاده مي شود که فرمول آن عبارت است  
از:

مثال: برای داده های 12 و 8 و 7 و 7 و 5 و 4 و 3 و 2 و 1 و 0 داریم:

$$IQR = Q_3 - Q_1 = 7 - 2 = 5$$

$$SIQR = \frac{Q_3 - Q_1}{2} = \frac{7 - 2}{2} = 2.5$$

$$IQ = \frac{Q_1 + Q_3}{2} = \frac{2 + 7}{2} = 4.5$$



## هـ) ضریب تغییرات :

نسبت انحراف استاندارد به میانگین یعنی:  $C.V = \frac{s}{x}$

را که اغلب به صورت درصد بیان می‌شود را ضریب تغییرات

می‌نامند ، این معیار که به واحد اندازه‌گیری بستگی ندارد ؛ در

عمل برای مقایسه بکار می‌رود.



مثال : کارخانه ای دو نوع لاستیک تولید می کند. برای نوع A میانگین

عمر 1000 کیلومتر با انحراف استاندارد 200 کیلومتر برای نوع B

میانگین عمر 1100 کیلومتر با انحراف استاندارد 1000 کیلومتر می باشد،

کدام نوع لاستیک بهتر است؟

$$C.V_A = \frac{S_A}{X_A} = \frac{200}{1000} = 0.2 \quad C.V_B = \frac{S_B}{X_B} = \frac{1000}{1100} = 0.91$$

لاستیک نوع A بهتر است زیرا دارای ضریب تغییرات کمتری است.

تأثیر عدد ثابت بر میانگین و واریانس :

هرگاه داده های  $x_1, x_2, \dots, x_n$  را در عدد حقیقی  $a$  ضرب و با عدد حقیقی  $b$  جمع کنیم، داده ای جدید  $ax_1 + b, ax_2 + b, \dots, ax_n + b$  حاصل می شوند که آنها را با  $y_1, y_2, \dots, y_n$  نشان می دهیم

در این صورت میانگین داده های جدید عبارت است از:

$$\bar{y} = \frac{\sum_{i=1}^n y_i}{n} = \frac{\sum_{i=1}^n (ax_i + b)}{n} = a\bar{x} + b$$



وبه همین شیوه واریانس داده های جدید عبارت است از:

$$S_y^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n} = \frac{\sum_{i=1}^n ((ax_i + b) - (a\bar{x} + b))^2}{n}$$
$$= a^2 \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n} = a^2 s_x^2$$



ضریب تغییرات چارکی :

در توزیعهایی که دارای اندکی داده نامعلوم یا پرت در

ابتدا یا انتهای منحنی فراوانی هستند از معیار ضریب  
تغییرات چارکی استفاده می شود که فرمول آن عبارت  
است از:

$$C.V.Q = \frac{Q_3 - Q_1}{Q_3 + Q_1}$$



مثال: توزیع فراوانی سن افراد جویای کار در جدول زیر آمده است.

	فراوانی	فراوانی تجمعی
کمتر از 12 سال	2	2
12-16	18	20
16-20	40	60
20-30	60	120
30-40	50	170
40-60	20	190
بیش از 60 سال	10	200

$$Q_1 = 19$$

$$Q_3 = 36$$

در این صورت

$$C.V.Q = \frac{Q_3 - Q_1}{Q_3 + Q_1} = \frac{36 - 19}{36 + 19} = \frac{17}{55} = 0.31$$



گشتاور و گشتاور مركزي داده ها :

فرض كنيد  $n$  داده به صورت  $x_1, x_2, \dots, x_k$  با فراواني  $f_1, f_2, \dots, f_k$  هاي باشند ، ميانه (توان  $\bar{x}$ )  $r$ ام  $x_i$  ها و هاي يعني:

$$m_r = \frac{\sum_{i=1}^k f_i (x_i - \bar{x})^r}{n}, \quad m_r^1 = \frac{\sum_{i=1}^k f_i x_i^r}{n}$$

را به ترتيب گشتار  $r$ ام و گشتاور مركزي  $r$ ام داده ها مي نامند كه در آن  $r$  يك عدد طبيعي است واضح است كه

$$m_1^1 = \bar{x}$$

$$m_1 = 0$$

$$m_2 = s^2$$

مثال: گشتاور مرکزی مرتبه اول تا چهارم برای داده‌های 10 و 9 و 8 و 7 و 6 عبارتند از:

$$\bar{x} = \frac{40}{5} = 8$$

$$m_1 = \frac{0}{5} = 0$$

$$m_2 = \frac{10}{5} = 2$$

$$m_3 = \frac{0}{5} = 0$$

$$m_4 = \frac{34}{5} = 6/8$$

$x_i$	$(x_i - \bar{x})$	$(x_i - \bar{x})^2$	$(x_i - \bar{x})^3$	$(x_i - \bar{x})^4$
6	-2	4	-8	16
7	-1	1	-1	1
8	0	0	0	0
9	1	1	1	1
10	2	4	8	16
40 مجموع	0	10	0	34

## تبدیل یا کد گذاری داده ها :

تبدیل یا کد گذاری مجموعه ای از داده های ناجور (خیلی کوچک یا خیلی بزرگ) عبارت از عملیاتی است که طی آن از هر مشاهده ، عدد ثابتی را کم (یا به آن اضافه) کرده و نتیجه را

بر عدد ثابتی تقسیم (یا در آن ضرب) می نمایند تا اعداد مناسب جهت محاسبات ساده تر حاصل آید.



اگر داده اولیه را با  $x_i$  نشان دهیم با انتخاب مقادیر ثابت  $a$  و  $b$  داده تبدیلی متناظر به صورت  $u_i = \frac{x_i - a}{b}$  تعریف می شود.

بنابر این

$$\bar{x} = b\bar{u} + a$$

$$s_x^2 = \frac{\sum (x_i - \bar{x})^2}{n} = b^2 \frac{\sum (u_i - \bar{u})^2}{n} = b^2 s_u^2$$

روش کد گذاری را در مورد داده های طبقه بندی شده  
نیز به

کاربرد . نماینده طبقه ای که در وسط جدول فراوانی قرار  
داد را به

عنوان مقدار  $a$  و طول طبقه را به جای  $b$  فرض می کنیم.  
در جدول های

فراوانی که تعداد طبقات زوج می باشد نماینده یکی از دو طبقه  
وسط

(به دلخواه) را بعنوان  $a$  برمی گزینیم.



مثال: در جدول فراوانی زیر با انتخاب  $a = 17500$  و  $b = 5000$

حدود طبقات	$f_i$	$x_i$	$u_i$	$f_i u_i$	$f_i u_i^2$
5000-10000	10	7500	-2	-20	40
10000-15000	12	12500	-1	-12	12
15000-20000	35	17500	0	0	0
20000-25000	30	22500	1	30	30
25000-30000	13	27500	2	26	52

$$\sum f_i u_i = 24 \quad \sum f_i u_i^2 = 134$$

$$\bar{U} = \frac{\sum f_i u_i}{n} = \frac{24}{100} = 0.24$$

$$S_u^2 = \frac{\sum f_i u_i^2}{n} - \bar{u}^2 = \frac{134}{100} - (0.24)^2$$

$$\bar{x} = b\bar{u} + a = 500(0.24) + 7500 = 8700$$

$$s_x^2 = b^2 s_u^2 = (500)^2 \times \frac{1}{28}$$

# فصل هفتم

## شاخصهای شکل توزیع

راهنمای طلایی  
تسلی طلایی  
پیشگامی

انتشارات طلایی  
پویندگان دانشگاه



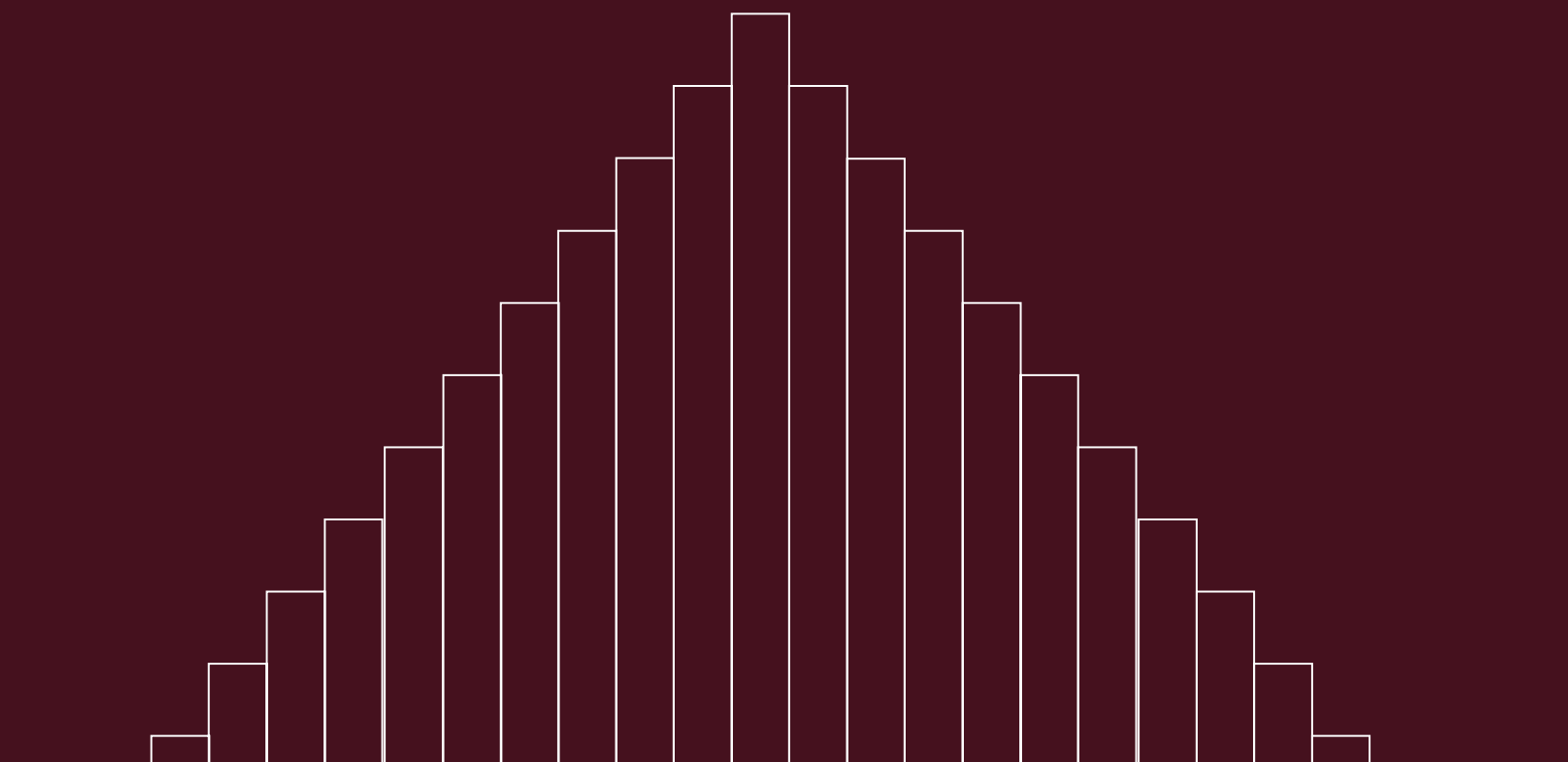
[www.bookgolden.com](http://www.bookgolden.com)

## منحنی متقارن زنگوله‌ای (نرمال)

اغلب در مواجهه با انبوهی از داده های آماری، آنها را در يك جدول فراوانی خلاصه کرده و پس از آنکه هیستوگرام مربوطه را رسم کردیم بر آن منحنی برآزش می‌دهیم ،

مشاهده خواهیم کرد منحنی قرینه زنگوله‌ای شکل بدست خواهد که به آن منحنی طبیعی یا نرمال گفته می شود . منحنی نرمال و مهمی برای توصیف خیلی از مجموعه های آماری است .





## هیستوگرام مقارن

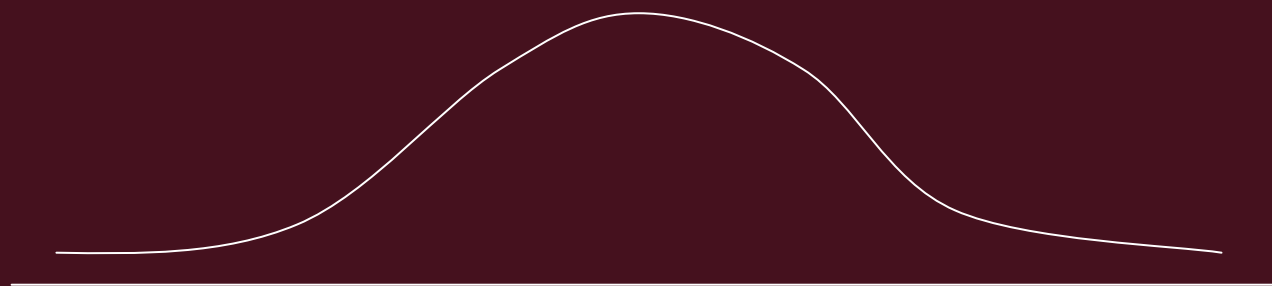
راهنمای طلایی  
تست طلایی  
پیک طلایی

انتشارات طلایی  
پویندگان دانشگاه



[www.bookgolden.com](http://www.bookgolden.com)

منحني فراواني متناظر آن به شكل زير است .



منحني متقارن زنگوله اي ( نرمال )

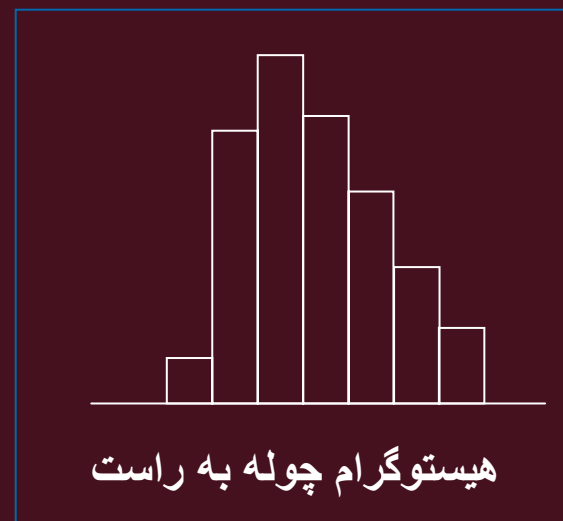
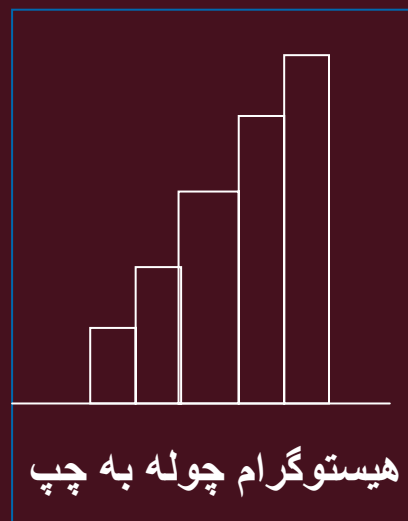
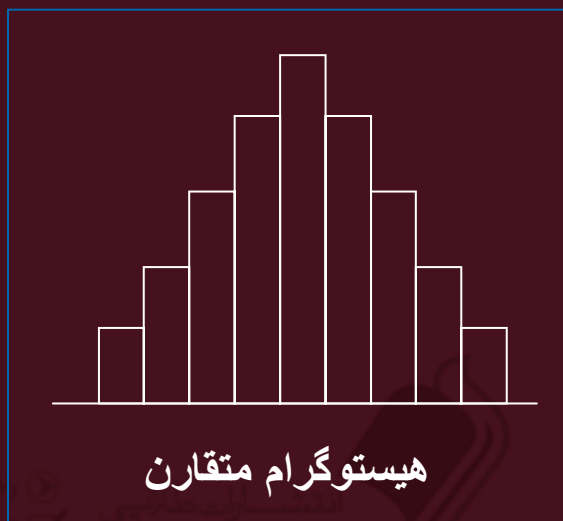


مفهوم توزیع طبیعی در آمار يك مفهوم اساسي و بنيادي است. دليل اهميت آن اين است كه توزيع فراواني هيستوگرام بسياري از حوادث طبيعي شكل منحنی طبيعي يا نرمال را دارند يعني فراواني آنها ابتدا به يك مقدار پيشينه صعود مي‌كند و سپس به طور متقارن کاهش مي‌يابد.



## چولگی:

میزان عدم تقارن منحنی را چولگی می‌گویند. در توزیعهای چوله معیارهای گرایش مرکزی بر هم منطبق نیستند ولی همیشه میانه بین نما و میانگین قرار دارد. اگر نما از میانگین بزرگتر باشد منحنی چوله به چپ است یعنی دم بلند منحنی به طرف چپ می‌باشد و اگر میانه از میانگین کوچکتر باشد منحنی چوله به راست می‌باشد.



تفاوت میانگین و نما  $(\bar{x} - M)$  یا تفاوت آن با میانگین  $(\bar{x} - m)$  می‌تواند مشخص کننده چولگی باشد حال این اختلافها با تقسیم بر یک معیار پراکندگی هم‌بعد مانند S از قید واحد خارج می‌شوند .

با این انگیزه پیرسون برای اندازه گیری میزان و جهت چولگی ضرائب چولگی زیر را تعریف نمود.



ضریب اول چولگی پیرسون	$sk_1 = \frac{\overline{x} - M}{s}$
ضریب اول چولگی پیرسون	$sk_2 = \frac{3(\overline{x} - m)}{s}$
ضریب اول چولگی پیرسون	$sk_3 = \frac{m_3}{s^3}$

که در آن  $m_3$  ، گشتاور مرکزی مرتبه سوم می باشد

در توزیعهایی که چولگی زیاد نباشد رابطه تجربی  $\overline{x} - m \approx 3(\overline{x} - M)$  که به رابطه پیرسون شهرت دارد برقرار است.

## تفسير ضرايب چولگي

اگر  $sk=0$  باشد توزيع چولگي ندارد.

اگر  $sk > 0$  باشد توزيع چوله به راست است .

اگر  $sk < 0$  باشد توزيع چوله به چپ است .

در عمل ضرائب چولگي با اعداد  $0/1$  و  $0/5$  مقايسه مي شوند  
بدین ترتيب که :

اگر  $|sk| < 0/1$  باشد توزيع چوله نيست و مي تواند متقارن  
باشد.

اگر  $0/5 < |sk| < 0/1$  باشد توزيع چوله است و هرچه  $sk$   
به  $0/1$  نزديك باشد چولگي كم و هر چه به  $0/5$  نزديك باشد  
چولگي زياد است .

اگر  $|sk| > 0/5$  باشد توزيع خيلي چوله است.

کرده‌اند و به شرح زیر است. ضریب چولگی اول پیرسون را محاسبه و تفسیر کنید .

3	4	5	6	7
40	24	20	8	4

$x_i$  2

$f_i$  4

ابندا معیارهای گرایش مرکز را محاسبه می‌کنیم

$x_i$	$f_i$	$F_i$	$f_i x_i$	$f_i x_i^2$
2	4	4	8	16
3	40	44	120	360
4	24	68	96	384
5	20	88	100	500
6	8	96	48	288
7	4	100	28	196

$$n = 100 \quad \sum f_i x_i = 400 \quad \sum f_i x_i^2 = 1744$$



میانگین

نما

$$\bar{x} = \frac{\sum f_i x_i}{n} = \frac{400}{100} = 4$$

$$M = 3$$

$$s = \sqrt{\frac{\sum f_i x_i^2}{n} - \bar{x}^2}$$

$$s = \sqrt{\frac{1744}{100} - (4)^2} = 1/2$$

$$sk_1 = \frac{\bar{x} - M}{s} = \frac{4 - 3}{1/2} = +2$$

بنابراین

چولگی زیاد و مثبت است.



نمودار چوله به راست است

مثال: طول عمر 100 باطري اتومبيل داراي ميانگين،ميانه و انحراف استاندارد  $3/5$  و  $3/48$  و  $1/65$  سال مي باشد ضريب چولگي با توجه به اطلاعات داده شده عبارت است از :

$$\bar{x} = 2/5 \quad m = 3/48 \quad s = 1/65$$

$$sk_1 = \frac{3(\bar{x} - m)}{s}$$

$$= \frac{3(3/5 - 3/48)}{1/65} = 0/036$$

## ضریب چولگی باولی :

در توزیع‌هایی که دارای تعداد اندکی داده نا معلوم یا پرت در ابتدا یا انتهای منحنی فراوانی هستند از این معیار که براساس انحرافات چارکی می‌باشد برای اندازه‌گیری میزان چولگی استفاده می‌شود .

که به آن ضریب چولگی باولی یا ضریب چولگی می‌گویند و با  $sk_B$  نمایش می‌دهند.

$$sk_B = \frac{(Q_3 - m) - (m - Q_1)}{(Q_3 - m) + (m - Q_1)}$$

پس :

$$sk_B = \frac{Q_3 + Q_1 - 2m}{Q_3 - Q_1}$$

تفسیر چولگی باولی :

اگر  $sk_B > 0$  توزیع چوله به راست است.

اگر  $sk_B < 0$  توزیع چوله به سمت چپ است.

اگر  $sk_B = 0$  توزیع می‌تواند متقارن باشد یا متقارن فرض شود.

ضریب چولگی باولی در فاصله‌ی ( +1 و -1 ) تغییر می‌کند و حال  
ن که ضریب چولگی پیرسون در فاصله‌ی ( +3 و -3 ) تغییر می‌کند و  
از این جهت است که ضریب باولی گاهی اوقات ترجیح داده می‌شود .



گاهی اوقات نیز از ضریب چولگی صدکی نیز برای  
اندازه‌گیری میزان چولگی استفاده می‌شود که فرمول آن  
عبارت است از :

$$sk_P = \frac{P_{90} - 2P_{50} + P_{10}}{P_{90} - P_{10}}$$

که در آن  $P_{10}$  و  $P_{50}$  و  $P_{90}$  به ترتیب صدک‌های دهم ،

پنجاهم و نودم توزیع می‌باشند  
مکان : در  $P_{10}$  و  $P_{50}$  و  $P_{90}$  توزیع مقدار ضریب چولگی صد

برابر است با :  $P_{10} = 123/33$  و  $P_{50} = 162/25$  و  $P_{90} = 201/5$

$$sk_P = \frac{201/5 - 2 \times 122/25 + 123/33}{201/5 - 123/33} = 0/004$$

کشیدگی :

دو توزیع ممکن است دارای میانگین یکسان باشند و کاملاً نیز قرینه باشند ولی از لحاظ کشیدگی با هم تفاوت زیادی داشته باشند .

میزان برجستگی یا پخی منحنی فراوانی نسبت به منحنی مقایسه را کشیدگی گویند و منحنی مقایسه را منحنی نرمال در نظر می گیرند .

$$m_4 = \frac{1}{4} \sum f_i (x_i - \bar{x})^4$$

باشد . لذا در توزیع نرمال  $\frac{m_4}{s^4} - 3 = 0$  و در توزیع های غیر نرمال این مقدار برابر صفر نیست و آن را ضریب کشیدگی می نامند و با  $K$  نشان می دهند .

## تفسیر ضریب کشیدگی :

توزیع نرمال است هرگاه  $k = 0$  باشد.

توزیع از توزیع نرمال کشیده‌تر است هرگاه  $k > 0$  باشد.

توزیع از توزیع نرمال پخ‌تر است هرگاه  $k < 0$  باشد.

توزیع تقریباً نرمال است هرگاه  $|k| < 0.1$  توزیع یا

توزیع نرمال تفاوت فاحش دارد هرگاه  $|k| > 0.5$

$|k| > 0.1$  باشد.

تفاوت توزیع با توزیع نرمال فاحش است هرگاه  $|k| > 0.5$

$|k|$

مثال : براي جدول توزيع فراواني زير ، ضريب کشيدگي را محاسبه کنيد .

1	2	3	4	5
20	10	40	10	20

$$K = \frac{m_4}{s^4} - 3$$



$x_i$	$f_i$	$f_i \times x_i$	$(x_i - \bar{x})$	$f_i(x_i - \bar{x})^2$	$f_i(x_i - \bar{x})^4$
1	20	20	-2	80	320
2	10	20	-1	10	10
3	40	120	0	0	0
4	10	40	1	10	10
5	20	100	2	80	320

$$n=100 \quad \sum f_i x_i = 300$$

$$\sum f_i(x_i - \bar{x})^2 = 180 \quad \sum f_i(x_i - \bar{x})^4 = 660$$

$$\bar{x} = \frac{\sum f_i x_i}{n} = \frac{300}{100} = 3$$

$$s = \sqrt{\frac{\sum f_i(x_i - \bar{x})^2}{n}} = \sqrt{\frac{180}{100}} = 1/34$$

$$m_4 = \frac{\sum f_i (x_i - \bar{x})^4}{n} = \frac{660}{100} = 6.6$$

$$K = \frac{m_4}{s^4} - 3 = \frac{6.6}{(1.34)^4} - 3 = -0.952$$

تفاوت توزیع با توزیع نرمال فاحش است و منحنی توزیع از منحنی نرمال پخ‌تر (کوتاه‌تر) است .

ضریب کشیدگی صدکی :

برای تحلیل توزیع‌هایی که با استفاده از گشتاورها قابل توصیف نیستند و با چنډک‌ها توصیف می‌شوند از معیار ضریب کشیدگی چنډکی برای اندازه‌گیری میزان کشیدگی استفاده می‌شود .



در توزیع منحنی نرمال ثابت می‌شود که مقدار  $\frac{Q_3 - Q_1}{2(P_{90} - P_{10})}$  که در آن  $P_{90}$  و  $P_{10}$  به ترتیب صدک‌های نودم و دهم می‌باشند برابر 0/263 است .

با این وصف می‌توان معیار ضریب کشیدگی صدکی را به صورت  $K_P = \frac{Q_3 - Q_1}{2(P_{90} - P_{10})} - 0/263$  برای اندازه‌گیری میزان کشیدگی تعریف نمود هر چه این معیار به صفر نزدیک‌تر باشد کشیدگی توزیع به نرمال نزدیک‌تر می‌باشد .



مثال : ضریب کشیدگی صدکي را براي داده‌ها طبقه‌بندی شده جدول زیر :

حدود طبقات	فراواني تجمعي	فراواني
<5	2	2
5 – 8	10	12
8 – 11	25	37
11 – 14	23	60
14 – 17	15	75
$\geq 17$	5	80

$$Q_3 = 14 \text{ و } Q_1 = 8/96$$

$$P_{90}=16/4 \text{ و } P_{10}=6/8$$

$$K_P = \frac{Q_3 - Q_1}{2(P_{90} - P_{10})} - 0/263$$

$$K_P = \frac{\frac{14 - 8/96}{2}}{16/4 - 6/8} - 0/263$$

$$K_P = -/0005$$

ضریب شکستگی نشان می‌دهد که توزیع منحنی نرمال است.

# فصل هشتم

## همبستگی

راهنمای طلایی  
تسلی طلایی  
پیشگام طلایی

انتشارات طلایی  
پویندگان دانشگاه



[www.bookgolden.com](http://www.bookgolden.com)

# همبستگی:

در بسیاری از مطالعات، منظور از بررسی متغیرها، مشخص کردن روابط موجود بین این متغیرهاست. مثلاً سالها این سوال مطرح بود که آیا بین استعمال دخانیات و سرطان‌های دستگاه تنفس رابطه‌ای وجود دارد یا خیر؟

سرانجام با بررسی‌های آماری معلوم شد که رابطه‌ای شدید بین این دو متغیر وجود دارد و تبلیغ منع استعمال سیگار شدت بیشتری یافت. مطالعه رابطه بین متغیرها را در اصطلاح آماری مطالعه همبستگی می‌گویند.

ساده‌ترین نوع همبستگی بررسی همبستگی بین  
دو متغیر می‌باشد متغیری که تغییرات آن تحت کنترل  
نباشد و تغییرات آن دستخوش تغییرات متغیر دیگر  
باشد را متغیر وابسته

یا پاسخ و متغیری که تغییرات آن توسط آزمایشگر قابل  
کنترل باشد و دستخوش تغییرات متغیر دیگر  
نباشد را متغیر مستقل یا کنترل می‌نامند.



[www.bookgolden.com](http://www.bookgolden.com)

دربررسی رابطه بین هزینه تبلیغات و افزایش میزان فروش  
هزینه تبلیغات متغیر مستقل و افزایش میزان فروش يك  
متغیر وابسته است

و در بررسی مقدار ماده افزودنی به بنزین و مقدار کاهش  
اکسید ازت، مقدار ماده افزودنی متغیر کنترل یا مستقل و  
مقدار کاهش اکسید ازت يك متغیر وابسته یا پاسخ است .  
متغیر کنترل یا مستقل را با  $X$  و متغیر وابسته یا پاسخ  
را با  $y$  نمایش می‌دهند.





## انواع همبستگی :

الف) در مواردی که رابطه بین دو دسته متغیر رابطه‌ای مستقیم باشد یعنی این که افزایش داده ها يك دسته متغیر با افزایش داده های نظیر آنها در متغیر دیگر همراه باشد، اصطلاحاً گفته می شود که میان این دو متغیر رابطه مثبت وجود دارد.



ب) گاهی رابطه بین دو متغیر به گونه‌ای است که تغییر یکی از متغیرها

با تغییر متغیر دیگر اما در جهت عکس آن همراه باشد یعنی کاهش یکی با افزایش دیگری همراه است. و بالعکس. در چنین مواردی گفته می‌شود بین این دو دسته متغیر رابطه معکوس یا همبستگی منفی وجود دارد.

ج) بعضی وقتها بین دو دسته متغیر هیچ نوع رابطه‌ای خواه مستقیم

و یا خواه غیرمستقیم وجود ندارد. در این صورت گفته می‌شود که همبستگی آنها صفر است.

## جمع آوری داده ها :

در بحث همبستگی ساده تنها دو یا چند متغیر  $X$  (متغیر مستقل) و  $Y$

داریم ، برای جمع آوری داده های مورد نیاز اگر متغیر  $X$  توسط آزمایشگر قابل کنترل باشد به آزادی مقادیر مختلف  $X$  مانند  $x_1, x_2, \dots, x_n$  مقادیر  $y$  را مشاهده و با  $y_1, y_2, \dots, y_n$  نشان می دهیم و اگر متغیر  $X$  قابل

کنترل نباشد باید نمونه ای تصادفی از جامعه مورد مطالعه انتخاب و مقادیر  $X$

و  $y$  را برای هر مورد مشاهده و به صورت زوجهای

$x$	$x_1$	$x_2$	$\dots$	$x_n$
$y$	$y_1$	$y_2$	$\dots$	$y_n$

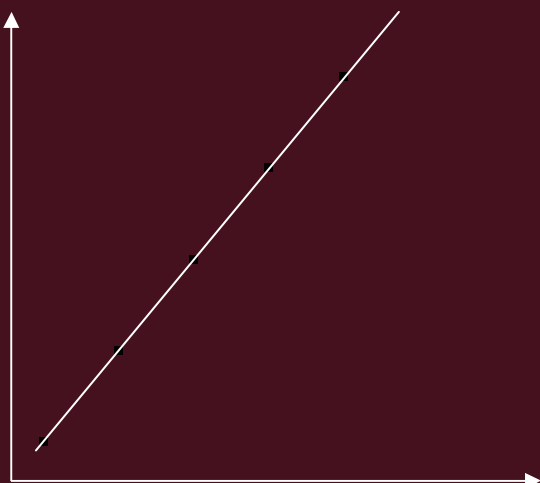
یا  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

## نمودار پراکنش :

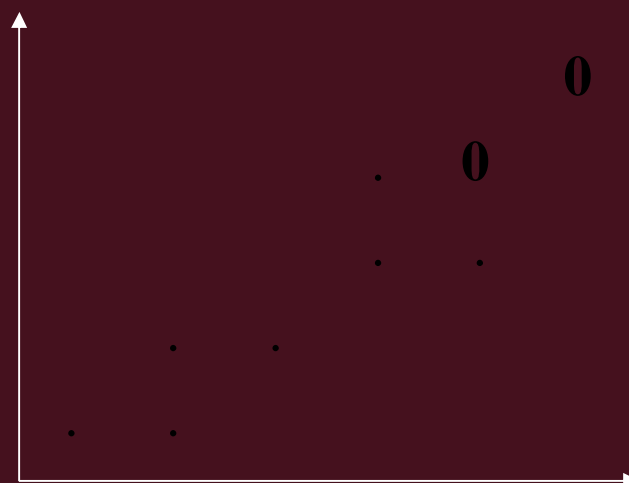
در مطالعه رابطه بین دو متغیر، اولین قدم منطقی، رسم داده‌ها به صورت نقاطی بر روی محور مختصات دو بعدی دکارتی است. نمودار پراکنش، بینشی درباره ماهیت رابطه‌ای که به وسیله داده‌ها نمایش داده می‌شود فراهم می‌کند

گاهی اوقات اطلاع منعکس در نمودار پراکنش برای جستجوی یک مدل ریاضی مناسب (خط راست یا یک نوع منحنی مشخص) مفید است.

در شکل‌های زیر نمایش ترسیمی همبستگی‌ها در حالت‌های مختلف نشان داده شده است.



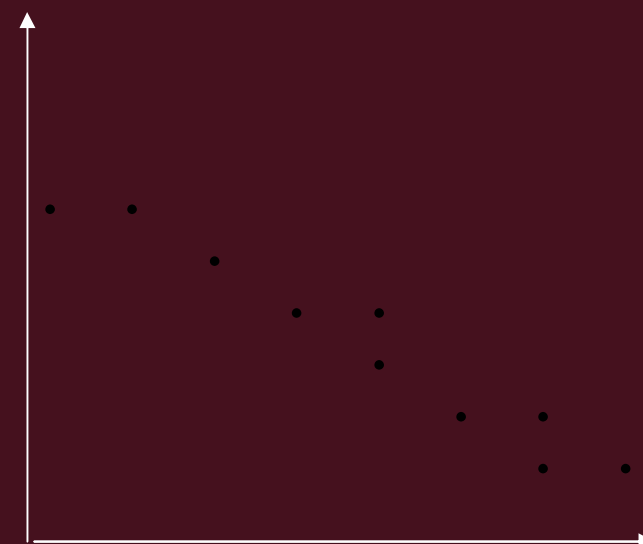
همبستگی کامل مثبت



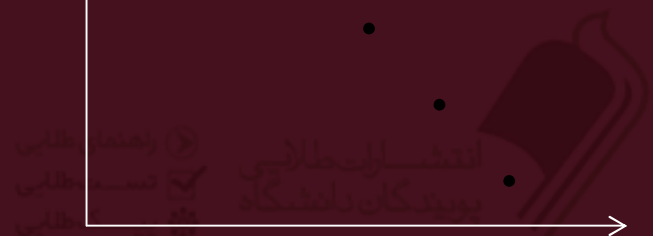
همبستگی ناقص و مثبت



همبستگی کامل و منفی



همبستگی ناقص و منفی



## ضرائب همبستگی :

فرض کنید زوجهای  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  مربوط  
به ساعات مطالعه قبل از آزمون ونمره آزمون باشند. می  
خواهیم بررسی کنیم که آیا بین  $x$  و  $y$  همبستگی خطی وجود  
دارد یا خیر؟

با رسم نمودار پراکنش تا حدودی می توان وجود همبستگی  
و نوع

آن را تشخیص داد اما ما به دنبال معیاری عددی هستیم که  
بتواند

جهت و میزان همبستگی بین دو متغیر را تعیین کند این  
معیار را

ضریب همبستگی گشتاوری پیرسون:

هنگامیکه داده ها دارای مقیاس فاصله ای یا نسبی هستند از این ضریب استفاده می کنند که توسط کارل پیرسون تهیه تنظیم گردیده است و به نام ضریب همبستگی پیرسون معروف است.

$$r_{x,y} = \frac{S_{xy}}{\sqrt{S_{xx}} \cdot \sqrt{S_{yy}}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

فرمول فوق را می توان با محاسبات ساده ریاضی به صورت

$$r_{x,y} = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sqrt{\sum_{i=1}^n x_i^2 - n \bar{x}^2} \times \sqrt{\sum_{i=1}^n y_i^2 - n \bar{y}^2}}$$

## خواص ضرب همبستگی :

(الف)  $r$  باید در فاصله  $[-1, +1]$  قرار گیرد.

(ب) مقدار عددی  $r$  قدرت رابطه خطی و علامت  $r$  نمایانگر جهت این رابطه است.

(ج) اگر مقادیر  $x$  با  $ax+b$  و مقادیر  $y$  با  $cy+d$  عوض شوند و در آنها  $a$  و  $c$  دارای علامت یکسان باشند  $r$  تغییر نمی‌کند یا به عبارتی با تغییر مبدأ و واحد اندازه گیری ضرب همبستگی تغییر نمی‌کند.



مثال: به عنوان قسمتی از مطالعه همبستگی بین تن و روان در موفقیت ورزشکاران، اندازه های زیر از اعضاء تیم کشتی المپیک ایران بدست آمده است.

میزان خشم (x)	6	7	5	21	13	5	13	14
میزان نیرو مندی (y)	28	23	29	22	20	19	28	19

ضریب همبستگی را بیابید.

$x$	$y$	$x \cdot y$	$x^2$	$y^2$
6	28	168	36	784
7	23	161	49	529
5	29	145	25	841
21	22	462	441	484
13	20	260	169	400
5	19	95	25	361
13	28	364	169	784
14	19	266	196	361

$$\sum x_i = 84 \quad \sum y_i = 188 \quad \sum x_i y_i = 1921 \quad \sum x_i^2 = 1110 \quad \sum y_i^2 = 4544$$

$$s_{xy} = \sum x_i y_i - n \bar{x} \bar{y} = 19248 - 108 \times 235 = -53$$

$$s_{xx} = \sum x_i^2 - n \bar{x}^2 = 11108 - 108 \times (105)^2 = 228$$

$$s_{yy} = \sum y_i^2 - n \bar{y}^2 = 45448 - 108 \times (235)^2 = 126$$

$$r_{x.y} = \frac{s_{x.y}}{\sqrt{s_{xx}} \sqrt{s_{yy}}} = \frac{-53}{\sqrt{228} \sqrt{126}} = -0.31$$

همبستگی معکوس و منفی است.

## ضریب تعیین :

ضریب همبستگی اندازه همبستگی بین دو متغیر را نشان می دهد اما این شاخص در باره ماهیت این همبستگی اطلاعات زیادی به ما نمی دهد. آن چه این روش مشخص می کند وجود همبستگی و یا بالا و پایین بودن نسبی آن است.

اطلاعات کامل تر در باره همبستگی وقتی امکان پذیر است که شاخص دیگری به نام ضریب تعیین محاسبه می شود که فرمول آن عبارت است از:

$$R^2 = r_{xy}^2 \times 100$$

## خواص ضریب تعیین :

$$0 \leq R^2 \leq 1$$

(الف)

ب) هرچه  $R^2$  به يك نزديك باشد میان متغیرها يك رابطه خیلی قوي برقرار است.

هرچه  $R^2$  به يك نزديك باشد بهتر مي توان  $y$  را به كمك يك رابطه خطي با داشتن  $x$  پيش بيني يا تعيين كرد از اين رو آنرا ضريب تعيين مي گويند.

مثال : برای محاسبه ضریب تعیین بین دو متغیر هزینه مسکن (y) و درآمد (x)، نمونه ای به حجم  $n=10$ ، انتخاب و براساس داده‌ها ، مقادیر زیر محاسبه شده است. مقدار  $R^2$  را محاسبه کنید.

$$\sum x_i y_i = 180 \quad \sum y_i^2 = 210 \quad \sum y_i = 40 \quad \sum x_i^2 = 300 \quad \sum x_i = 50$$

$$S_{xx} = \sum x_i^2 - n\bar{x}^2 = 50$$

$$S_{yy} = \sum y_i^2 - n\bar{y}^2 = 50$$

$$S_{xy} = \sum x_i y_i - n\bar{x}\bar{y} = -20$$

$$r_{xy} = \frac{S_{xy}}{\sqrt{S_{xx}} \sqrt{S_{yy}}} = \frac{-20}{\sqrt{50} \times \sqrt{50}} = -0.4$$

$$R^2 = r_{x,y}^2 \times 100 = (-0.4)^2 \times 100 = 16$$

یعنی 16 درصد تغییر پذیری کل به علت رگرسیون روی X و بقیه به علت خطای پیش بینی (سایر عوامل غیر از درآمد) می باشد.

## ضریب همبستگی رتبه ای اسپرمن :

ضریب همبستگی پیرسون، برای توصیف همبستگی بین دو متغیر که با استفاده از مقیاس فاصله ای یا نسبی اندازه گیری شده باشند به کار برده می شود.  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

فرض کنید یک نمونه تصادفی زوجی

باشد در این روش رتبه  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  را با  $(R_1, S_1), (R_2, S_2), \dots, (R_n, S_n)$  می دهیم. اینک به جای نمونه بالا نمونه زوجی رتبه ای

بدست می آید اینک ضریب همبستگی ساده پیرسون را با زوجهای رتبه ای محاسبه می کنیم این ضریب همبستگی را که اسپرمن در سال

1904 معرفی کرده است و با  $r_{sp}$  نشان داده ضریب همبستگی

اسپرمن یا همبستگی رتبه ای می نامیم بنابراین:

$$r_{sp} = \frac{\sum_{i=1}^n (R_i - \bar{R})(S_i - \bar{S})}{\sqrt{\sum_{i=1}^n (R_i - \bar{R})^2} \sqrt{\sum_{i=1}^n (S_i - \bar{S})^2}}$$

حال اگر تعریف کنیم  $D_i = R_i - S_i$  آنگاه می توان فرمول ضریب

همبستگی را به صورت زیر بازنویسی نمود:

$$r_{sp} = 1 - \frac{\sum_{i=1}^n D_i^2}{n(n^2 - 1)}$$

مثال : داده های زیر تعداد ساعات مطالعه ده دانشجو و نمره های آنها در یکی از دروس آمار می باشد. یافته  $R_s$  را پیدا کنید.



x: تعداد ساعات	y: نمرة	x رتبه: r <sub>i</sub>	y رتبه: s <sub>i</sub>	D <sub>i</sub>	D <sub>i</sub> <sup>2</sup>
8	56	4 / 5	4	0 / 5	0 / 25
5	44	2 / 5	2	0 / 5	0 / 25
11	79	7	8	- 1	1
13	72	8	7	1	1
10	70	6	6	0	0
5	54	2 / 5	3	- 0 / 5	0 / 25
18	94	10	10	0	0
18	85	9	9	0	0
2	33	1	1	0	0
8	65	4 / 5	5	- 0 / 5	0 / 25

$$r_{sp} = 1 - \frac{\sum_{i=1}^n D_i^2}{n(n^2 - 1)} = 1 - \frac{18}{990} = 0/98$$

www.bookgolden.com



## ضریب توافق یا سازگاری چوپروف :

یک نمونه 20 تایی از کارگران یک کارخانه را انتخاب و آنها را بر حسب جنس (X) و به تفکیک مهارت (y) مورد مطالعه قرار داده ایم که مشاهدات عبارتند از:

(ساده ، زن)(ماهر ، مرد)(نیمه ماهر ، زن)(نیمه ماهر ، زن)

(نیمه ماهر ، مرد)(ساده،مرد)(ساده ، مرد)(ساده ، زن)(ساده ، زن)

(ماهر ، مرد)(ماهر ، مرد)(ماهر ، زن)(ساده ، زن)(ساده ، زن)

(ساده ، زن)(نیمه ماهر ، مرد)(ساده ، مرد)(ساده ، مرد)(ساده ، مرد)  
(مرد ،

$x \backslash y$	ماهر	نیمه ماهر	ساده
زن	2	2	5
مرد	4	2	5

برای اینکه بدانیم بین دو صفت  $x$  و  $y$  بستگی و یا به عبارتی دیگر توافق یا سازگاری وجود دارد یا خیر، از ضریب توافق یا سازگاری چوپروف استفاده می کنیم که فرمول آن عبارت است از:



$$r_{Tch} = \sqrt{\frac{x^2}{n + x^2}}$$

که در آن

$$x^2 = \sum_i^r \sum_j^c \frac{(f_{ij} - e_{ij})^2}{e_{ij}}$$

که در آن  $f_{ij}$  فراوانی سطر  $i$ ام و ستون  $j$ م می باشد

$$e_{ij} = \frac{f_{i0} \times f_{0j}}{n}$$

با استفاده از محاسبات جبری می توان عبارت  $x^2$  را به صورت زیر بازنویسی نمود

$$x^2 = \left( \sum_{i=1}^r \sum_{j=1}^c \frac{f_{ij}^2}{e_{ij}} \right) - n$$

مثال: جدول زیر نتایج يك بررسی در مورد رضایت  
از زندگی (x) به تفکیک جنس (y) را نشان می دهد،  
ضریب

توافق چوپروف را محاسبه کنید.

$x \backslash y$	زن	مرد	$f_{i0}$
رضایت مالی	60	80	140
رضایت روحی	70	100	170
رضایت شهرداری	70	120	190
$f_{0j}$	200	300	$n = 500$

$$e_{11} = \frac{140 \times 200}{500} = 56$$

$$e_{12} = \frac{140 \times 300}{500} = 84$$

$$e_{21} = \frac{170 \times 200}{500} = 68$$

$$e_{22} = \frac{170 \times 300}{500} = 102$$

$$e_{31} = \frac{190 \times 200}{500} = 76$$

$$e_{32} = \frac{190 \times 300}{500} = 114$$

$$x^2 = \left( \sum_{i=1}^3 \sum_{j=1}^2 \frac{f_{ij}^2}{e_{ij}} \right) - n = \left( \frac{60^2}{56} + \frac{80^2}{84} + \frac{70^2}{68} + \frac{100^2}{102} + \frac{70^2}{76} + \frac{120^2}{114} \right) - 500 = 1/36$$

$$r_{Tch} = \sqrt{\frac{x^2}{n + x^2}} = \sqrt{\frac{1/36}{500 + 1/36}} = 0/05$$

ضریب توافق پیرسون :

اندازه درجه رابطه ، بستگی یا ارتباط طبقه‌بندیها  
در يك جدول توافقي عبارت است از:

$$c = \sqrt{\frac{x^2}{n + x^2}}$$

## ضریب توافق کرامر :

ضریب توافق پیرسون را معمولاً برای جداول مربعی به کار می‌برند در صورتی که جدول توافق به صورت مستطیل باشد آن گاه از ضریب توافق کرامر که آن را با  $V$  نمایش می‌دهند استفاده می‌شود که فرمول آن عبارت است از:

$$V = \sqrt{\frac{x^2}{n \cdot \min(k-1, l-1)}}$$

که در آن  $k$  و  $l$  به ترتیب تعداد سطرها و ستونها می‌باشد.

## ضریب کروسکال:

این ضریب اعم از اینکه متغیر کمی یا کیفی باشد به کار برده می‌شود زیرا فقط از فراوانیهای جدول تبعیت می‌کند و به مقادیر متغیر ارتباطی ندارد و فرمول آن عبارت است از:

$$\tau_y = \frac{E_2}{E_1}$$

که در آن

$$E_1 = \sum_{j=1}^l \left( \frac{N - n_{0j}}{N} \right) n_{0j} \quad , \quad E_2 = \sum_{i=1}^k \sum_{j=1}^l \left( \frac{n_{i0} - j_{ij}}{n_{i0}} \right) n_{ij}$$

لازم به یادآوری است که  $n_{i0}$  فراوانی توزیع حاشیه ای  $x$ ها و  $n_{0j}$  فراوانی توزیع حاشیه ای  $y$ ها می‌باشد.



## ضریب گاتمن :

این ضریب نیز از فراوانی جدول توافقی تبعیت می کند و اغلب در مواردی به کار برده می شود که نتوان صفات مورد مطالعه را به صورت رتبه ای طبقه بندی کرد و فرمول آن عبارت است از:

$$\lambda = \frac{\sum m_y - M_y}{N - M_y}$$

که در آن  $M_y$  عبارت است از  $\max n_j$  یعنی بزرگترین فراوانی توزیع حاشیه ای  $y$  ها و  $m_y$  نیز عبارت است از بزرگترین فراوانی توزیع های شرطی  $y$  ها و  $N$  حجم کل می باشد.

واضح است که همواره  $0 \leq \lambda \leq 1$  هرچقدر به يك نزدیکتر باشد میزان وابستگی صفات مورد مطالعه بیشتر است.

مثال: جدول زیر نتیجه آمارگیری نیروی انسانی  
در یک شهر بزرگ را نشان می دهد مطلوب  
است ضریب کروسکال و ضریب گاتمن و تفسیر

آنها  
علل بیرونی

$n_0$	عدم علاقه	نبودن مدرسه	مخالفت والدین	فقر مالی	جنس
88	12	11	8	57	مرد
92	7	10	48	27	زن
$n=180$	19	21	56	84	$n_j$

**حل:**  $E_1 = \sum_{j=1}^4 \left( \frac{n - n_{.j}}{n} \right) n_{.j} = 118/5$

و به همین شیوه

$$E_2 = \sum_{i=1}^2 \sum_{j=1}^4 \left( \frac{n_{i0} - n_{ij}}{n_{i0}} \right) n_{ij} = 104/7$$

$$\tau_y = \frac{E_1 - E_2}{E_1} = \frac{118/5 - 104/7}{118/5} = 0/12$$

**0/12** میزان شدت وابستگی بین علل بیسوادی و جنس را نشان می

دهد، حال ضریب گاتمن را محاسبه می کنیم:

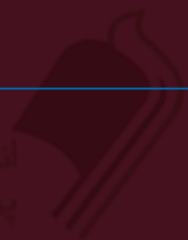
$$\lambda = \frac{\sum m_y - M_y}{N - M_y} = \frac{(57+48) - 84}{180 - 84} = 0/22$$

عدد 0/22 میزان شدت وابستگی علل بیسوادی و جنس را نشان می دهد

تا حدودی می توان گفت ، که علل بیسوادی به جنس مربوط است.

# فصل نهم

## رگرسیون



# رگرسیون

از نظر تاریخی کلمه رگرسیون با زمینه فنی کنونی اش اثر برای اولین بار توسط فرانسیس گالتن ، کسی که قد پسران و متوسط قد والدینشان را تجزیه و تحلیل کرد بکار برده شد.

گالتن از روی مشاهداتش نتیجه گرفت که قد پسران متعلق به والدین خیلی بلند ( یا کوتاه ) به طور کلی بلندتر ( کوتاهتر ) از قد متوسط بوده است ولی البته نه بلندی یا کوتاهی والدینشان. این نتیجه گیری در سال 1855 تحت عنوان « رگرسیون به طرف میانگین » منتشر شد.

در این متن اصطلاح رگرسیون اشاره بر این داشت که قد پسران به جای میل به مقادیر فریخته تر به سوی مقدار متوسط میل و بازگشت دارد.

يك مسئله ساده رگرسيون :

داده‌هاي جدول زير مقدار آب داده شده بر حسب سانتي متر و مقدار محصول يونجه بر حسب تن را نشان مي‌دهد .  
مسئله پيش بيني محصول يونجه به عنوان تابعي از مقدار آبياري است.

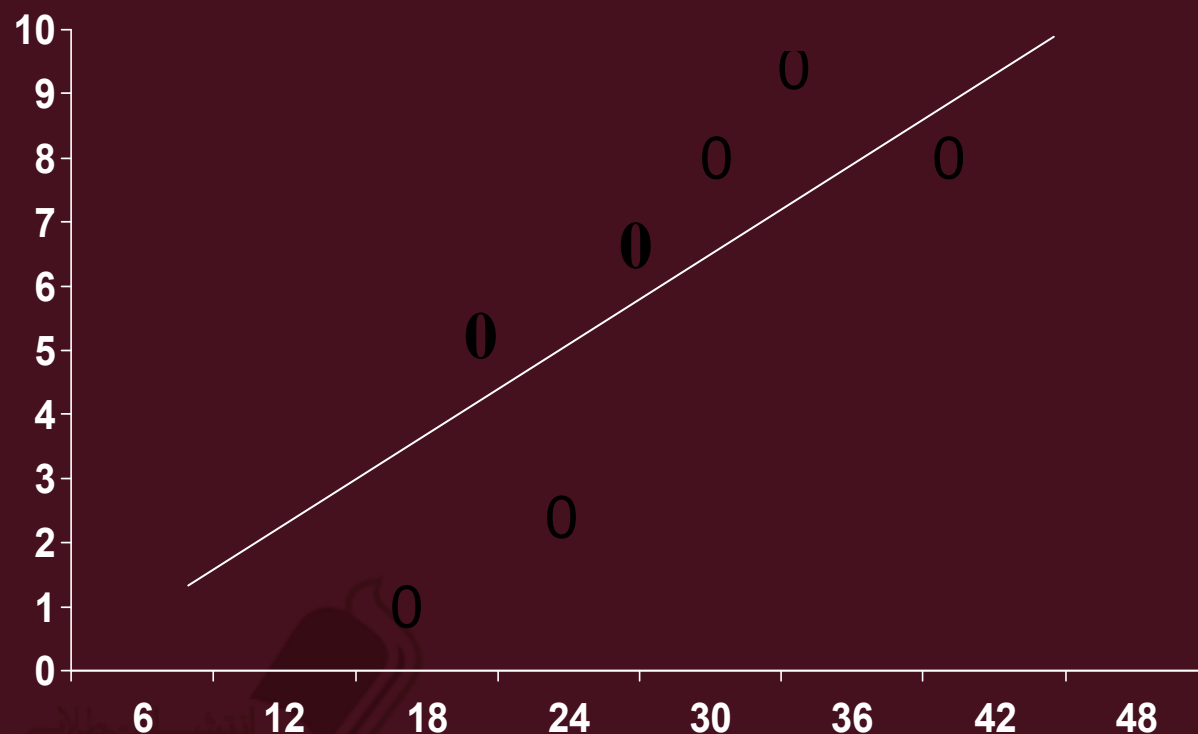
	12	18	24	30	36	42	48
میزان محصول ( )	5/27	5/68	6/25	7/21	8/02	8/71	8/42

در مطالعه‌ي رابطه بين دو متغير، اولين قدم منطقي رسم داده‌ها به صورت نقاطي بر روي صفحه‌ي مختصات دو بعدي دکارتی می‌باشد که شکل حاصل نمودار پراکنش نامیده می‌شود.

که چگونگی پراکنش نقاط در اطراف يك خط راست یا يك نوع منحنی مشخص را نشان می‌دهد.



و نیز برداشتن عینی از میزان پراکندگی داده‌ها پیرامون خط خط یا منحنی را فراهم می‌کند .  
 بنابراین اطلاع منعکس در نمودار پراکنش برای جستجوی یک مدل ریاضی مناسب ( خطی یا منحنی ) مفید است .





نمودار پراکنش فوق بازگو کننده‌ی این موضوع است که نقاط رسم شده روی نمودار پراکنش داده‌ها تقریباً در امتداد يك خط راست پراکنده شده‌اند .

بنابر این می‌توان به منظور پیشگویی مقدار  $y$  از روی مقادیر  $x$  ، خطی راست به این مجموعه نقاط برازاند .  
چنین خطی در شکل برازش داده شده است ، که مقدار پیشگویی شده‌ی  $y$  بعنوان فاصله عمودی نقطه‌ی  $x$  تا خط انتخاب می‌شود .



اصل حداقل مربعات :

ملاحظه نمودیم که رسم نمودار پراکنش هرگاه نقاط در امتداد  
یک خط راست پراکنده شده باشند می‌توان به این نقاط خط  
راستی به معادله  $y=a+bx$  برازش داد به ازای

$a$  و  $b$  های متفاوت می‌توان بی‌نهایت خط به این نقاط برازش  
داد اما درصدد هستیم خط راستی با شیب  $b$  و عرض از  
مبدأ  $a$  را با این ویژگی بیابیم



که به تمام نقاط رسم شده در نمودار پراکنش از هر خط دیگری نزدیکتر باشد به عبارتی مجموع کل مربعات خطاهای پیش بینی کمترین مقدار را داشته باشد که به آن اصل حداقل مربعات انحراف می‌گویند

$$\left( \sum_{i=1}^n \xi_i^2 \right)$$

شیب خط و عرض از مبدأ چنین خطی که به ازای آن  $\sum_{i=1}^n \xi_i^2$  می‌نی‌م شود از رابطه زیر محاسبه می‌شود.

مجموع مربعات خطاهای پیش بینی

$$= R = \sum_{i=1}^n \xi_i^2 = \sum (y_i - (a + bx))^2$$

با گرفتن مشتق جزئی نسبت به  $a$  و  $b$  و حل دستگاه زیر

شیب خط

عرض از مبدأ

$$\frac{\delta R}{\delta a} = -2 \sum_{i=1}^n (y_i - (a + bx_i)) = 0$$

$$\frac{\delta R}{\delta b} = -2 \sum_{i=1}^n x_i (y_i - (a + bx_i)) = 0$$

$$\sum_{i=1}^n y_i = na + b \sum_{i=1}^n x_i$$

$$\sum_{i=1}^n x_i y_i = a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2$$

$$b = \frac{S_{xy}}{S_{xx}}$$

$$a = \bar{y} - b\bar{x}$$

داریم

# که در آن

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}$$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n \bar{x}^2$$

$$S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n y_i^2 - n \bar{y}^2$$

مثال : معادله‌ی خط رگرسیون داده‌های مربوط به مقدار آب (  $x$  ) و میزان محصول (  $y$  ) یونجه را تعیین کنید.

$x_i$	$y_i$	$x_i^2$	$x_i y_i$
12	5 / 27	144	63 / 24
18	5 / 68	324	102 / 24
24	6 / 25	572	150
30	7 / 21	900	216 / 3
36	8 / 02	1296	288 / 72
42	8 / 71	1764	364 / 82
48	8 / 42	2304	404 / 16
210	49 / 56	7308	1590 / 48

مجموع

$$s_{xy} = \sum x_i y_i - n \bar{x} \bar{y} = 1590 / 48 - 7 \times (30) \times (7 / 08) = 103 / 68$$

$$s_{xx} = \sum x_i^2 - n \bar{x}^2 = 7308 - 7 \times (30)^2 = 1008$$

$$b = \frac{s_{xy}}{s_{xx}} = \frac{103/68}{1008} = 0/103$$

$$a = \bar{y} - b \bar{x} = 7 / 08 - 0/103 \times 30 = 3 / 99$$



بنابر این معادله‌ی خط رگرسیون عبارت است از :

$$y = 3/99 + 0/103 x$$

لذا با استفاده از معادله رگرسیون می‌توان به ازاء هر مقدار دلخواه  $x$  اندازه متغیر وابسته  $y$  را پیش بینی کرد.

با استفاده از روابط ریاضی می‌توان با استفاده از ضریب همبستگی ، شیب خط رگرسیون را محاسبه نمود .

$$b = \frac{s_{xy}}{\sqrt{s_{xx}} \sqrt{s_{xx}}}$$

$$b = \frac{s_{xy}}{\sqrt{s_{xx}} \sqrt{s_{yy}}} \cdot \frac{\sqrt{s_{yy}}}{\sqrt{s_{xx}}}$$

$$b = r_{xy} \times \frac{s_y}{s_x}$$



مثال : فرض کنید بخواهیم نمره‌های دستور زبان انگلیسی را از طریق نمره‌های درک و فهم پیش بینی کنیم و همبستگی بین این دو متغیر 0/5 و انحراف استاندارد های نمرات دستور زبان و درک فهم به ترتیب 3 و 2 باشد شیب خط رگرسیون را بیابید.  
حل :

$$b_{xy} = r_{x.y} \cdot \frac{s_y}{s_x}$$

$$b_{xy} = 0/5 \times \frac{3}{2} = 0/75$$

