

آمار توصیفی

ارائه دهنده: امیرحسین سبحانی

محاسبه چندك براي داده‌هاي گسسته

فرض كنيد y_1, y_2, \dots داده‌هاي ما باشند و شكل مرتب شده آنها را با $y_{(1)}, y_{(2)}, \dots$ نمايش دهيم. براي محاسبه چندك

$(n+1)p$

صحيح باشد $k = (n+1)p, \quad Q_p = y_{(k)}$

صحيح نباشد $k = [(n+1)p], \quad \alpha = (n+1)p - k \longrightarrow Q_p = (1-\alpha)y_{(k)} + \alpha y_{(k+1)}$

محاسبه چندك براي داده‌هاي پيوسته

با توجه به ستون فراواني تجمعي در جدول فراواني، ردیفی را که چندك (مثال $p=0.25$) در آن قرار دارد را مشخص مي‌کنیم.

رده	x_i	f_i	g_i
1/35-1/55	1/45	4	4
1/55-1/75	1/65	6	10
1/75-1/95	1/85	12	22
1/95-2/15	2/05	9	31
2/15-2/35	2/25	8	39
2/35-2/55	2/45	6	45
2/55-2/75	2/65	2	47
2/75-2/95	2/85	3	50
جمع		50	—

$Q_{0.25}$

$$Q_p = L_{Q_p} + \left(\frac{np - g_{np}}{f_{Q_p}} \right) I$$

$$(n) \times (p)$$

$$50 \times 0.25 = 12.5$$

L_{Q_p} : کران پایین رده چندك

g_{np} : فراواني تجمعی رده بلافاصله قبل از رده چندك

f_{Q_p} : فراواني رده چندك

I : طول رده

$$Q_{0.25} = 1.75 + \left(\frac{50 \times 0.25 - 10}{12} \right) \times 0.2 = 1.79$$

میانگین قدر مطلق انحراف ها

- میانگین قدر مطلق انحرافات: میانگین قدر مطلق انحراف که به آن انحراف متوسط می‌گویند به صورت زیر تعریف می‌شود.

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

- تذکر: گاهی ممکن است از میانه به جای میانگین برای محاسبه متوسط انحراف استفاده شود.

میانگین قدر مطلق انحراف ها

مثال: با توجه به جدول فراوانی زیر میانگین قدر مطلق انحراف ها را محاسبه نمایید:

x_i	f_i	$f_i x_i$	$f_i x_i - \bar{x} $
18.8	2	37.6	62.4
21.9	7	153.3	196.7
25	10	250	250
28.1	17	477.7	372.3
31.2	3	93.6	56.4
34.3	6	205.8	94.2
37.4	5	187	63
Σ	50	1405	1095

$$\bar{x} = \frac{\sum f_i x_i}{n} = \frac{1405}{50} = 28.1$$

$$AD = \frac{\sum_{i=1}^k f_i |x_i - \bar{x}|}{n} = \frac{1095}{50} = 21.9$$

میانگین توان دوم انحراف ها (واریانس)

واریانس داده های یک جامعه به صورت زیر تعریف می گردد:

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

و واریانس نمونه ای به صورت زیر تعریف می گردد

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

همچنین انحراف معیار جامعه به صورت جذر مثبت واریانس و انحراف معیار نمونه ای به صورت جذر مثبت واریانس نمونه ای تعریف می شود:

$$\sigma = \sqrt{\sigma^2} , S = \sqrt{S^2}$$

میانگین توان دوم انحراف ها (واریانس)

• قضیه: واریانس را می توان به کمک فرمول زیر محاسبه نمود:

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

x_i	f_i	$f_i x_i$	$f_i x_i^2$
18.8	2	37.6	706.88
21.9	7	153.3	3357.27
25	10	250	6250
28.1	17	477.7	13423.37
31.2	3	93.6	2920.32
34.3	6	205.8	7058.94
37.4	5	187	6993.8
Σ	50	1405	40710.58

$$\bar{x} = \frac{\sum f_i x_i}{n} = \frac{1405}{50} = 28.1$$

$$S^2 = \frac{1}{n-1} \left(\sum_{i=1}^k f_i x_i^2 - n \bar{x}^2 \right) = \frac{1}{49} (40710.58 - 50(28.1)^2) = 25.103$$

میانگین توان دوم انحراف ها (واریانس)

- اگر به کمک تبدیل $y_i = ax_i + b$ داده‌های x_1, x_2, \dots, x_n را به y_1, y_2, \dots, y_n تبدیل کنیم رابطه زیر بین واریانس x_i ها و y_i برقرار است:

$$S_y^2 = a^2 S_x^2, S_y = a S_x \quad a > 0$$

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

روش تبدیل داده ها جهت محاسبه واریانس

روش تبدیل داده ها جهت محاسبه واریانس: برای ساده تر کردن محاسبات جهت محاسبه میانگین و واریانس خصوصاً وقتی داده ها بزرگ باشند می توان از روش تبدیل داده ها استفاده نمود. در این روش به کمک تبدیل

$$y_i = \frac{x_i - a}{b} \quad i = 1, 2, \dots, n$$

داده های x_1, x_2, \dots, x_n را به y_1, y_2, \dots, y_n تبدیل می کنیم. در اینجا a و b به ترتیب برای تغییر مبدا و تغییر واحد اندازه گیری استفاده میشود. معمولاً برای داده های پیوسته a را نماینده رده نمایی و b را طول رده یعنی I می گیرند. سپس میانگین و واریانس برای داده های جدید y_1, y_2, \dots, y_n محاسبه می شود و به کمک فرمول های زیر میانگین و واریانس داده های قبلی محاسبه می گردد:

$$S_x^2 = b^2 S_y^2, \quad S_x = b S_y$$

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

روش تبدیل داده ها جهت محاسبه واریانس

رده	x_i	f_i	y_i	$f_i y_i$	$f_i y_i^2$
17.25-20.35	18.8	2	-3	-6	18
20.35-23.45	21.9	7	-2	-14	28
23.45-26.55	25	10	-1	-10	10
26.55-29.65	28.1	17	0	0	0
29.65-32.75	31.2	3	1	3	3
32.75-35.85	34.3	6	2	12	24
35.85-38.95	37.4	5	3	15	45
Σ		50		0	128

$$y_i = \frac{x_i - 28.1}{3.1} \quad i=1,2,\dots,7$$

$$\bar{y} = \frac{0}{50} = 0$$

$$S_y^2 = \frac{1}{49} (128 - 0) = 2.612$$

$$\bar{x} = 28.1 + 3.1(0) = 28.1$$

$$S_x^2 = (3.1)^2 (2.612) = 25.103$$

$$S_x = \sqrt{S_x^2} = 5.01$$

ضریب تغییرات

ضریب تغییر: واریانس و انحراف استاندارد به واحد اندازه گیری بستگی دارد. برای مقایسه دو سری داده باید از شاخص هایی استفاده کنیم که به واحد اندازه گیری بستگی نداشته باشد یکی از این شاخص ها ضریب تغییر است که به کمک فرمول زیر تعریف می شود:

$$CV = \frac{S}{\bar{x}}, CV = \frac{\sigma}{\mu}$$

مثال: کارخانه ای دو نوع لامپ تولید می کند، لامپ نوع اول دارای میانگین طول عمر 200 ساعت و انحراف معیار 11 ساعت است و لامپ دوم دارای میانگین کارکرد 240 و انحراف معیار 12 ساعت است. کدام نوع لامپ بهتر است؟

جواب: لامپ دوم مناسب تر است. چون میانگین بیشتر و ضریب تغییر کمتری دارد.

$$CV_1 = \frac{S_1}{\bar{x}_1} = \frac{11}{200} = 0.055$$

$$CV_2 = \frac{S_2}{\bar{x}_2} = \frac{12}{240} = 0.05$$

Z-نمره (Z-SCORE)

Z-نمره (مقدار استاندارد): تعداد انحراف استاندارد هایی است که یک مقدار معین بالاتر و یا پایین تر از میانگین است.

Z-نمره به کمک فرمول زیر قابل محاسبه است:

$$z = \frac{x - \bar{x}}{S} , \quad z = \frac{x - \mu}{\sigma}$$

✓ Z-نمره فاقد مقیاس اندازه گیری است.

✓ یک داده، داده پرت است اگر Z-نمره آن بیش از 3 یا کمتر از -3 شود.

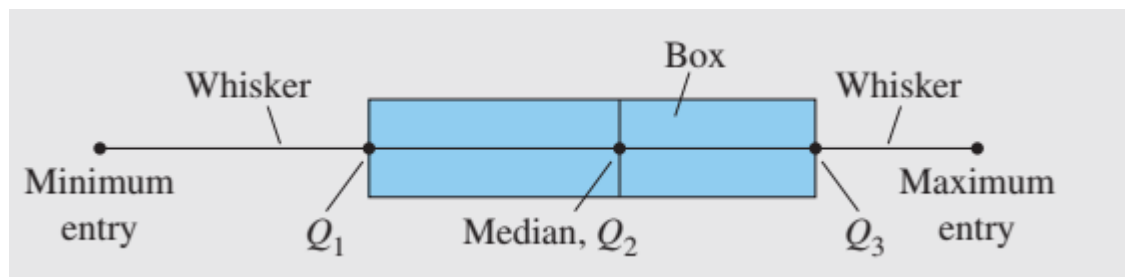
✓ اگر Z-نمره یک داده بیش از 2 یا کمتر از -2 شود غیر عادی است.

مثال: داده های زیر مربوط به تعداد ضربان های قلب یک نمونه است. مقدار پرت را به کمک Z-نمره بیابید.

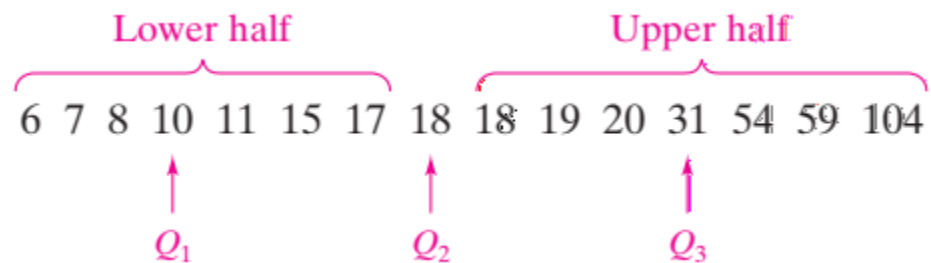
AGE	GENDER (1=M)	PULSE	Z-score		
43	0	80	0.475079	Mean	74.04082
38	0	94	1.59119	Std	12.54356
69	0	58	-1.27881		
44	0	66	-0.64103		
72	0	56	-1.43825		
37	0	82	0.634524		
29	0	78	0.315635		
44	0	86	0.953413		
51	0	88	1.112857		
66	0	56	-1.43825		
46	0	36	-3.0327		
31	0	66	-0.64103		
35	0	84	0.793968		
55	0	76	0.15619		
22	0	78	0.315635		
51	0	64	-0.80048		
26	0	66	-0.64103		
62	0	78	0.315635		
43	0	60	-1.11936		
56	0	64	-0.80048		
54	0	84	0.793968		
20	0	82	0.634524		
63	0	70	-0.32214		

نمودار جعبه ای

از نمودار جعبه ای برای نمایش پراکندگی داده ها استفاده می شود. برای رسم این نمودار از پنج شاخص استفاده می کنیم: چارک اول، چارک سوم، میانه، مینیموم و ماکسیموم داده ها. برای پیدا کردن مینیموم و ماکسیموم، داده های پرت را کنار میگذاریم. مقدار ماکسیموم را بزرگترین داده ای که بیشتر از $Q_3 + 1.5IQR$ نباشد و مقدار مینیمم را کوچکترین داده ای که کمتر از $Q_1 - 1.5IQR$ نباشد در نظر می گیریم. برای رسم این نمودار، مستطیلی بین چارک اول و سوم رسم کرده و میانه را با خطی در جعبه مشخص می کنیم. همچنین مینیموم و ماکسیموم را با پاره خطی نمایش می دهیم و این پاره خط ها را به مستطیل مطابق شکل زیر وصل می کنیم. همچنین داده های پرت را نیز با دایره یا ستاره مشخص می کنیم.



7 18 11 6 59 17 18 54 104 20 31 8 10 15 19

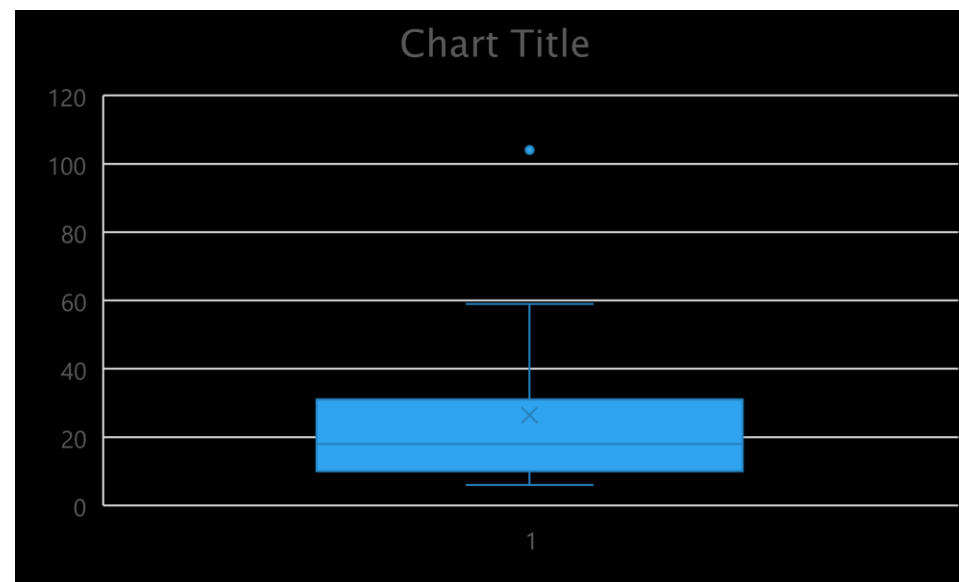


$$IQR = Q_3 - Q_1 = 31 - 10 = 21$$

$$QD = \frac{1}{2} IQR = 10.5$$

$$Q_1 - 1.5 IQR = -21.5$$

$$Q_3 + 1.5 IQR = 62.5$$



died	Survived
1.175	1.13
1.23	1.575
1.31	1.68
1.5	1.76
1.6	1.93
1.72	2.015
1.75	2.09
1.77	2.6
2.275	2.7
2.5	2.95
1.03	3.16
1.1	3.4
1.185	3.64
1.23	2.83
1.26	1.41
1.30	1.715
1.30	1.72
1.55	2.04
1.82	2.2
1.89	2.4
1.94	2.55
2.20	2.57
2.27	3.005
2.44	
2.56	
2.73	
1.05	

died	Survived
1.03	1.13
1.05	1.41
1.1	1.575
1.175	1.68
1.185	1.715
1.23	1.72
1.23	1.76
1.26	1.93
1.30	2.015
1.30	2.04
1.31	2.09
1.5	2.2
1.55	2.4
1.6	2.55
1.72	2.57
1.75	2.6
1.77	2.7
1.82	2.83
1.89	2.95
1.94	3.005
2.20	3.16
2.27	3.4
2.275	3.64
2.44	
2.5	
2.56	
2.73	

داده های روبرو مربوط به وزن هنگام تولد نوزادان دچار مشکل ریه می باشد. این داده ها به دو گروه متفاوت بر اساس نوزادان نجات یافته و نوزادان از دست رفته تقسیم شده اند. به کمک نمودار جعبه ای تحلیل کنید که آیا وزن نوزادان رابطه ای با مرگ آنها داشته است؟



DIED

● $n = 27, p = 0.25, (n + 1)p = 7$

$$Q_1 = x_7 = 1.23$$


$$n = 27, p = 0.5, (n + 1)p = 14$$

$$Q_2 = x_{14} = 1.6$$

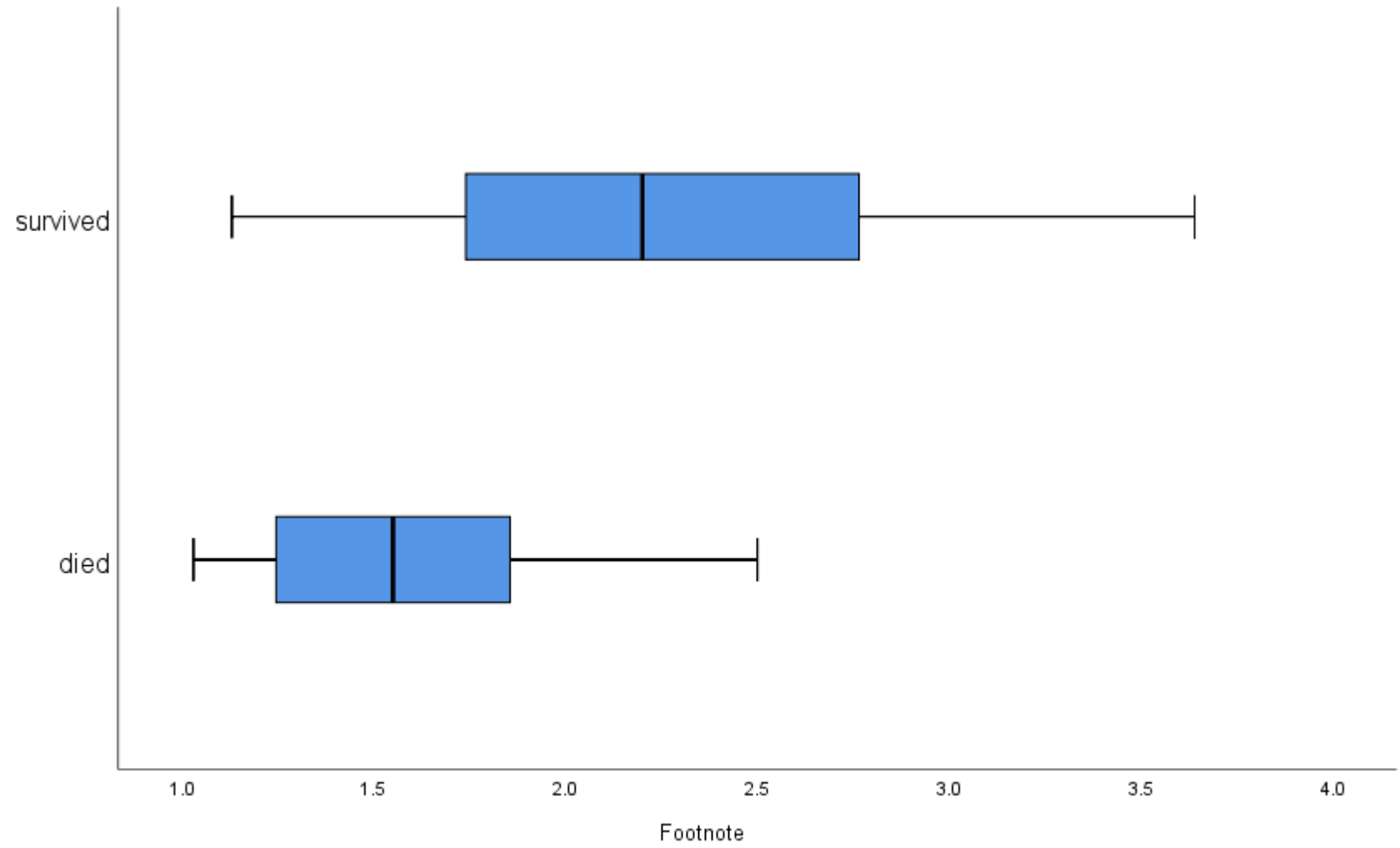
$$n = 27, p = 0.75, (n + 1)p = 21$$

$$Q_3 = x_{21} = 2.2$$

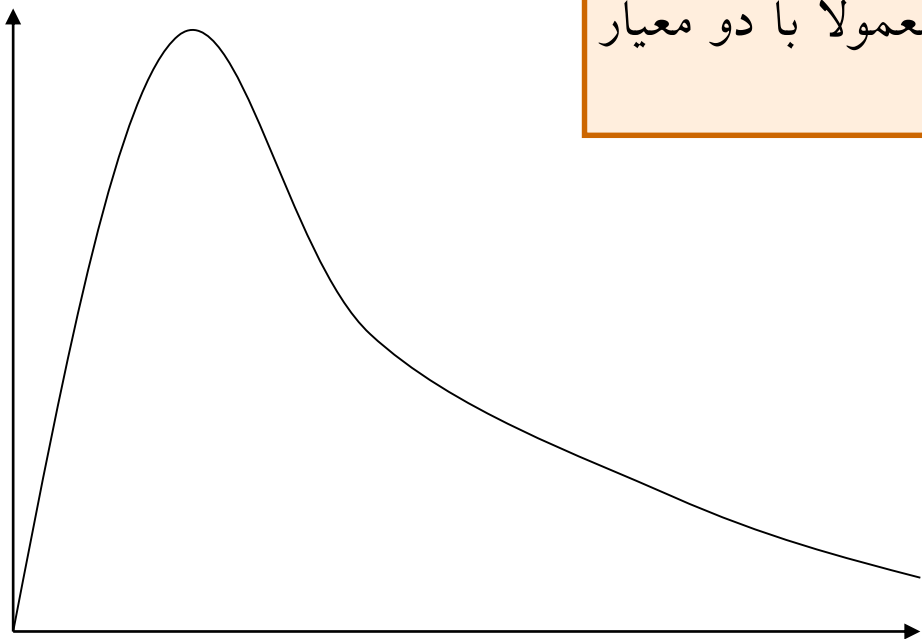
$$IQR = Q_3 - Q_1 = 2.2 - 1.23 = 0.97$$

$$QD = \frac{1}{2} IQR = 0.485$$


died	Survived
1.175	1.13
1.23	1.575
1.31	1.68
1.5	1.76
1.6	1.93
1.72	2.015
1.75	2.09
1.77	2.6
2.275	2.7
2.5	2.95
1.03	3.16
1.1	3.4
1.185	3.64
1.23	2.83
1.26	1.41
1.30	1.715
1.30	1.72
1.55	2.04
1.82	2.2
1.89	2.4
1.94	2.55
2.20	2.57
2.27	3.005
2.44	
2.56	
2.73	
1.05	

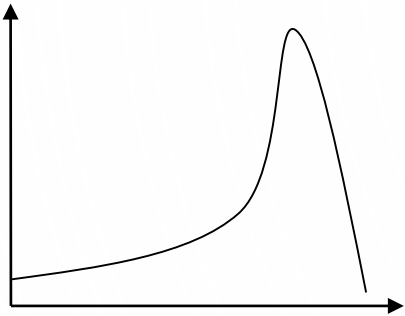


منحنیهای فراوانی در طبیعت تنوع زیادی دارند، اما بسیاری از منحنیهای فراوانی یا متقارن هستند یا چوله و یا برجسته و یا پخ. ایده آلترین منحنی فراوانی متقارن، منحنی فراوانی نرمال استاندارد است. در طبیعت، عموماً منحنی فراوانی متقارن ایده آل کمتر یافت می‌شود میزان انحراف از تقارن ایده آل را معمولاً با دو معیار چولگی و برجستگی می‌سنجند.

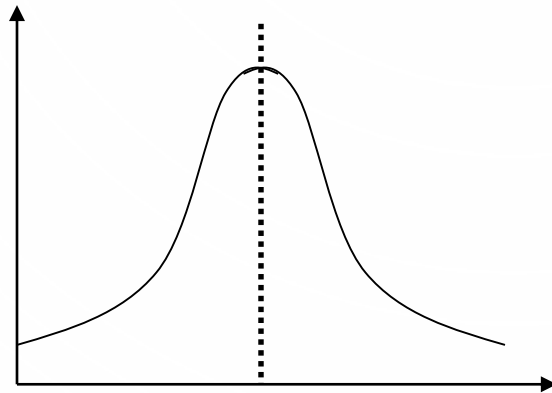


منحنیهای فراوانی:

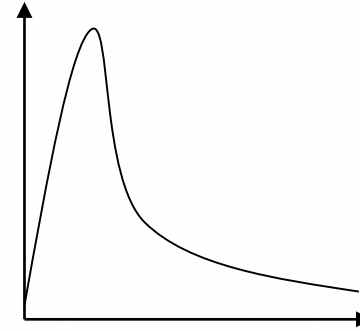
نامتقارن / چوله به چپ
Skewed to left



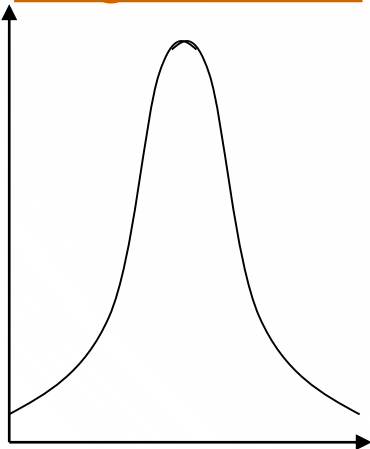
منحنی نرمال استاندارد
Standard normal



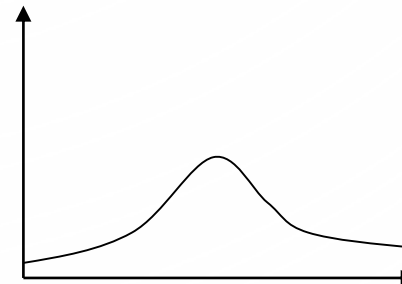
نامتقارن / چوله به راست
Skewed to right



برجسته
High kurtosis

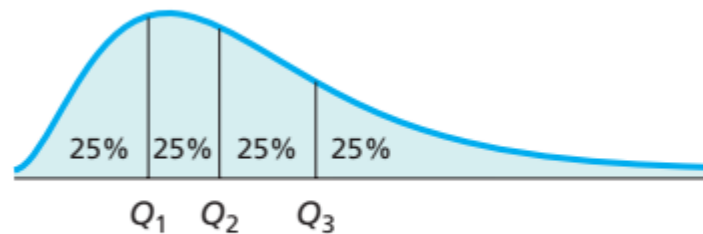


پخ
low kurtosis

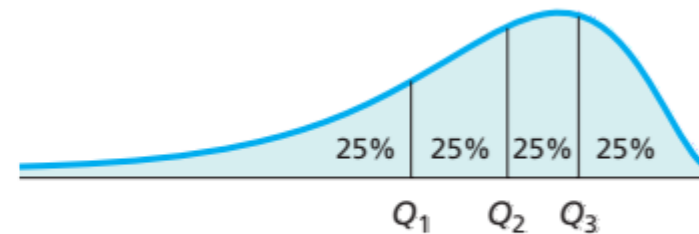


چولگی

چولگی معیاری از سنجش تقارن و یا عدم تقارن توزیع داده ها در منحنی فراوانی به حساب می آید. توزیع داده ها را چوله گوئیم اگر متقارن نباشد و در یک سمت بیش از سمت دیگر گسترش پیدا کرده باشد.



(c) Right skewed



(d) Left skewed

چولگی منفی (چولگی به چپ): اگر داده‌های با فراوانی بالا در سمت راست منحنی فراوانی باشند و دم منحنی به سمت چپ کشیده شده باشد منحنی را چوله به چپ میگویند

چولگی مثبت (چولگی راست): اگر داده‌های با فراوانی بالا در سمت چپ منحنی فراوانی قرار داشته باشند و دم منحنی فراوانی به سمت راست کشیده شده باشد و گسترش پیدا کرده باشد در این صورت منحنی را چوله به راست میگویند

تذکر: در توزیع متقارن میانه نما و میانگین بر هم منطبق است

چولگی

تذکر: توزیع های چوله به راست از توزیع های چوله به چپ بیشتر ظاهر می شوند زیرا معمولا داده های غیر عادی بزرگ هستند.

باتوجه به اینکه میانگین به سادگی تحت تاثیر داده های پرت قرار می گیرد، در توزیع های چوله به چپ مد از میانگین بزرگتر است و در توزیع های چوله به راست، مد از میانگین کوچکتر است.

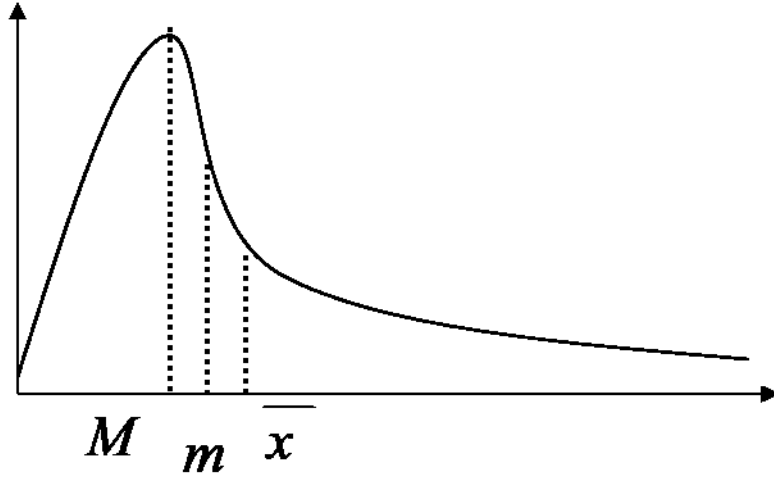
در حالت کلی در بیشتر مواقع اگر توزیع چوله به راست باشد:

$$M < m < \bar{x}$$

و اگر توزیع چوله به چپ باشد معمولا داریم:

$$\bar{x} < m < M$$

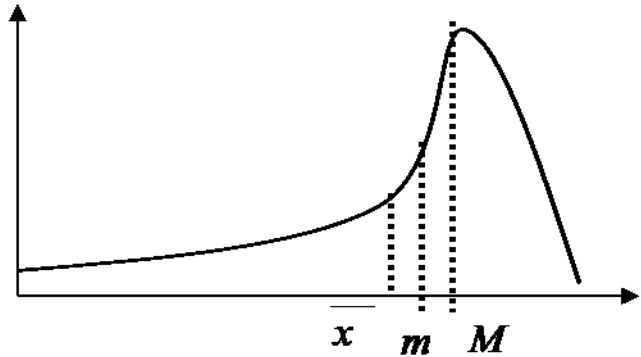
چولگی



$$M < m < \bar{x}$$

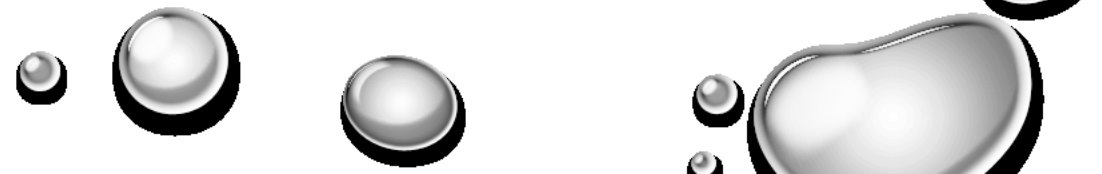
به راست یا مثبت

چولگی:



$$M > m > \bar{x}$$

به چپ یا منفی



چولگی

معیار چولگی کارل پیرسون: ضریب چولگی کارل پیرسون بر اساس این واقعیت تعریف شده است که در توزیع چوله به راست مد کمتر از میانگین و در توزیع چوله به چپ مد بیشتر از میانگین است:

$$SK_1 = \frac{\bar{x} - M}{S}$$

وقتی از مقدار مد اطلاعی نداشته باشیم، از فرمول دوم پیرسون استفاده می‌کنیم:

$$SK_2 = \frac{3(\bar{x} - m)}{S}$$

اگر ضریب پیرسون منفی باشد توزیع چولگی منفی (به چپ) و اگر مثبت باشد چولگی به راست دارد.

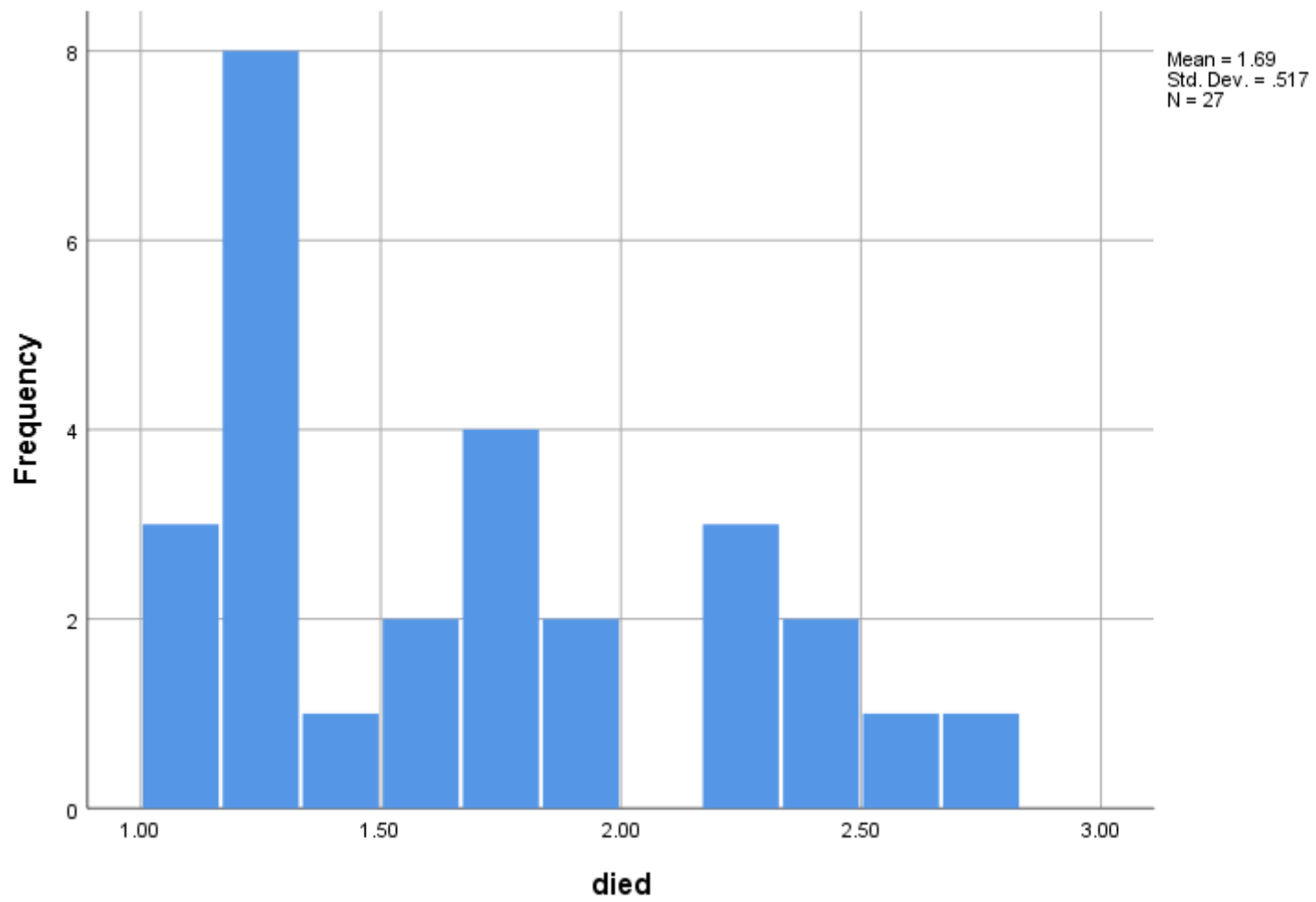
چولگی

ضریب چولگی گشتاوری:

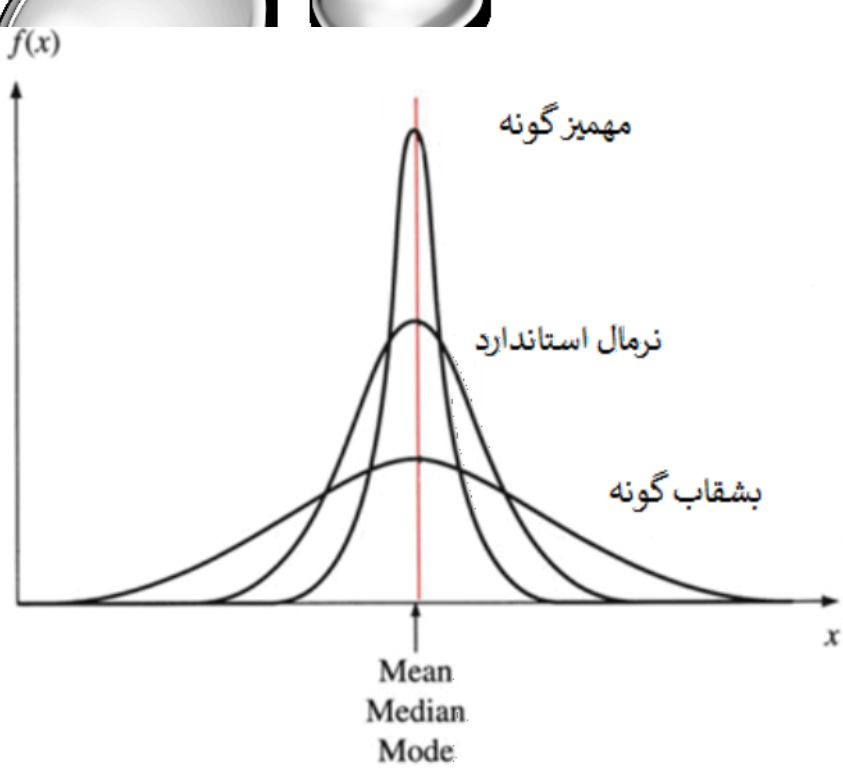
$$\alpha_3 = \frac{\mu_3}{\sigma^3}, \mu_3 = \frac{\sum f_i (x_i - \mu)^3}{n}$$

ضریب چولگی چارکی:

$$SQ = \frac{Q_3 - 2Q_2 + Q_1}{Q_3 - Q_1}$$



کشیدگی



کشیدگی یک توزیع معیار سنجش برآمدگی توزیع نسبت به توزیع نرمال استاندا،
توزیع ها از نظر کشیدگی سه نوع هستند:

توزیع نرمال استاندارد: که کشیدگی آن را معیار قرار داده و توزیع های دیگر را با آن می سنجیم.
توزیع با کشیدگی بالا (مهمیز گونه): که کشیدگی آن از کشیدگی توزیع نرمال استاندارد بیشتر است.
توزیع با کشیدگی منفی (بشقاب گونه): که کشیدگی آن از توزیع نرمال استاندارد کمتر است.

ضرایب کشیدگی

ضریب کشیدگی گشتاوری:

$$\alpha_4 = \frac{\mu_4}{\sigma^4} - 3$$

ضریب کشیدگی چندکی:

$$E = \frac{QD}{Q_{0.9} - Q_{0.1}} - 0.263$$