

A NOTE ON ADJUSTMENTS FOR FIRST AND SECOND MOMENTS IN A GROUPED FREQUENCY DISTRIBUTION SPLIT UP INTO SUB-SECTIONS

By J. M. SEN GUPTA
Statistical Laboratory, Calcutta

Object of the Study : Sheppard's adjustments for grouping in the moments about the mean for frequency curves with high kurtosis at both the ends are well known. These moments refer to the whole of the frequency distribution. In this note the whole range of the variate has been divided up into several sections, and studies have been made as to how grouping affects the first two moments for each section separately. Incidentally, Sheppard's Correction for the second moment has been verified. The normal frequency distribution with mean zero and unit standard deviation has been taken up for the study and the Normal Probability Table published by Sheppard in Biometrical Tables Part I has been used which gives the accumulated probability for the normal distribution at intervals of .01 of the deviate.

Experimental Procedure : The normal frequency distribution for $N = 100,000$ was constructed with the help of Sheppard's Tables, and one half of the curve was used in the computing process. The total range for half the curve, which really extends to positive infinity in terms of the deviates, was assumed to terminate for practical purposes at $+4\sigma$, a point after which the contribution to the total frequency is less than 0.5. The accumulated frequencies read off from Sheppard's Table at intervals of .01 of the deviate were often used. The frequencies that occurred were added up in groups of two or three and so on, and the deviate sub-totals were obtained for the different frequency distributions with $-45, -10, -15, -20, -30, -35, -40, -50, -60, -65$, and -1.20 as class intervals. The sums and sums of squares $S(x)$ and $S(x^2)$ with the above class intervals were then calculated directly. The deviate x being measured from the origin, these raw sums and sums of squares required no further correction with reference to the general mean. The sums and sums of squares were however recalculated in four sub-totals, corresponding to the four sections into which the entire half-range was arbitrarily split up. Each of these sections represented the portion of the curve intercepted between ordinates at intervals of 1.20 of the deviate. The variation in error, introduced by grouping at different regions in the curve could thus be observed separately. Table I gives the sums and sums of squares for one half of the normal curve in four sub-sections for each of the twelve class intervals worked with.

TABLE I. ABSOLUTE SUMS AND SUMS OF SQUARES FOR ONE HALF OF THE NORMAL CURVE.

Sections A	$I(-60-1.20)$ $n=3493$		$II(1.21-2.40)$ $n=10687$		$III(2.41-3.60)$ $n=804$		$IV(3.61-4.80)$ $n=16$		Total $(-60-4.80)$ $n=50,000$	
	$S(x)$	$S(x^2)$	$S(x)$	$S(x^2)$	$S(x)$	$S(x^2)$	$S(x)$	$S(x^2)$	$S(x)$	$S(x^2)$
.05	20490	15149	17182	29243	2178	5901	81	237	32980	30011
.10	20403	15184	17103	29352	2180	5969	82	237	32924	30015
.15	—	—	—	—	—	—	—	—	32959	30053
.20	20544	15184	17236	29760	2196	6097	82	234	34024	30187
.25	—	—	—	—	—	—	—	—	40102	30259
.30	20630	15128	17307	29864	2195	6042	82	230	40191	30373
.35	—	—	—	—	—	—	—	—	40205	30408
.40	20731	15079	17408	29239	2204	6108	82	244	40429	30670
.50	—	—	—	—	—	—	—	—	40729	31039
.60	21009	14925	17004	30005	2242	6200	83	248	41098	31495
.80	—	—	—	—	—	—	—	—	42046	32068
1.20	20806	13847	16237	31426	2142	7256	67	262	44812	36002

Corrections for μ_2 : Table 2 shows the observed μ_2 estimates separately for each of the sections and for the entire curve. The sectional μ_2 's however, refer to the general mean and not to the respective sectional means. This has been given for each of the class intervals that have been used.

The true value of μ_2 with an infinitesimal class interval Δx , is obviously unity for the entire curve. The true sectional values of μ_2 were available in the Incomplete Normal Moments Table in Biometrical Tables Part I. The deviations D of the estimated μ_2 values for each of the class intervals by sections and for the total is also being shown in Table 2.

The second differences of $D(\mu_2)$ was found to be constant, which at once suggested the form of the fit, as $D(\mu_2) = c\Delta x$. Least square fittings were also tried. The sectional fits were bad, except in the form $D = c\Delta x$ where c was practically -0.833 for the entire half curve. The graduated deviations accordingly have been shown in Table 2. It will be seen that the total fit is extremely unsatisfactory, demonstrating that Sheppard's correction term $4\sigma/12$ for μ_2 is inappropriate and exact for all practical purposes. Even for the highest class interval of 1.20 , which gives us only 4 classes for half the curve, the agreement is quite good, and we conclude that, with class intervals upto the order of the standard deviation, Sheppard's adjustments are sufficient, so far as a normal frequency distribution is concerned.

It may be noted that the error in the estimates of sectional μ_2 , due to grouping changes its sign, being negative within the region $0-1.20$ and positive for the rest.

Corrections for the sectional μ_2 : We shall now consider errors in the estimates of the mean referring to portions of a full curve contained within stated intervals of the deviate. Table 3, analogous to Table 2, is showing the mean values for one half of the normal curve in four sub-sections calculated with

different class intervals 'h'. It will be seen, that the means have increasingly been overestimated with increasing 'h'. The true value for these sectional means with infinitesimal class interval can however be theoretically obtained by direct integration. We find that the deviations of the different estimates for the entire half based on different sizes of class-intervals from the true Mean fits excellently in the form $D = h/15$ except for the largest class intervals used. As regards sectional fits we find that a correction term of the form h/α (where α depends on the distance from the origin of the region in which the moment is being calculated) gives good fits.

TABLE 2. DEVIATION IN VALUES OF F_1 BY SECTIONS OF A FULL CURVE [TWO SIDES POOLED]

Section	II(00-1-20)		II(1-21-2-40)		III(2-41-3-60)		IV(3-61-4-80)		Total (-00-4-80)	
	N=7606	N=21374	N=1608	N=32	N=32	N=100000	N=32	N=100000	N=32	N=100000
$\mu_1(0)=34464$	$\mu_1(0)=2.6778$	$\mu_1(0)=7.4104$	$\mu_1(0)=14750$							
A	D	$-A^2/54$	D	$+A^2/57$	D	$+A^2/60$	D	$+A^2/63$	D	$+A^2/12$
.05	-5	-5	8	0	38	28	62	21	21	21
.10	-18	-19	35	37	139	111	63	84	83	83
.15	-	-	-	-	-	-	-	168	168	168
.20	-70	-74	145	148	486	444	125	333	333	333
.25	-	-	-	-	-	-	-	519	519	519
.30	-163	-167	329	333	1045	1000	188	747	750	750
.35	-	-	-	-	-	-	-	1018	1021	1021
.40	-291	-296	584	593	1866	1778	500	1341	1333	1333
.45	-	-	-	-	-	-	-	2979	2963	2963
.50	-891	-867	1329	1333	4130	4000	750	2907	2900	2900
.55	-	-	-	-	-	-	-	5336	5333	5333
.60	-	-	-	-	-	-	-	12003	12000	12000
1-20	-3465	-2667	5825	5333	15800	16000	2875	-	-	-

TABLE 3. ESTIMATED MEANS WITH DIFFERENT CLASS INTERVALS (A).

A	N=38483		II(1-21-2-40)		III(2-41-3-60)		IV(3-61-4-80)		Total (-00-4-80)	
	I(-00-1-20)	N=10087	$\mu_1(0)=53101$	$\mu_1(0)=1.6076$	$\mu_1(0)=2.7090$	$\mu_1(0)=3.8125$	$\mu_1(0)=3.8125$	$\mu_1(0)=3.8125$	N=50,000	$\mu_1(0)=3.8125$
	D	$n^2/22$	D	$n^2/7.5$	D	$n^2/4.5$	D	$n^2/15$	D	$n^2/13$
.05	-00013	-00011	-0001	-00012	-00012	-00006	-00015	-00016	-00017	-00017
.10	47	47	12	13	24	22	213	66	67	67
.15	-	-	-	-	-	-	-	150	160	160
.20	180	182	55	53	99	90	375	266	267	267
.25	-	-	-	-	-	-	-	415	417	417
.30	403	409	118	120	211	200	500	598	600	600
.35	-	-	-	-	-	-	-	816	817	817
.40	718	710	213	213	373	358	875	1068	1067	1067
.45	-	-	-	-	-	-	-	1869	1867	1867
.50	1622	1636	481	480	796	800	1250	2407	2400	2400
.55	-	-	-	-	-	-	-	4304	4288	4288
1-20	8810	6545	1024	1020	2010	2300	3875	9834	9800	9800

Table 4 shows the observed and graduated deviations due to grouping for sections starting from the origin and extended up to the limits 1-20, 2-40, 3-60 and 4-80. It will be seen that the corrections required is maximum for the first and falls rapidly to the limiting value of $A/15$ for the entire half.

TABLE 4. DEVIATIONS IN F_1 UP TO STATED LIMITS FROM THE ORIGIN.

Truncation	-00-1-20		-00-2-40		-00-3-60		-00-4-80	
	A	D	A/22	D	A/15-23	D	A/15	D
.05	-00013	11	-00014	16	-00016	17	-00018	17
.10	47	45	83	65	67	67	66	67
.20	180	182	254	262	266	267	296	297
.30	403	406	573	590	599	600	594	600
.40	718	710	1025	1049	1069	1067	1068	1067
.50	1622	1636	3214	3281	3403	3400	2407	2400
1-20	8810	6545	9512	9442	9828	9600	9434	9600

The error due to grouping in the sectional mean as percentage of the true mean of the half-curve, measured from the general mean, may be expressed in the form $100(2k)^2/(15n\sqrt{12}) = 100(2k)^2/n/12^{3/2}$ approximately where the range is $2k$, the mean $\pm k$, n being the number of cells used and σ the standard deviation. In routine work it is often a practice to split up the range contained within $\pm 3\sigma$ from the general mean into 12 intervals. The error in the estimate of the mean of the half-curve in such a case expressed as percentage to the true mean corresponding reducing to $100(9.5)^2/(12 \times 12) = 2.5$ per cent with σ equal to unity. In table 3 the per cent deviation for $A=1-20$ where the full range is covered by 8 cells with $\sigma=1$ is -12 per cent against an observed percentage deviation of $-0.9834 \times 100/7.8789 = 12.5$ per cent.

Paper received: 6th June, 1943