

A Crash Course to Descriptive Statistics by Jason Schneider

I give credit, with much gratitude, to Dr. Jeffrey Witzel who taught the content of descriptive statistics, intended for linguistics, in which I use the notes from his lectures while I attended graduate school at UTA.

Measures of Central Tendency

Measures of central tendency

- indicate the “typical” behavior of the data

Measures of Central Tendency

1. Mean

-aka, average

-usually symbolized as M

$$M = \frac{\sum X}{N}$$

M = mean

X = individual scores

N = number of scores

\sum = sum of

What is the mean for the values below?

$$(77+75+70+55+4+65+61)/7$$
$$= \text{AVERAGE}(\text{range})$$

Measures of Central Tendency

2. Mode

- the most frequent score

- note that there could be more than one most frequent score.

 - bimodal (if there are two peaks)

 - trimodal (if there are three peaks)....

 - What is the mode for the values from the data?

 - 2, 4, 6, 2, 5, 2, 67, 2, 8, 2, 4, 9, 10, 22, 45

 - =MODE(range)

Measures of Central Tendency

3. Median

-the value that half of the values fall below and the other half fall above

e.g. the median of 100, 95, 83, 71, 61, 57, 30 is...

71

e.g. the median of 100, 95, 83, 75, 71, 61, 57, 30 is...

73

What is the median for the values from the data?

2, 4, 6, 2, 5, 2, 67, 2, 8, 2, 4, 9, 10, 22, 45

=MEDIAN(range)

Measures of Central Tendency

4. Midpoint

-the value halfway between the highest value and the lowest value

$$\text{Midpoint} = \frac{\text{High} + \text{Low}}{2}$$

e.g. if the highest score was 100 and the lowest was 30, then the midpoint is....65

What is the midpoint for 3, 5, 8, 6, 3?
=(MAX(range)+MIN(range))/2

Measures of Dispersion

Measures of dispersion

- indicate how data varied with the typical distribution of data

Measures of Dispersion

1. Range

-the number of points between the highest value and the lowest value +1

What is the range for 3, 5, 8, 6, 3?
 $=\text{MAX}(\text{range}) - \text{MIN}(\text{range}) + 1$

Measures of Dispersion

2. High and low values

- the highest value and the lowest value
- shows where the range falls on the continuum of possible values

What are the high and low values 3, 5, 8, 6, 3?

=MAX(range)

=MIN(range)

Measures of Dispersion

3. Variance

-the average of the squared difference between each value and the mean

$$\text{Variance} = \frac{\sum (X-M)^2}{N}$$

The difference is squared so that the values above and below the mean do not cancel each other out.

$$(77-69)^2 =$$

$$(8)^2 = 64$$

$$(61-69)^2 =$$

$$(-8)^2 = 64$$

What is the variance for values below?

$$\begin{aligned} &77+75+70+55+4+65+61 \\ &= \text{VARP}(\text{range}) \end{aligned}$$

Measures of Dispersion

	A	B	C	D	E
1		X	$(X-M)^2$ $=(X-SE$2)^2$		
2		77	355.59183673	mean $=AVERAGE(B2:B8)$	58.14285714
3		75	284.16326531	variance $=VARP(B2:B8)$	539.5510204
4		70	140.59183673	manual variance $=C9/E5$	539.5510204
5		55	9.8775510204	length $=ROWS(B2:B8)$	7
6		4	2931.4489796		
7		65	47.020408163		
8		61	8.1632653061		
9	Sum $=SUM(B2:B8)$ $=SUM(C2:C8)$	407	3776.8571429		

Let's pretend that we have 30 or more data points to use VARP. However, the calculated variance is wrong here because we should use VAR, not VARP.

Measures of Dispersion

3. Variance

-In addition to

$$\text{Variance} = \frac{\sum(X-M)^2}{N} \quad =\text{VARP}(\text{range})$$

you can also use the following formula

$$\text{Variance} = \frac{\sum(X-M)^2}{N-1} \quad =\text{VAR}(\text{range})$$

→ The first one (the N formula) is for the whole population, while the second one ($N-1$ formula) is for a sample.

→ Population – whole group; sample – portion of the group

→ As a rule of thumb... 30+ values (the N formula);
 under 30 values ($N-1$ formula)

Measures of Dispersion

4. Standard deviation

-aka *SD*

-generally considered the most reliable estimate of dispersion.

$$SD = \sqrt{\frac{\sum(X-M)^2}{N}}$$

Does this look familiar?

SD = square root of variance

The standard deviation shows how much the data deviates from the mean of the dataset.

What is the SD for the values below?

2, 4, 6, 2, 5, 2, 67, 2, 8, 2, 4, 9, 10, 22, 45
=STDEVP(range)

Measures of Dispersion

Think what a large standard deviation might indicate. If there is considerable distance between the data points from the mean of the dataset, this might indicate that our data significantly varies from each other and might be the result of differing factors, making our data the result of different things and not the one thing we might expect. This typically indicates that we have outliers, too.

Measures of Dispersion

	A	B	C	D	E	F
1		77	77	74	69	69
2		75	75	75	68	68
3		70	70	70	70	68
4		55	55	67	67	67
5		4	60	60	60	63
6		65	65	65	65	65
7		61	61	61	61	61
8	Standard Deviation	25.08936409	8.173709305	5.912053869	3.903600292	2.968084199

Notice that the data points become more similar and closer as the standard deviation becomes smaller because the differences in data from the mean of the dataset lessons as the data points become more alike.

Measures of Dispersion

4. Standard deviation

-In addition to

$$SD = \sqrt{\frac{\sum(X-M)^2}{N}} \quad = \text{STDEVP}(\text{range})$$

you can also use the following formula

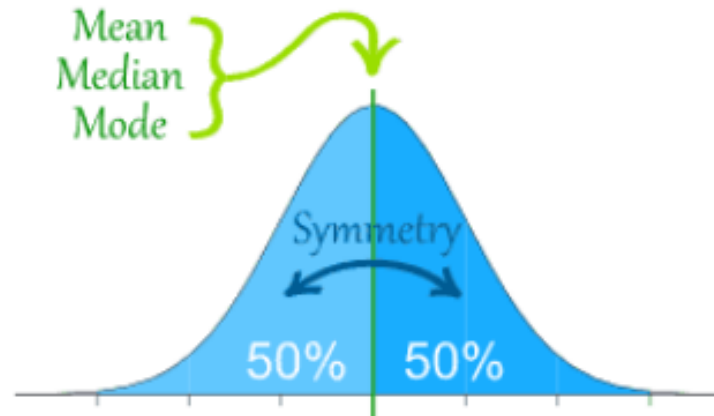
$$SD = \sqrt{\frac{\sum(X-M)^2}{N-1}} \quad = \text{STDEV}(\text{range})$$

→ The first one (the N formula) is for the whole population, while the second one ($N-1$ formula) is for a sample.

→ As a rule of thumb... 30+ values (the N formula);
under 30 values ($N-1$ formula)

Normal Distributions

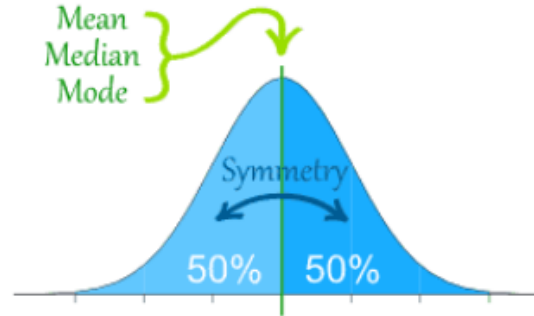
Putting it all together: Descriptive stats and the normal distribution



(from mathisfun.com)

Normal Distributions

Putting it all together: Descriptive stats and the normal distribution



(from mathisfun.com)

Central Tendency

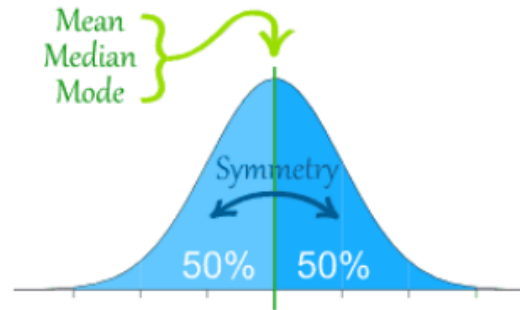
In a real normal distribution ...

- the mean, mode, median, and midpoint will all be the same.

→ Most of the time, however, you have something approximating normal distribution. If it is approximating a normal distribution, then these measures of central tendencies should be similar.

Normal Distributions

Normal distributions



(from mathisfun.com)

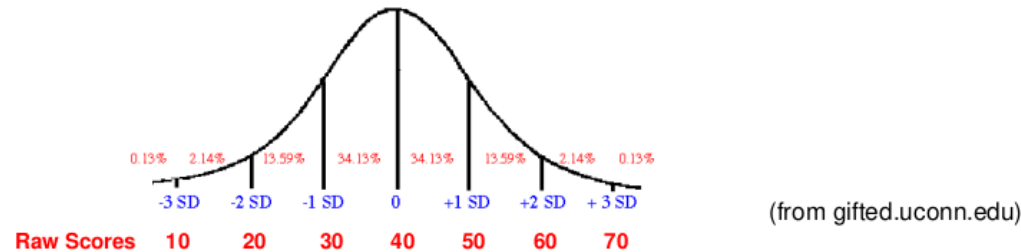
Dispersion

- In a real normal distribution ...
- the lowest and highest values should be exactly the same in distance from the mean.
i.e., the range is symmetrical.
- also, standard deviations (SDs) conform to a certain pattern.

Example of Dispersion

Percentiles

The total percentage of students who scored equal to or below a given point in the normal distribution.



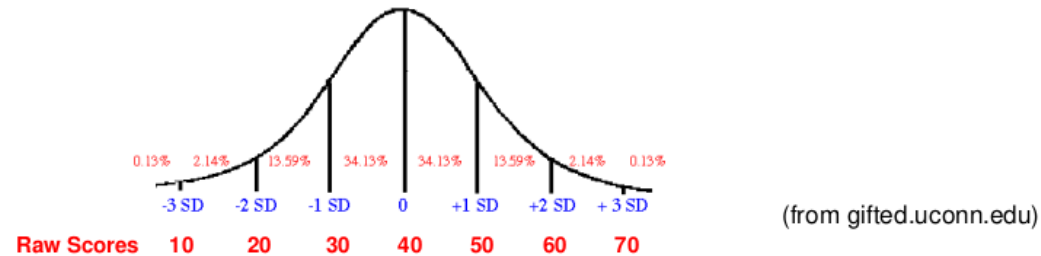
Let's say the mean is 40, and the SD is 10.

- What percentile is the score of 40? ... 50th percentile
- What percentile is the score of 30? ... $0.13 + 2.14 + 13.59 = 15.86 \rightarrow 16^{\text{th}}$
- What percentile is the score of 50?
- What percentile is the score of 60?

Example of Dispersion

Percentiles

The total percentage of students who scored equal to or below a given point in the normal distribution.



Let's say the mean is 40, and the SD is 10.

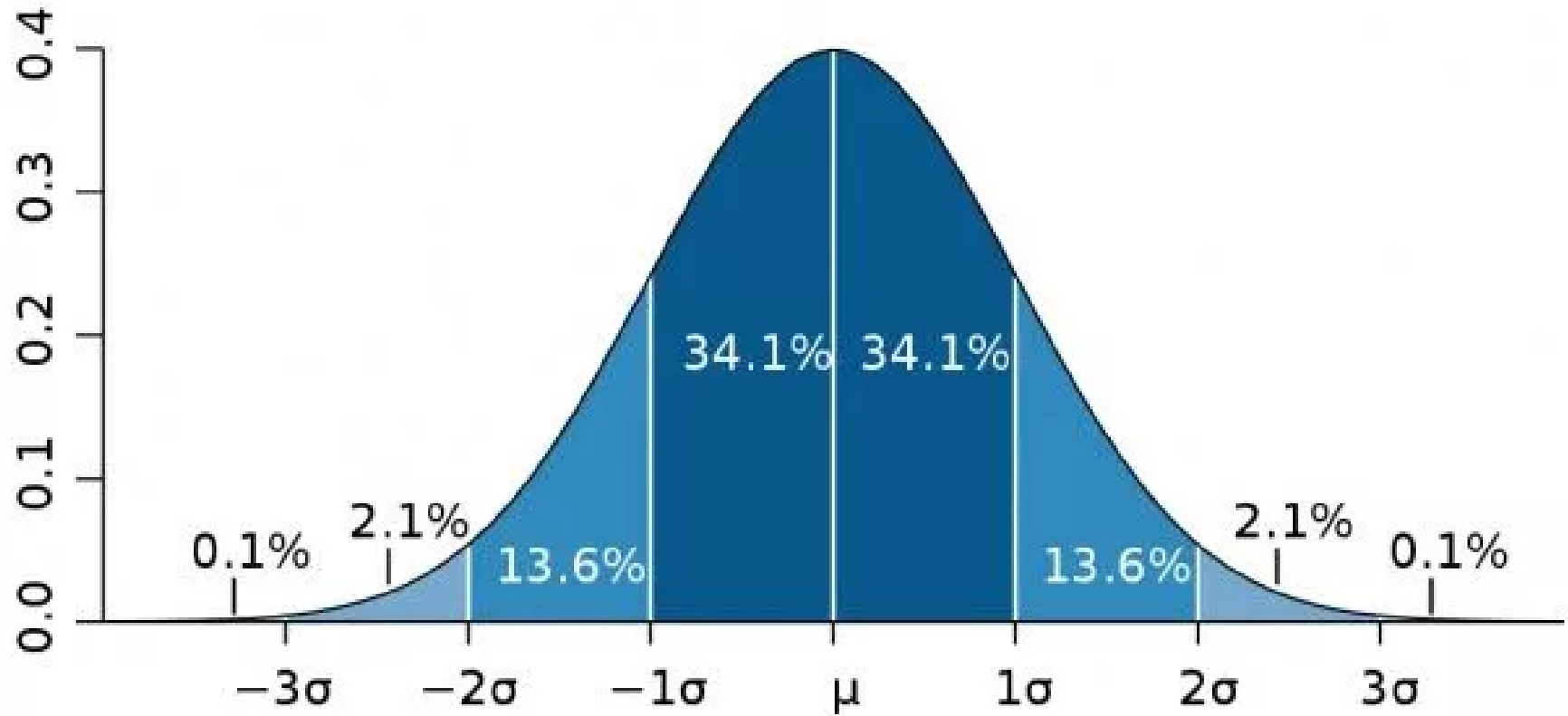
- What percentile is the score of 40? ... $15.86 + 34.13 = 49.99 \rightarrow 50^{\text{th}}$
- What percentile is the score of 30? ... $0.13 + 2.14 + 13.59 = 15.86 \rightarrow 16^{\text{th}}$
- What percentile is the score of 50? ... $50 + 34.13 = 84.13 \rightarrow 84^{\text{th}}$
- What percentile is the score of 60? ... $84.13 + 13.59 = 97.72 \rightarrow 98^{\text{th}}$

Normal Distributions

When a dataset that is normally distributed, the normal distribution will follow the below rules:

- The center of the bell curve is the mean of the data point (also the highest point in the bell curve).
- 68.2% of the total data points lie in the range (Mean – Standard Deviation to Mean + Standard Deviation).
- 95.5% of the total data points lie in the range (Mean – 2*Standard Deviation to Mean + 2*Standard Deviation)
- 99.7% of the total data points lie in the range (Mean – 3*Standard Deviation to Mean + 3*Standard Deviation)

Normal Distributions



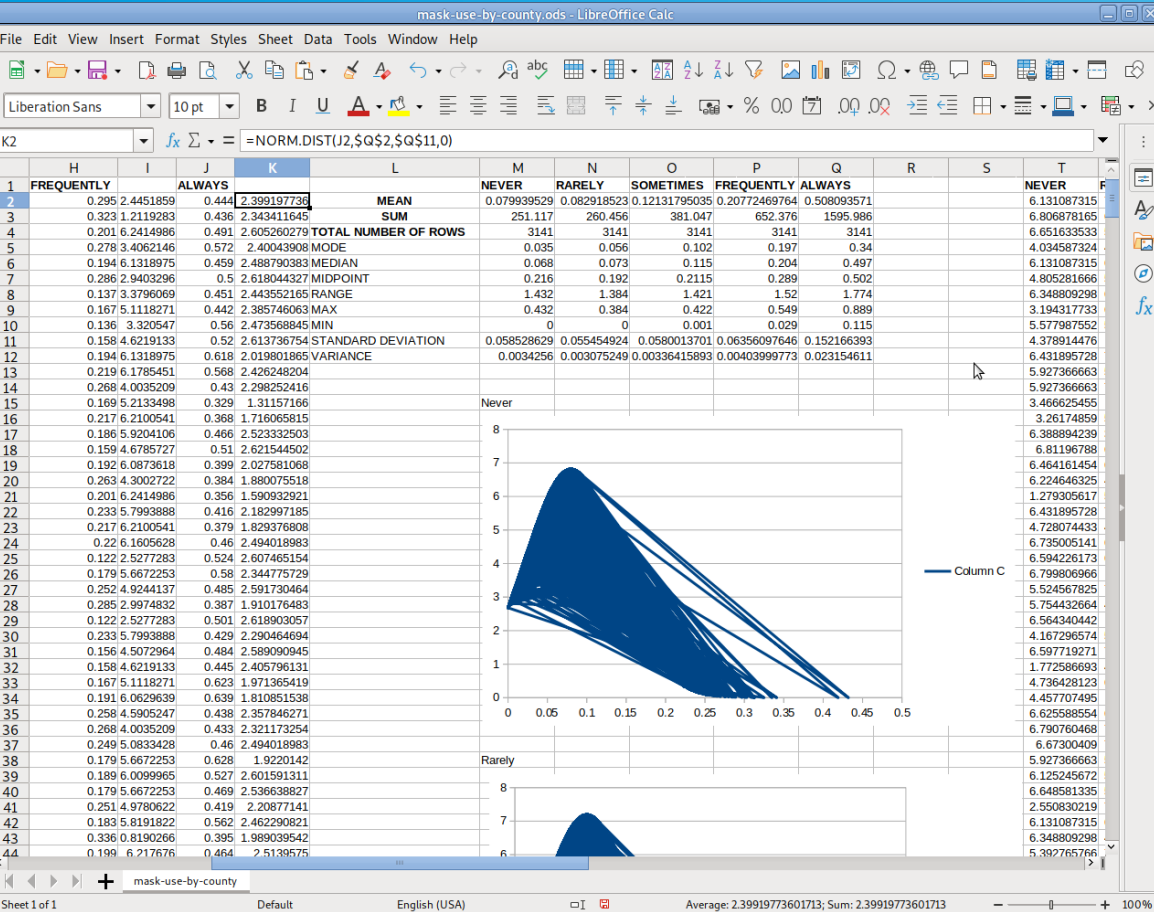
Normal Distributions

Use the formula **Mean – 3* Standard Deviation** to develop the y-axis. The mean and standard deviation should be created from the column that will be the x-axis.

- In Excel or Libre Calc, use the function **=NORM.DIST(the cell from the spread sheet, the mean of the column that contains that cell, the standard deviation of the column that contains that cell, False)**
- For example, **=NORM.DIST(J2,\$Q\$2,\$Q\$11,0)**

The fourth input can be False or 0, the same thing. Use the dollar signs to prevent Excel or Libre Calc from altering the cell if dragging or copying. Lastly, do not forget using the equal sign before the function; otherwise, the spreadsheet will not recognize this as a function.

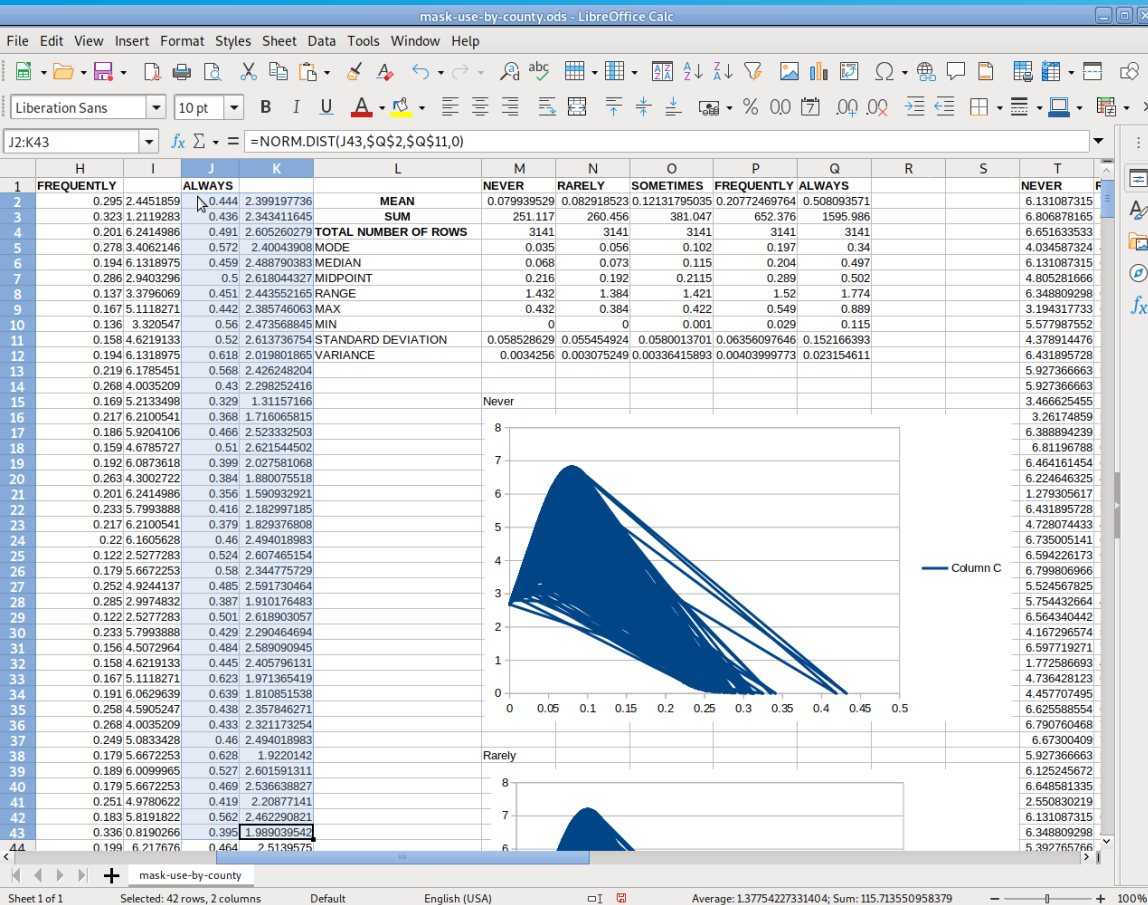
Normal Distributions



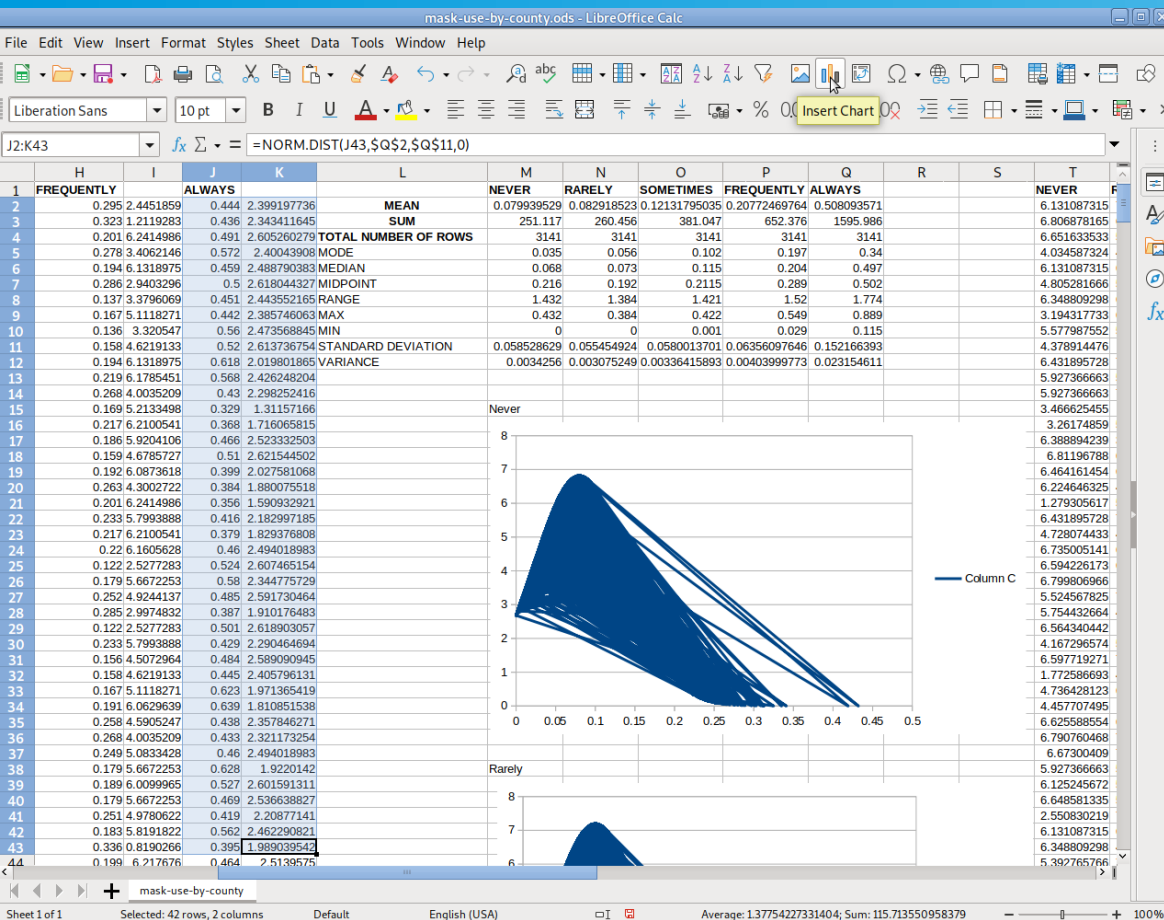
For example, use the norm.dist function in the cell k2, applying to j2. Then this same function can be applied to all the cells in k as shown in the picture to the left.

Normal Distributions

Then select all cells in both j and k columns.

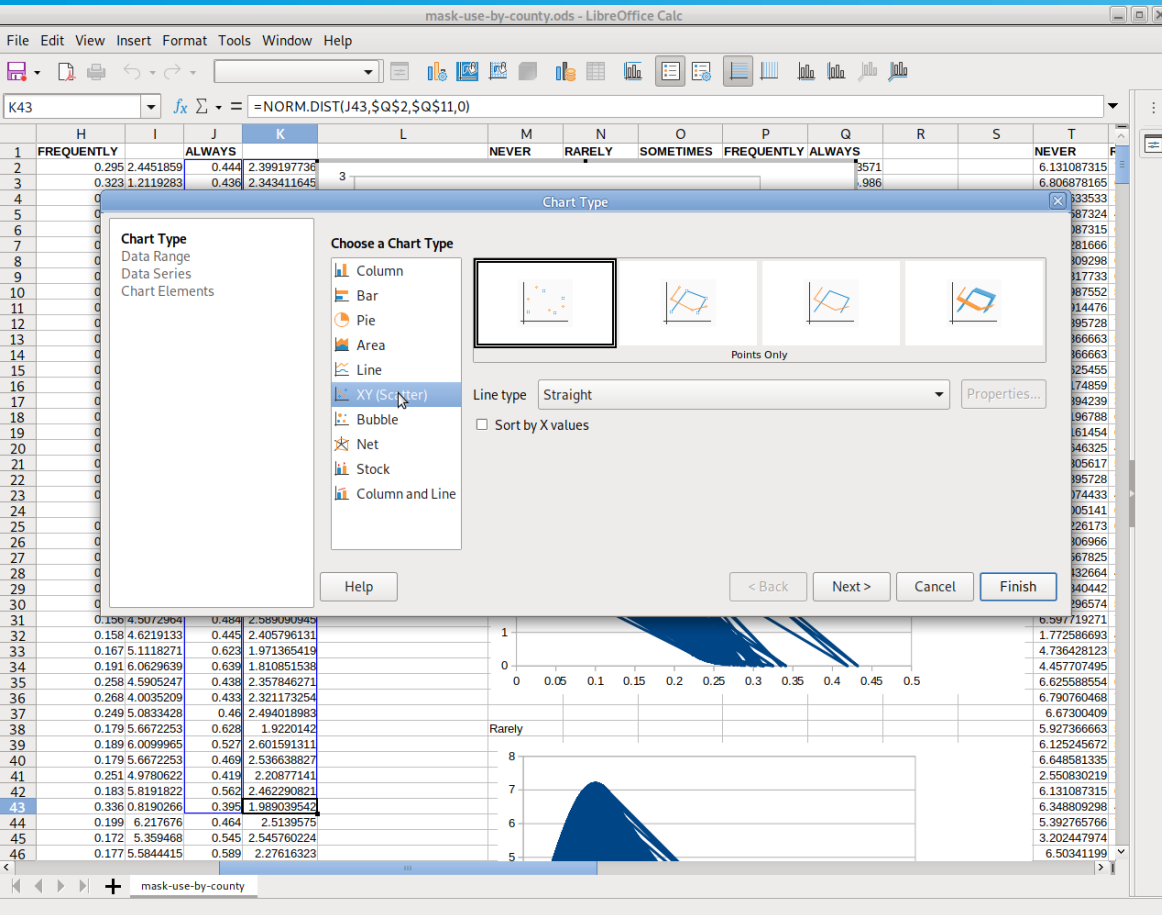


Normal Distributions



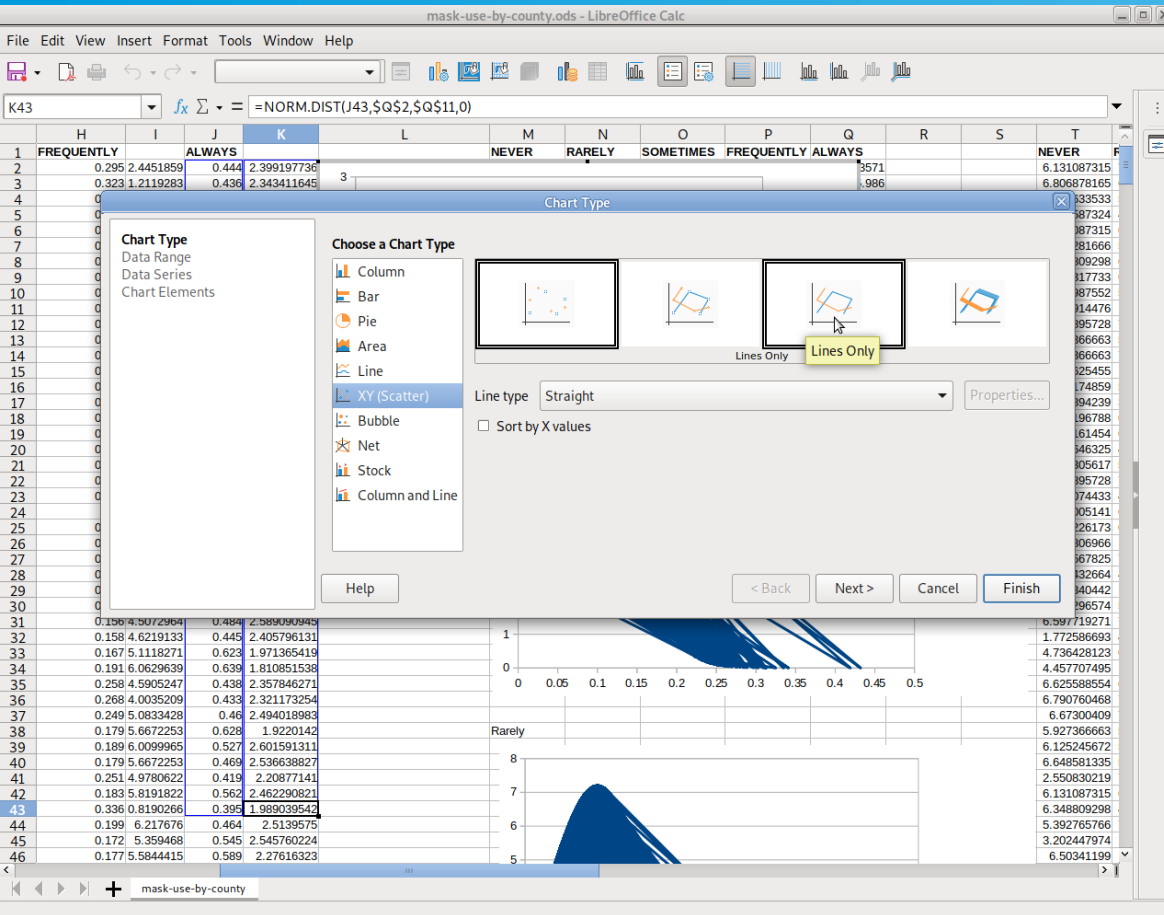
Click on Insert Chart as shown in the picture to the left as columns j and k are selected, indicated by j and k being highlighted blue.

Normal Distributions



After clicking on Insert Chart, select XY (Scatter), or Scatter, depending on the office program used for opening spreadsheets like Excel.

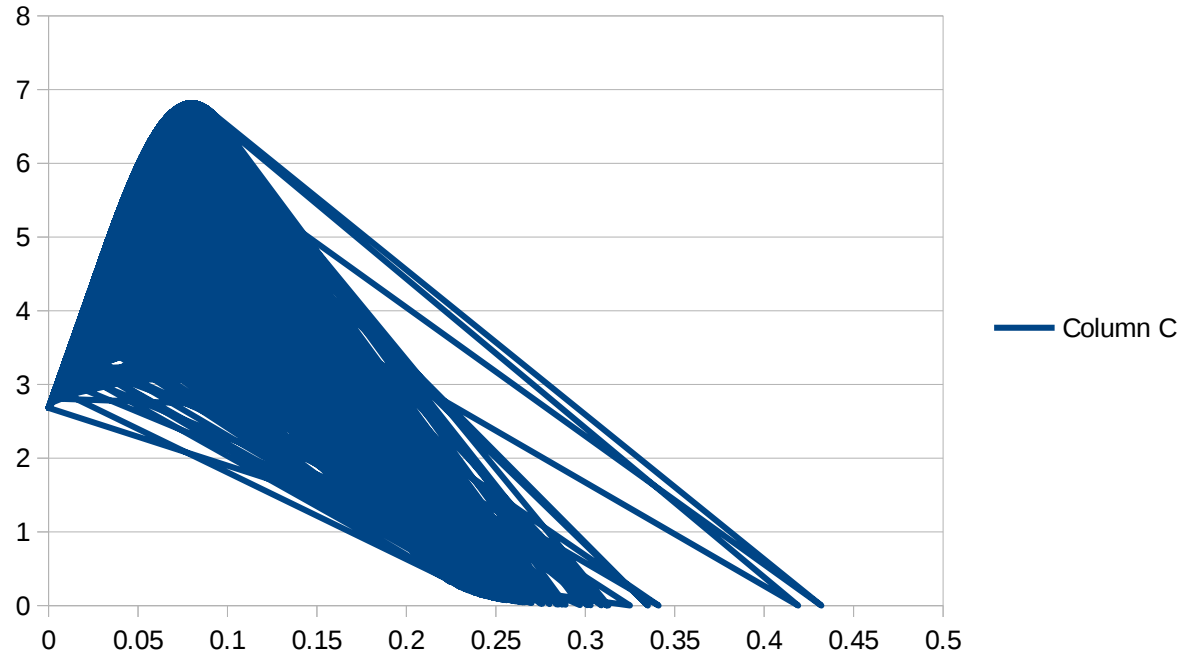
Normal Distributions



Lastly, select Lines Only, or Smooth Lines if using Excel, before clicking on Finish to create a graph or figure for a normal distribution as shown of the skewed distribution that has been shown in the pictures illustrating these steps.

Normal Distributions

Then the result should look like the figure or graph shown to the right, illustrating a normal distribution varying according to how skewed it might be.



Summary of Functions

Excel cheat sheet

Number of students (N)	<code>=COUNT (range)</code>
Mean	<code>=AVERAGE (range)</code>
Mode	<code>=MODE (range)</code>
Median	<code>=MEDIAN (range)</code>
Midpoint	<code>= (MAX (range) + MIN (range)) / 2</code>
High	<code>=MAX (range)</code>
Low	<code>=MIN (range)</code>
Range	<code>=MAX (range) - MIN (range) + 1</code>
Variance (N formula)	<code>=VARP (range)</code>
Variance ($N-1$ formula)	<code>=VAR (range)</code>
SD (N formula)	<code>=STDEVP (range)</code>
SD ($N-1$ formula)	<code>=STDEV (range)</code>

Example table

Table 4.7: Fall 1986 (First Administration) ELIPT Results

Statistics	Subtests			
	Listening	Reading	Vocabulary	Writing
<i>N</i>	153.00	153.00	154.00	153.00
Total items (<i>k</i>)	55.00	60.00	100.00	100.00
Mean (\bar{X})	34.76	40.64	69.34	75.08
Mode	32.00	43.00	86.00	77.00
Median	34.45	41.00	71.67	75.50
Midpoint	34.50	39.00	59.50	69.00
Low-High	17-52	21-57	20-99	44-94
Range	36.00	37.00	80.00	51.00
<i>S</i>	7.29	7.48	16.08	8.94

(Brown, 1996)

Example table

Table 4.8: Fall 1986 ELIPT Results

Subtest	<i>N</i>	<i>k</i>	Central Tendency				Dispersion		
			\bar{X}	Mode	Median	Midpoint	Low-High	Range	<i>S</i>
Listening	153	55	54.76	52	34.45	54.50	17-52	36	7.29
Reading	153	60	40.64	43	41.00	39.00	21-57	37	7.48
Vocabulary	154	100	69.34	86	71.67	59.50	29-99	80	16.08
Writing	153	100	75.08	77	75.50	69.00	44-94	51	8.94

(Brown, 1996)